# SC'14 Tutorial:
# Hands-on Practical Hybrid Parallel Application Performance Engineering

**Markus Geimer**
Jülich Supercomputing Centre

**Sameer Shende**
University of Oregon

**Bert Wesarg**
Technische Universität Dresden

**Brian Wylie**
Jülich Supercomputing Centre

# Agenda

| Time | Topic | Presenter |
|------|-------|-----------|
| 08:30 | Introduction to VI-HPS & parallel performance engineering | Wylie |
| 09:15 | VI-HPS Linux Live-ISO and MPI+OpenMP example code | Wylie / all |
| 09:30 | Instrumentation & measurement with **Score-P** | Wesarg |
| 10:00 | *Break* | |
| 10:30 | Profile examination with **CUBE** | Geimer |
| 11:00 | Configuration & customization of Score-P measurements | Geimer |
| 11:30 | Profile examination with **TAU** ParaProf | Shende |
| 12:00 | *Lunch* | |
| 13:30 | Automated trace analysis with **Scalasca** | Geimer |
| 14:15 | Interactive trace analysis with **Vampir** | Wesarg |
| 15:00 | *Break* | |
| 15:30 | Specialized Score-P measurements & analysis | Wesarg |
| 16:00 | Performance data management with **TAU** PerfExplorer | Shende |
| 16:15 | Finding typical parallel performance bottlenecks | Wesarg |
| 16:45 | Review & conclusion | Wylie |
| 17:00 | *Adjourn* | |

# Introduction to VI-HPS

## Brian Wylie
## Jülich Supercomputing Centre

**Mission**: Improve the quality and accelerate the development process of complex simulation codes running on highly-parallel computer systems

- Start-up funding (2006–2011) by Helmholtz Association of German Research Centres

**HELMHOLTZ** | **ASSOCIATION**

- Activities
  - Development and integration of HPC programming tools
    - diagnose programming errors and optimization opportunities
  - Training & support to apply these tools
  - Academic workshops

## http://www.vi-hps.org

# Forschungszentrum Jülich

- Jülich Supercomputing Centre

# RWTH Aachen University

- Centre for Computing & Communication

# Technische Universität Dresden

- Centre for Information Services & HPC

# University of Tennessee (Knoxville)

- Innovative Computing Laboratory

# Barcelona Supercomputing Center

- Centro Nacional de Supercomputación

# German Research School

- Laboratory of Parallel Programming

# Lawrence Livermore National Lab.

- Centre for Applied Scientific Computing

# Technical University of Munich

- Chair for Computer Architecture

# University of Oregon

- Performance Research Laboratory

# University of Stuttgart

- HPC Centre

# University of Versailles St-Quentin

- LRC ITACA

# Allinea Software Ltd

# MUST

- MPI usage correctness checking

# PAPI

- Interfacing to hardware performance counters

# Periscope

- Automatic analysis via an on-line distributed search

# Scalasca

- Large-scale parallel performance analysis

# TAU

- Integrated parallel performance system

# Vampir

- Interactive graphical trace visualization & analysis

# Score-P

- Community instrumentation & measurement infrastructure

# Productivity tools (cont.)

DDT/MAP/PR
- Parallel debugging & profiling

KCachegrind
- Callgraph-based cache analysis [x86 only]

MAQAO
- Assembly instrumentation & optimization [x86-64 only]

mpiP/mpiPview
- MPI profiling tool and analysis viewer

Open MPI
- Integrated memory checking

Open|Speedshop
- Integrated parallel performance analysis environment

Paraver/Dimemas/Extrae
- Event tracing and graphical trace visualization & analysis
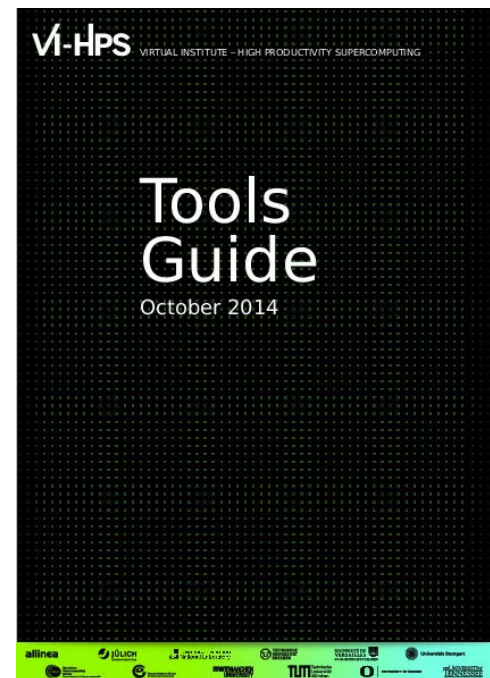
Rubik
- Process mapping generation & optimization [BG only]

SIONlib/Spindle
- Optimized native parallel file I/O & library loading

STAT
- Stack trace analysis tools

For a brief overview of tools consult the VI-HPS Tools Guide:

# Technologies and their integration



KCACHEGRIND

LWM2 / MAP / MPIP / O|SS / MAQAO

TAU

SCORE-P

PAPI

Hardware monitoring

Automatic profile & trace analysis

PERISCOPE

MUST

SCALASCA

DDT

Debugging, error & anomaly detection

Visual trace analysis

VAMPIR / PARAVER

STAT

Execution

Optimization

SYSMON / SPINDLE / SIONLIB / OPENMPI

RUBIK / MAQAO

Tools will ***not*** automatically make you, your applications or computer systems more *productive*.

However, they can help you understand ***how*** your parallel code executes and ***when / where*** it's necessary to work on *correctness* and *performance* issues.

- Goals
  - Give an overview of the programming tools suite
  - Explain the functionality of individual tools
  - Teach how to use the tools effectively
  - Offer hands-on experience and expert assistance using tools
  - Receive feedback from users to guide future development

- For best results, bring & analyze/tune your own code(s)!

- VI-HPS Hands-on Tutorial series
  - SC'08, ICCS'09, SC'09, Cluster'10, SC'10, SC'11, EuroMPI'12, XSEDE'13, SC'13, *SC'14 (New Orleans)*

- VI-HPS Tuning Workshop series
  - 2008 (Aachen & Dresden), 2009 (Jülich & Bremen), 2010 (Garching & Amsterdam/NL), 2011 (Stuttgart & Aachen), 2012 (St-Quentin/F & Garching), 2013 (Saclay/F & Jülich) 2014 (Barcelona/Spain, Kobe/Japan, Saclay/France, Edinburgh/UK)

- ## 17th VI-HPS Tuning Workshop (23-27 February 2015)
  - Hosted by HLRS, Stuttgart, Germany
  - Using PRACE Tier-0 *Hornet* Cray XC40 system
  - VI-HPS and Cray tools to be presented

- ## Further events to be determined
  - (one-day) tutorials
    - With guided exercises usually using a Live-ISO
  - (multi-day) training workshops
    - With your own applications on actual HPC systems

- ## Check www.vi-hps.org/training for announced events
- ## Contact us if you might be interested in hosting an event

- Bootable Linux installation on DVD (or USB memory stick)
- Includes everything needed to try out our parallel tools on an 64-bit x86-architecture notebook computer
  - VI-HPS tools: MUST, PAPI, Score-P, Periscope, Scalasca, TAU, Vampir*
  - Also: Eclipse/PTP, DDT*, TotalView*
    - \* time/capability-limited evaluation licences provided for commercial products
  - GCC (w/ OpenMP), OpenMPI
  - Manuals/User Guides
  - Tutorial exercises & examples
- Produced by U. Oregon PRL
  - Sameer Shende

**Parallel Productivity Tools Live DVD**
**Also includes:** TotalView, DyninstAPI, PDT, Eclipse PTP, Berkeley UPC, ptoolsrte, Chapel, and much more…

**Partners:**
ParaTools, Inc.
University of Florida
University of Oregon
Totalview Technologies
RWTH Aachen University
HLRS / University of Stuttgart
Jülich Supercomputing Centre
Technische Universität Dresden
Technische Universität München
University of Wisconsin at Madison
Pittsburgh Supercomputing Center
University of Tennessee at Knoxville
National Center for Supercomputing Applications

November 2011

**POINT VI-HPS**

http://nic.uoregon.edu/point   http://www.vi-hps.org

- ISO image approximately 10GB
  - download latest version from website
  - http://www.vi-hps.org/training/live-iso/
  - optionally create bootable DVD or USB drive

- Boot directly from disk
  - enables hardware counter access and offers best performance, but no save/resume

- Boot within virtual machine (e.g., VirtualBox)
  - faster boot time and can save/resume state, but may not allow hardware counter access

- Boots into Linux environment for HPC
  - supports building and running provided MPI and/or OpenMP parallel application codes
  - and experimentation with VI-HPS (and third-party) tools

Difference Engine

"The most constant difficulty in contriving the engine has arisen from the desire to reduce the time in which the calculations were executed to the shortest which is possible."

Charles Babbage
1791 – 1871

- # Moore's law is still in charge, but
    - ## Clock rates no longer increase
    - ## Performance gains only through increased parallelism
- # Optimizations of applications more difficult
    - ## Increasing application complexity
        - ### Multi-physics
        - ### Multi-scale
    - ## Increasing machine complexity
        - ### Hierarchical networks / memory
        - ### More CPUs / multi-core

☞ Every doubling of scale reveals a new bottleneck!

- "Sequential" performance factors
  - Computation
    - ☞ Choose right algorithm, use optimizing compiler
  - Cache and memory
    - ☞ Tough! Only limited tool support, hope compiler gets it right
  - Input / output
    - ☞ Often not given enough attention

- "Parallel" performance factors
  - Partitioning / decomposition
  - Communication (i.e., message passing)
  - Multithreading
  - Synchronization / locking
    - ☞ More or less understood, good tool support

- Successful engineering is a combination of
    - The right algorithms and libraries
    - Compiler flags and directives
    - Thinking !!!

- Measurement is better than guessing
    - To determine performance bottlenecks
    - To compare alternatives
    - To validate tuning decisions and optimizations
        - ☞After each step!

"We should forget about small efficiencies, say 97% of the time: premature optimization is the root of all evil."

Charles A. R. Hoare

■ It's easier to optimize a slow correct program than to debug a fast incorrect one

☞ *Nobody cares how fast you can compute a wrong answer...*

# Performance engineering workflow

- Prepare application with symbols
- Insert extra code (probes/hooks)

- Collection of performance data
- Aggregation of performance data

**Preparation** | **Measurement**

**Optimization** | **Analysis**

- Modifications intended to eliminate/reduce performance problem

- Calculation of metrics
- Identification of performance problems
- Presentation of results

- Programs typically spend 80% of their time in 20% of the code

- Programmers typically spend 20% of their effort to get 80% of the total speedup possible for the application
  - ☞ *Know when to stop!*

- Don't optimize what does not matter
  - ☞ *Make the common case fast!*

"If you optimize everything, you will always be unhappy."

Donald E. Knuth

- # What can be measured?

  - ## A **count** of how often an event occurs
    - ### E.g., the number of MPI point-to-point messages sent

  - ## The **duration** of some interval
    - ### E.g., the time spent these send calls

  - ## The **size** of some parameter
    - ### E.g., the number of bytes transmitted by these calls

- # Derived metrics

  - ## E.g., rates / throughput
  - ## Needed for normalization

- ## Execution time

- ## Number of function calls

- ## CPI
  - ### CPU cycles per instruction

- ## FLOPS
  - ### Floating-point operations executed per second

"math" Operations?
HW Operations?
HW Instructions?
32-/64-bit? …

- # Wall-clock time
  - ## Includes waiting time: I/O, memory, other system activities
  - ## In time-sharing environments also the time consumed by other applications
- # CPU time
  - ## Time spent by the CPU to execute the application
  - ## Does not include time the program was context-switched out
    - ### Problem: Does not include inherent waiting time (e.g., I/O)
    - ### Problem: Portability? What is user, what is system time?
- # Problem: Execution time is non-deterministic
  - ## Use mean or minimum of several runs

- ## Inclusive
  - ### Information of all sub-elements aggregated into single value
- ## Exclusive
  - ### Information cannot be subdivided further

- # How are performance measurements triggered?
  - Sampling
  - Code instrumentation

- # How is performance data recorded?
  - Profiling / Runtime summarization
  - Tracing

- # How is performance data analyzed?
  - Online
  - Post mortem

```
int main()
{
  int i;

  for (i=0; i < 3; i++)
    foo(i);

  return 0;
}

void foo(int i)
{

  if (i > 0)
    foo(i - 1);

}
```

- ■ Running program is periodically interrupted to take measurement
  - ■ Timer interrupt, OS signal, or HWC overflow
  - ■ Service routine examines return-address stack
  - ■ Addresses are mapped to routines using symbol table information
- ■ **Statistical** inference of program behavior
  - ■ Not very detailed information on highly volatile metrics
  - ■ Requires long-running applications
- ■ Works with unmodified executables

```
int main()
{
  int i;
  Enter("main");
  for (i=0; i < 3; i++)
    foo(i);
  Leave("main");
  return 0;
}

void foo(int i)
{
  Enter("foo");
  if (i > 0)
    foo(i – 1);
  Leave("foo");
}
```

- **Measurement code is inserted such that every event of interest is captured directly**
  - Can be done in various ways
- Advantage:
  - Much more detailed information
- Disadvantage:
  - Processing of source-code / executable necessary
  - Large relative overheads for small functions

- ## Static instrumentation
  - Program is instrumented prior to execution

- ## Dynamic instrumentation
  - Program is instrumented at runtime

- # Code is inserted
  - Manually
  - Automatically
    - By a preprocessor / source-to-source translation tool
    - By a compiler
    - By linking against a pre-instrumented library / runtime system
    - By binary-rewrite / dynamic instrumentation tool

# Accuracy

- Intrusion overhead
  - Measurement itself needs time and thus lowers performance
- Perturbation
  - Measurement alters program behaviour
  - E.g., memory access pattern
- Accuracy of timers & counters

# Granularity

- How many measurements?
- How much information / processing during each measurement?

☞ *Tradeoff: Accuracy vs. Expressiveness of data*

- How are performance measurements triggered?
  - Sampling
  - Code instrumentation

- How is performance data recorded?
  - Profiling / Runtime summarization
  - Tracing

- How is performance data analyzed?
  - Online
  - Post mortem

- Recording of aggregated information
  - Total, maximum, minimum, …

- For measurements
  - Time
  - Counts
    - Function calls
    - Bytes transferred
    - Hardware counters

- Over program and system entities
  - Functions, call sites, basic blocks, loops, …
  - Processes, threads

☞ *Profile = summarization of events over execution interval*

- # Flat profile
  - Shows distribution of metrics per routine / instrumented region
  - Calling context is not taken into account

- # Call-path profile
  - Shows distribution of metrics per executed call path
  - Sometimes only distinguished by partial calling context (e.g., two levels)

- # Special-purpose profiles
  - Focus on specific aspects, e.g., MPI calls or OpenMP constructs
  - Comparing processes/threads

- Recording detailed information about significant points (events) during execution of the program
    - Enter / leave of a region (function, loop, …)
    - Send / receive a message, …
- Save information in event record
    - Timestamp, location, event type
    - Plus event-specific information (e.g., communicator, sender / receiver, …)
- Abstract execution model on level of defined events

☞ *Event trace = Chronologically ordered sequence of event records*

# Event tracing

**Process A**

```
void foo() {
  trc_enter("foo");
  ...
  trc_send(B);
  send(B, tag, buf);
  ...
  trc_exit("foo");
}
```

**instrument**

**Process B**

```
void bar() {
  trc_enter("bar");
  ...
  recv(A, tag, buf);
  trc_recv(A);
  ...
  trc_exit("bar");
}
```

**MONITOR**

**MONITOR**

synchronize(d)

**Local** trace A

| ... | |
|-----|-----|
| 58 | ENTER foo |
| 62 | SEND to B |
| 64 | EXIT foo |
| ... | |

**Local** trace B

| ... | |
|-----|-----|
| 60 | ENTER bar |
| 68 | RECV from A |
| 69 | EXIT bar |
| ... | |

**Global** trace view

| ... | | |
|-----|-----|-----|
| 58 | A | ENTER foo |
| 60 | B | ENTER bar |
| 62 | A | SEND to B |
| 64 | A | EXIT foo |
| 68 | B | RECV from A |
| 69 | B | EXIT bar |
| ... | | |

**merge**

- # Tracing advantages

    - Event traces preserve the **temporal** and **spatial** relationships among individual events (☞ context)
    - Allows reconstruction of **dynamic** application behaviour on any required level of abstraction
    - Most general measurement technique
        - Profile data can be reconstructed from event traces

- # Disadvantages

    - Traces can very quickly become extremely large
    - Writing events to file at runtime may causes perturbation

- **How are performance measurements triggered?**
  - Sampling
  - Code instrumentation

- **How is performance data recorded?**
  - Profiling / Runtime summarization
  - Tracing

- **How is performance data analyzed?**
  - Online
  - Post mortem

- ## Performance data is processed during measurement run

  - ### Process-local profile aggregation

  - ### More sophisticated inter-process analysis using

    - "Piggyback" messages

    - Hierarchical network of analysis agents

- ## Inter-process analysis often involves application steering to interrupt and re-configure the measurement

- Performance data is stored at end of measurement run

- Data analysis is performed afterwards

  - Automatic search for bottlenecks

  - Visual trace analysis

  - Calculation of statistics

**Global** trace view

| | | |
|---|---|---|
| **...** | | |
| 58 | A | ENTER foo |
| 60 | B | ENTER bar |
| 62 | A | SEND to B |
| 64 | A | EXIT foo |
| 68 | B | RECV from A |
| 69 | B | EXIT bar |
| **...** | | |

**Post-Mortem Analysis** →

Legend:
- main
- foo
- bar

☞*A combination of different methods, tools and techniques is typically needed!*

- Analysis
    - Statistics, visualization, automatic analysis, data mining, ...
- Measurement
    - Sampling / instrumentation, profiling / tracing, ...
- Instrumentation
    - Source code / binary, manual / automatic, ...

- ## Do I have a performance problem at all?
  - ### Time / speedup / scalability measurements
- ## What is the key bottleneck (computation / communication)?
  - ### MPI / OpenMP / flat profiling
- ## Where is the key bottleneck?
  - ### Call-path profiling, detailed basic block profiling
- ## Why is it there?
  - ### Hardware counter analysis, trace selected parts to keep trace size manageable
- ## Does the code have scalability problems?
  - ### Load imbalance analysis, compare profiles at various sizes function-by-function

# Hands-on example code:
# NPB-MZ-MPI / BT
# (on Live-ISO/DVD)

VI-HPS Team

1. Reference preparation for validation
2. Program instrumentation
3. Summary measurement collection
4. Summary analysis report examination
5. Summary experiment scoring
6. Summary measurement collection with filtering
7. Filtered summary analysis report examination
8. Event trace collection
9. Event trace examination & analysis

- The NAS Parallel Benchmark suite (MPI+OpenMP version)
  - Available from

    http://www.nas.nasa.gov/Software/NPB

  - 3 benchmarks in Fortran77
  - Configurable for various sizes & classes
- Move into the NPB3.3-MZ-MPI root directory

```
% cd Tutorial; ls
bin/      common/  jobscript/  Makefile  README.install  SP-MZ/
BT-MZ/  config/  LU-MZ/        README    README.tutorial  sys/
```

- Subdirectories contain source code for each benchmark
  - plus additional configuration and common code
- The provided distribution has already been configured for the tutorial, such that it's ready to "make" one or more of the benchmarks and install them into a (tool-specific) "bin" subdirectory

# Building an NPB-MZ-MPI benchmark

- ## Type "make" for instructions

```
% make

   =================================================
   =       NAS PARALLEL BENCHMARKS 3.3         =
   =       MPI+OpenMP Multi-Zone Versions      =
   =       F77                                 =
   =================================================

   To make a NAS multi-zone benchmark type

           make <benchmark-name> CLASS=<class> NPROCS=<nprocs>

   where <benchmark-name> is "bt-mz", "lu-mz", or "sp-mz"
         <class>            is "S", "W", "A" through "F"
         <nprocs>           is number of processes

   [...]

   ***********************************************************
   * Custom build configuration is specified in config/make.def  *
   * Suggested tutorial exercise configuration for LiveISO/DVD:   *
   *          make bt-mz CLASS=W NPROCS=4                         *
   ***********************************************************
```

Hint: the recommended build configuration is available via
```
% make suite
```

- Specify the benchmark configuration
  - benchmark name: **bt-mz**, lu-mz, sp-mz
  - the number of MPI processes: NPROCS=**4**
  - the benchmark class (S, W, A, B, C, D, E): CLASS=**W**

```
% make bt-mz CLASS=W NPROCS=4
cd BT-MZ; make CLASS=W NPROCS=4 VERSION=
make: Entering directory 'BT-MZ'
cd ../sys; cc  -o setparams setparams.c
../sys/setparams bt-mz 4 W
mpif77 -c  -O3 -fopenmp bt.f
  [...]
cd ../common;  mpif77 -c  -O3 -fopenmp timers.f
mpif77 –O3 -fopenmp -o ../bin/bt-mz_W.4 \
bt.o initialize.o exact_solution.o exact_rhs.o set_constants.o \
adi.o rhs.o zone_setup.o x_solve.o y_solve.o exch_qbc.o \
solve_subs.o z_solve.o add.o error.o verify.o mpi_setup.o \
../common/print_results.o ../common/timers.o
Built executable ../bin/bt-mz_W.4
make: Leaving directory 'BT-MZ'
```

- ## What does it do?
  - Solves a discretized version of unsteady, compressible Navier-Stokes equations in three spatial dimensions
  - Performs 200 time-steps on a regular 3-dimensional grid

- ## Implemented in 20 or so Fortran77 source modules

- ## Uses MPI & OpenMP in combination
  - 4 processes with 4 threads each should be reasonable
    - don't expect to see speed-up when run on a laptop!
  - bt-mz_W.4 should run in around 5 to 12 seconds on a laptop
  - bt-mz_B.4 is more suitable for dedicated HPC compute nodes
    - Each class step takes around 10-15x longer

- ## Launch as a hybrid MPI+OpenMP application

Alternatively execute script:
```
% sh ../jobscript/ISO/run.sh
```

```
% cd bin
% OMP_NUM_THREADS=4  mpiexec -np 4 ./bt-mz_W.4
 NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark
 Number of zones:    4 x    4
 Iterations: 200   dt:   0.000800
 Number of active processes:     4
 Total number of threads:     16  (  4.0 threads/process)

 Time step    1
 Time step   20
 Time step   40
  [...]
 Time step  160
 Time step  180
 Time step  200
 Verification Successful

 BT-MZ Benchmark Completed.
 Time in seconds = 5.57
```

Hint: save the benchmark output (or note the run time) to be able to refer to it later

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

Markus Geimer[2], Bert Wesarg[1], Brian Wylie[2]

With contributions from
Andreas Knüpfer[1] and Christian Rössel[2]
[1]ZIH TU Dresden , [2]FZ Jülich

- Several performance tools co-exist
- Separate 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

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/SHMEM, OpenMP/Pthreads, 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/OpenACC/…

- 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

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

- 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

# Score-P hands-on:
# NPB-MZ-MPI / BT

1. Reference preparation for validation
2. Program instrumentation
3. Summary measurement collection
4. Summary analysis report examination
5. Summary experiment scoring
6. Summary measurement collection with filtering
7. Filtered summary analysis report examination
8. Event trace collection
9. Event trace examination & analysis

- Change back to directory containing NPB BT-MZ

```
% cd ..
```

- Edit config/make.def to adjust build configuration
  - Modify specification of compiler/linker: MPIF77

```
...
#-----------------------------------------------------------------
# The Fortran compiler used for MPI programs
#-----------------------------------------------------------------
#MPIF77 = mpif77

# Alternative variants to perform instrumentation
...
MPIF77 = scorep mpif77

# This links MPI Fortran programs; usually the same as ${MPIF77}
FLINK  = $(MPIF77)
...
```

Uncomment the Score-P compiler wrapper specification

- ## Return to root directory and clean-up

```
% make clean
```

- ## Re-build executable using Score-P instrumenter

```
% make bt-mz CLASS=W NPROCS=4
cd BT-MZ; make CLASS=W NPROCS=4 VERSION=
make: Entering directory 'BT-MZ'
cd ../sys; cc  -o setparams setparams.c -lm
../sys/setparams bt-mz 4 W
scorep mpif77 -c  -O3 -fopenmp bt.f
 [...]
cd ../common;  scorep mpif77 -c  -O3 -fopenmp timers.f
scorep mpif77 -O3 -fopenmp -o ../bin.scorep/bt-mz_W.4 \
bt.o initialize.o exact_solution.o exact_rhs.o set_constants.o \
adi.o rhs.o zone_setup.o x_solve.o y_solve.o exch_qbc.o \
solve_subs.o z_solve.o add.o error.o verify.o mpi_setup.o \
../common/print_results.o ../common/timers.o
Built executable ../bin.scorep/bt-mz_W.4
make: Leaving directory 'BT-MZ'
```

1. Reference preparation for validation
2. Program instrumentation
3. Summary measurement collection
4. Summary analysis report examination
5. Summary experiment scoring
6. Summary measurement collection with filtering
7. Filtered summary analysis report examination
8. Event trace collection
9. Event trace examination & analysis

- Score-P measurements are configured via environment variables:

```
% scorep-info config-vars --full
SCOREP_ENABLE_PROFILING
  Description: Enable profiling
 [...]
SCOREP_ENABLE_TRACING
  Description: Enable tracing
 [...]
SCOREP_TOTAL_MEMORY
  Description: Total memory in bytes for the measurement system
 [...]
SCOREP_EXPERIMENT_DIRECTORY
  Description: Name of the experiment directory
 [...]
SCOREP_FILTERING_FILE
  Description: A file name which contain the filter rules
 [...]
SCOREP_METRIC_PAPI
  Description: PAPI metric names to measure
 [...]
SCOREP_METRIC_RUSAGE
  Description: Resource usage metric names to measure
 [... More configuration variables ...]
```

- Change to the directory containing the new executable adjust configuration and run application

```
% cd bin.scorep
% export SCOREP_EXPERIMENT_DIRECTORY=scorep_bt-mz_W_4x4_sum
% OMP_NUM_THREADS=4  mpiexec -np 4 ./bt-mz_W.4
 NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark
 Number of zones:   4 x   4
 Iterations: 200   dt:   0.000800
 Number of active processes:     4
 Use the default load factors with threads
 Total number of threads:     16  (  4.0 threads/process)
 Use the default load factors with threads

 Time step    1
 Time step   20
  [...]
 Time step  180
 Time step  200
 Verification Successful

 BT-MZ Benchmark Completed.
 Time in seconds = 54.39
```

- Creates experiment directory ./scorep_bt-mz_W_4x4_sum containing
  - a record of the measurement configuration (scorep.cfg)
  - the analysis report that was collated after measurement (profile.cubex)

```
% ls
...  scorep_bt-mz_W_4x4_sum
% ls scorep_bt-mz_W_4x4_sum
profile.cubex  scorep.cfg
```

- Interactive exploration with CUBE / ParaProf

```
% cube scorep_bt-mz_W_4x4_sum/profile.cubex

          [CUBE GUI showing summary analysis report]

% paraprof scorep_bt-mz_W_4x4_sum/profile.cubex

      [TAU ParaProf GUI showing summary analysis report]
```

# Analysis report examination
# with CUBE

Markus Geimer

Jülich Supercomputing Centre

- Parallel program analysis report exploration tools
  - Libraries for XML report reading & writing
  - Algebra utilities for report processing
  - GUI for interactive analysis exploration
    - requires Qt4
- Originally developed as part of Scalasca toolset
- Now available as a separate component
  - Can be installed independently of Score-P, e.g., on laptop or desktop
  - Latest release: CUBE 4.2.3 (June 2014)

- Representation of values (severity matrix) on three hierarchical axes
  - Performance property (metric)
  - Call path (program location)
  - System location (process/thread)

- Three coupled tree browsers

- CUBE displays severities
  - As value: for precise comparison
  - As colour: for easy identification of hotspots
  - Inclusive value when closed & exclusive value when expanded
  - Customizable via display modes

# Analysis presentation

1. Reference preparation for validation
2. Program instrumentation
3. Summary measurement collection
4. Summary analysis report examination
5. Summary experiment scoring
6. Summary measurement collection with filtering
7. Filtered summary analysis report examination
8. Event trace collection
9. Event trace examination & analysis

# Analysis report exploration (opening view)

# Metric selection

# Expanding the call tree

- ## Inclusive
  - Information of all sub-elements aggregated into single value
- ## Exclusive
  - Information cannot be subdivided further



```
int foo()
{
    int a;
    a = 1 + 1;

    bar();

    a = a + 1;
    return a;
}
```

Inclusive

Exclusive

# Selecting a call path

# Source-code view



/home/geimer/Projects/Tests/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f

```fortran
      subroutine binvcrhs( lhs,c,r )

c---------------------------------------------------------------
c---------------------------------------------------------------

c---------------------------------------------------------------
c
c---------------------------------------------------------------

      implicit none

      double precision pivot, coeff, lhs
      dimension lhs(5,5)
      double precision c(5,5), r(5)


c---------------------------------------------------------------
c
c---------------------------------------------------------------

      pivot = 1.00d0/lhs(1,1)
      lhs(1,2) = lhs(1,2)*pivot
      lhs(1,3) = lhs(1,3)*pivot
      lhs(1,4) = lhs(1,4)*pivot
      lhs(1,5) = lhs(1,5)*pivot
      c(1,1) = c(1,1)*pivot
      c(1,2) = c(1,2)*pivot
      c(1,3) = c(1,3)*pivot
      c(1,4) = c(1,4)*pivot
```

○ Read only     Save     Save as     Font...     Close

# Flat profile view

Box plot shows distribution across the system; with min/max/avg/median/quartiles

- ## Absolute
  - Absolute value shown in seconds/bytes/counts

- ## Selection percent
  - Value shown as percentage w.r.t. the selected node "on the left" (metric/call path)

- ## Peer percent (system tree only)
  - Value shown as percentage relative to the maximum peer value

# Context-sensitive help

- ## Extracting solver sub-tree from analysis report

```
% cube_cut  -r '<<ITERATION>>'  scorep_bt-mz_W_4x4_sum/profile.cubex
Writing cut.cubex... done.
```

- ## Calculating difference of two reports

```
% cube_diff  scorep_bt-mz_W_4x4_sum/profile.cubex  cut.cubex
Writing diff.cubex... done.
```

- ## Additional utilities for merging, calculating mean, etc.
  - Default output of cube_*utility* is a new report *utility*.cubex
- ## Further utilities for report scoring & statistics
- ## Run utility with "-h" (or no arguments) for brief usage info

- CUBE
  - Parallel program analysis report exploration tools
    - Libraries for XML report reading & writing
    - Algebra utilities for report processing
    - GUI for interactive analysis exploration
  - Available under New BSD open-source license
  - Documentation & sources:
    - http://www.scalasca.org
  - User guide also part of installation:
    - `cube-config --cube-dir`/share/doc/CubeGuide.pdf
  - Contact:
    - mailto: scalasca@fz-juelich.de

# Score-P hands-on:
# NPB-MZ-MPI / BT (filtered)

1. Reference preparation for validation
2. Program instrumentation
3. Summary measurement collection
4. Summary analysis report examination
5. Summary experiment scoring
6. Summary measurement collection with filtering
7. Filtered summary analysis report examination
8. Event trace collection
9. Event trace examination & analysis

- 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

- ## Report scoring as textual output

```
% scorep-score scorep_bt-mz_W_4x4_sum/profile.cubex
Estimated aggregate size of event trace:                    1025MB
Estimated requirements for largest trace buffer (max_buf):  265MB
Estimated memory requirements (SCOREP_TOTAL_MEMORY):        273MB
(hint: When tracing set SCOREP_TOTAL_MEMORY=273MB to avoid intermedi...
 or reduce requirements using USR regions

flt     type  max_buf[B]      visits  time[s]  time[%]  time/visit[us]  region
        ALL  277,799,918  41,157,533    91.76    100.0            2.23  ALL
        USR  274,792,492  40,418,321    11.38     12.4            0.28  USR
        OMP    6,882,860     685,952    51.42     56.0           74.96  OMP
        COM      371,956      45,944    15.20     16.6          330.81  COM
        MPI      102,286       7,316    13.76     15.0         1880.84  MPI
```
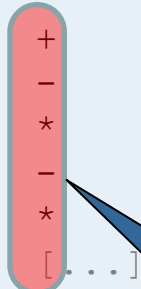
**1 GB total memory (265 MB per rank)!**

**High visit count, but very low time/visit ration!**

- ## Region/callpath classification
  - MPI (pure MPI library functions)
  - OMP (pure OpenMP functions/regions)
  - USR (user-level source local computation)
  - COM ("combined" USR + OpenMP/MPI)
  - ANY/ALL (aggregate of all region types)

- ## Score report breakdown by region

```
% scorep-score -r scorep_bt-mz_W_4x4_sum/profile.cubex
  [...]
flt     type   max_buf[B]      visits time[s] time[%] time/visit[us]   region
        USR   85,774,338 12,516,672    3.56     3.9            0.28   matmul_sub_
        USR   85,774,338 12,516,672    2.87     3.1            0.23   matvec_sub_
        USR   85,774,338 12,516,672    4.15     4.5            0.33   binvcrhs_
        USR    7,974,876  1,170,624    0.34     0.4            0.29   lhsinit_
        USR    7,974,876  1,170,624    0.32     0.3            0.27   binvrhs
        USR    3,473,912    526,848    0.14     0.1            0.26   exact_solution...
        OMP      410,040     25,728    0.01     0.0            0.50   !$omp parallel...
        OMP      410,040     25,728    0.01     0.0            0.49   !$omp parallel...
        OMP      410,040        728    0.01     0.0            0.48   !$omp parallel...
        OMP      410,040                         0.0            0.47   !$omp parallel...
        OMP      209,040                         0.0            0.98   !$omp do @exch...
        OMP      209,040                         0.0            0.97   !$omp do @exch...
        OMP      209,040                         0.3            9.69   !$omp implicit...
        OMP      209,040                         0.3            9.66   !$omp implicit...
        OMP      209,040                         0.0
        OMP      209,040     25,728    0.24     0.3
        OMP      209,040     25,728    0.02     0.0
  [...]
```

> More than 270 MB just for these 6 regions

- ## Summary measurement analysis score reveals
  - Total size of event trace would be ~1 GB
  - Maximum trace buffer size would be ~265 MB per rank
    - smaller buffer would require flushes to disk during measurement resulting in substantial perturbation
  - 99.8% 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 12.4% of total time
    - however, much of that is very likely to be measurement overhead for frequently-executed small routines (high visit count but very low time/visit ratio)

- ## Advisable to tune measurement configuration
  - Specify an adequate trace buffer size
  - Specify a filter file listing (USR) regions not to be measured

1. Reference preparation for validation
2. Program instrumentation
3. Summary measurement collection
4. Summary analysis report examination
5. Summary experiment scoring
6. Summary measurement collection with filtering
7. Filtered summary analysis report examination
8. Event trace collection
9. Event trace examination & analysis

- ## Report scoring with prospective filter listing 6 USR regions

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

% scorep-score -f ../config/scorep.filt scorep_bt-mz_W_4x4_sum/profile.cubex
Estimated aggregate size of event trace:                    23MB
Estimated requirements for largest trace buffer (max_buf):  8MB
Estimated memory requirements (SCOREP_TOTAL_MEMORY):        16MB
(hint: When tracing set SCOREP_TOTAL_MEMORY=16MB to avoid intermediate flushes
 or reduce requirements using USR regions filters.)
```

> 23 MB of memory in total,
> 8 MB per rank!

- ## Score report breakdown by region

```
% scorep-score -r -f ../config/scorep.filt \
> scorep_bt-mz_W_4x4_sum/profile.cubex
flt      type   max_buf[B]       visits  time[s]  time[%]  time/visit[us]  region
 -        ALL  277,799,918  41,157,533    91.76    100.0            2.23  ALL
 -        USR  274,792,492  40,418,321    11.38     12.4            0.28  USR
 -        OMP    6,882,860     685,952    51.42     56.0           74.96  OMP
 -        COM      371,956      45,944    15.20     16.6          330.81  COM
 -        MPI      102,286       7,316    13.76     15.0         1880.84  MPI

 *        ALL    7,357,804     739,321    80.38     87.6          108.72  ALL-FLT
 +        FLT  274,791,764  40,418,212    11.37     12.4            0.28  FLT
 -        OMP    6,882,860     685,952    51.42     56.0           74.96  OMP-FLT
 *        COM      371,956      45,944    15.20     16.6          330.81  COM-FLT
 -        MPI      102,286       7,316    13.76     15.0         1880.84  MPI-FLT
 *        USR          728         109     0.00      0.0            2.38  USR-FLT
[...]
```

Filtered routines marked with '+'

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

```
% export SCOREP_EXPERIMENT_DIRECTORY=scorep_bt-mz_W_4x4_sum_filtered
% export SCOREP_FILTERING_FILE=../config/scorep.filt
% OMP_NUM_THREADS=4  mpiexec -np 4 ./bt-mz_W.4
 NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark
 Number of zones:   4 x   4
 Iterations: 200   dt:   0.000800
 Number of active processes:     4
 Use the default load factors with threads
 Total number of threads:     16  (  4.0 threads/process)
 Use the default load factors with threads

 Time step    1
 Time step   20
  [...]
 Time step  180
 Time step  200
 Verification Successful

 BT-MZ Benchmark Completed.
 Time in seconds = 8.11
```

- ## Scoring of new analysis report as textual output

```
% scorep-score scorep_bt-mz_W_4x4_sum_filtered/profile.cubex
Estimated aggregate size of event trace:                     23MB
Estimated requirements for largest trace buffer (max_buf): 8MB
Estimated memory requirements (SCOREP_TOTAL_MEMORY):      16MB
(hint: When tracing set SCOREP_TOTAL_MEMORY=16MB to avoid intermediate flushes
 or reduce requirements using USR regions filters.)

flt     type max_buf[B]  visits time[s] time[%] time/visit[us]   region
        ALL  7,357,804 739,321   25.32   100.0          34.25   ALL
        OMP  6,882,860 685,952   16.64    65.7          24.26   OMP
        COM    371,956  45,944    3.90    15.4          84.87   COM
        MPI    102,286   7,316    4.78    18.9         653.21   MPI
        USR        728     109    0.00     0.0           2.41   USR
```

- ## Significant reduction in runtime (measurement overhead)
  - Not only reduced time for USR regions, but MPI/OMP reduced too!

- ## Further measurement tuning (filtering) may be appropriate
  - e.g., use "timer_*" to filter timer_start_, timer_read_, etc.

- Parallel performance framework and toolkit
  - Supports all HPC platforms, compilers, runtime system
  - Provides portable instrumentation, measurement, analysis

## TAU Architecture

### Instrumentation

**Source**
- C, C++, Fortran
- Python, UPC, Java
- Robust parsers (PDT)

**Wrapping**
- Interposition (PMPI)
- Wrapper generation

**Linking**
- Static, dynamic
- Preloading

**Executable**
- Dynamic (Dyninst)
- Binary (Dyninst, MAQAO)

*Measurement API*

### Measurement

**Events**
- static/dynamic
- routine, basic block, loop
- threading, communication
- heterogeneous

**Profiling**
- flat, callpath, phase, parameter, snapshot
- probe, sampling, hybrid

**Tracing**
- TAU / Scalasca tracing
- Open Trace Format (OTF)

**Metadata**
- system, user-defined

*Measured data*

### Analysis

**Profiles**
- *ParaProf* parallel profile analyzer / visualizer
- *PerfDMF* parallel profile database
- *PerfExplorer* parallel profile data mining

**Tracing**
- TAU trace translation
  - OTF, SLOG-2
- Trace analysis / visualizer
  - *Vampir, Jumpshot*

**Online**
- event unification
- statistics calculation

# TAU Performance System®

- Instrumentation
  - Fortran, C++, C, UPC, Java, Python, Chapel
  - Automatic instrumentation

- Measurement and analysis support
  - MPI, OpenSHMEM, ARMCI, PGAS, DMAPP
  - pthreads, OpenMP, hybrid, other thread models
  - GPU, CUDA, OpenCL, OpenACC
  - Parallel profiling and tracing
  - Use of Score-P for native OTF2 and CUBEX generation
  - Efficient callpath proflles and trace generation using Score-P

- Analysis
  - Parallel profile analysis (ParaProf), data mining (PerfExplorer)
  - Performance database technology (TAUdb)
  - 3D profile browser

- TAU supports both sampling and direct instrumentation
- Memory debugging as well as I/O performance evaluation
- Profiling as well as tracing
- Interfaces with Score-P for more efficient measurements
- TAU's instrumentation covers:
  - Runtime library interposition (tau_exec)
  - Compiler-based instrumentation
  - PDT based Source level instrumentation: routine & loop
  - Event based sampling (TAU_SAMPLING=1)
  - Callstack unwinding with sampling (TAU_EBS_UNWIND=1)
  - OpenMP Tools Interface (OMPT, tau_exec –T ompt)
  - CUDA CUPTI, OpenCL (tau_exec  -T cupti   -cupti)

- How much time is spent in each application routine and outer *loops*? Within loops, what is the contribution of each *statement*?

- How many instructions are executed in these code regions? Floating point, Level 1 and 2 *data cache misses*, hits, branches taken?

- What is the memory usage of the code? When and where is memory allocated/de-allocated? Are there any memory leaks?

- What are the I/O characteristics of the code?  What is the peak read and write *bandwidth* of individual calls, total volume?

- What is the contribution of each *phase* of the program? What is the time wasted/spent waiting for collectives, and I/O operations in Initialization, Computation, I/O phases?

- How does the application *scale*? What is the efficiency, runtime breakdown of performance across different core counts?

- TAU supports several measurement and thread options

    Phase profiling, profiling with hardware counters, MPI library, CUDA…

    Each measurement configuration of TAU corresponds to a unique stub makefile (configuration file) and library that is generated when you configure it

- To instrument source code automatically using PDT

    Choose an appropriate TAU stub makefile in <arch>/lib:

    **% export TAU_MAKEFILE=$TAU/Makefile.tau-mpi-pdt**

    **% export TAU_OPTIONS='-optVerbose …' (see tau_compiler.sh )**

    **% export PATH=$TAU_ROOT/x86_64/bin:$PATH**

    **% export TAU=$TAU_ROOT/x86_64/lib**

    Use tau_f90.sh, tau_cxx.sh, tau_upc.sh, or tau_cc.sh as F90, C++, UPC, or C compilers respectively:

    **% mpif90 foo.f90        changes to**

    **% tau_f90.sh foo.f90**

- Set runtime environment variables, execute application and analyze performance data:

    **% pprof   (for text based profile display)        % paraprof  (for GUI)**

% module load openmpi tau; ls $TAU/Makefile.*

Makefile.tau-icpc-papi-mpi-pdt

Makefile.tau-icpc-papi-mpi-pthread-pdt

Makefile.tau-icpc-papi-ompt-mpi-pdt-openmp

Makefile.tau-mpi-pdt

Makefile.tau-papi-mpi-pdt

Makefile.tau-papi-mpi-pdt-openmp-opari-scorep

Makefile.tau-papi-mpi-pdt-scorep

Makefile.tau-papi-mpi-pthread-pdt

Makefile.tau-papi-pthread-pdt

Makefile.tau-papi-shmem-mpi-pdt

- **For an MPI+F90 application with Intel MPI, you may choose**
Makefile.tau-mpi-pdt
  - Supports MPI instrumentation & PDT for automatic source instrumentation

`% export TAU_MAKEFILE=$TAU/Makefile.tau-mpi-pdt`

% tau_f90.sh matmult.f90 -o matmult

% mpirun -np 4 ./matmult

% paraprof

% export TAU=$TAU_ROOT/x86_64/lib

% export TAU_MAKEFILE=$TAU/Makefile.tau-papi-mpi-pdt-openmp-opari-scorep

% export OMP_NUM_THREADS=10

% make CC=tau_cc.sh CXX=tau_cxx.sh F90=tau_f90.sh


% mpirun -np 4  ./matmult


% cd score*; paraprof profile.cubex &

- Installing PDT:
  - wget http://tau.uoregon.edu/pdt_lite.tgz
  - ./configure –prefix=<dir>; make ; make install
- Installing TAU:
  - wget http://tau.uoregon.edu/tau.tgz
  - ./configure –arch=x86_64 -bfd=download -pdt=<dir> -papi=<dir> ...
  - For MIC:
  - ./configure –arch=mic_linux –pdt=<dir> -pdt_c++=g++ -papi=dir …
  - make install
- Using TAU:
  - export TAU_MAKEFILE=<taudir>/x86_64/
                                    lib/Makefile.tau-<TAGS>
  - make CC=tau_cc.sh   CXX=tau_cxx.sh   F90=tau_f90.sh

# Compile-Time Options

Optional parameters for the TAU_OPTIONS environment variable:
% tau_compiler.sh

| | |
|---|---|
| -optVerbose | Turn on verbose debugging messages |
| -optCompInst | Use compiler based instrumentation |
| -optNoCompInst | Do not revert to compiler instrumentation if source instrumentation fails. |
| -optTrackIO | Wrap POSIX I/O call and calculates vol/bw of I/O operations (Requires TAU to be configured with –*iowrapper*) |
| -optTrackGOMP | Enable tracking GNU OpenMP runtime layer (used without –opari) |
| -optMemDbg | Enable runtime bounds checking (see TAU_MEMDBG_* env vars) |
| -optKeepFiles | Does not remove intermediate .pdb and .inst.* files |
| -optPreProcess | Preprocess sources (OpenMP, Fortran) before instrumentation |
| -optTauSelectFile="<file>" | Specify selective instrumentation file for *tau_instrumentor* |
| -optTauWrapFile="<file>" | Specify path to *link_options.tau* generated by *tau_gen_wrapper* |
| -optHeaderInst | Enable Instrumentation of headers |
| -optTrackUPCR | Track UPC runtime layer routines (used with tau_upc.sh) |
| -optLinking="" | Options passed to the linker. Typically $(TAU_MPI_FLIBS) $(TAU_LIBS) $(TAU_CXXLIBS) |
| -optCompile="" | Options passed to the compiler. Typically $(TAU_MPI_INCLUDE) $(TAU_INCLUDE) $(TAU_DEFS) |
| -optPdtF95Opts="" | Add options for Fortran parser in PDT (f95parse/gfparse) … |

•Optional parameters for the TAU_OPTIONS environment variable:

% tau_compiler.sh

| | |
|---|---|
| -optMICOffload | Links code for Intel MIC offloading, requires both host and MIC TAU libraries |
| -optShared (default) | Use TAU's shared library (libTAU.so) instead of static library |
| -optPdtCxxOpts="" | Options for C++ parser in PDT (cxxparse). |
| -optPdtF90Parser="" | Specify a different Fortran parser |
| -optPdtCleanscapeParser | Specify the Cleanscape Fortran parser instead of GNU gfparser |
| -optTau="" | Specify options to the tau_instrumentor |
| -optTrackDMAPP | Enable instrumentation of low-level DMAPP API calls on Cray |
| -optTrackPthread | Enable instrumentation of pthread calls |

See tau_compiler.sh for a full list of TAU_OPTIONS.

…

# Compiling Fortran Codes with TAU

- If your Fortran code uses free format in .f files (fixed is default for .f), you may use:
  % export TAU_OPTIONS= '-optPdtF95Opts="-R free" -optVerbose '

- To use the compiler based instrumentation instead of PDT (source-based):
  % export TAU_OPTIONS= '-optCompInst -optVerbose'

- If your Fortran code uses C preprocessor directives (#include, #ifdef, #endif):
  % export TAU_OPTIONS= '-optPreProcess -optVerbose -optDetectMemoryLeaks'

- To use an instrumentation specification file:
  % export TAU_OPTIONS= '-optTauSelectFile=select.tau -optVerbose -optPreProcess'
  % cat select.tau
  BEGIN_INSTRUMENT_SECTION
  loops  routine="#"
  # this statement instruments all outer loops in all routines. # is wildcard as well as comment in first column.
  END_INSTRUMENT_SECTION

# Runtime Environment Variables

| Environment Variable | Default | Description |
|---|---|---|
| TAU_TRACE | 0 | Setting to 1 turns on tracing |
| TAU_CALLPATH | 0 | Setting to 1 turns on callpath profiling |
| TAU_TRACK_MEMORY_LEAKS | 0 | Setting to 1 turns on leak detection (for use with –optMemDbg or tau_exec) |
| TAU_MEMDBG_PROTECT_ABOVE | 0 | Setting to 1 turns on bounds checking for dynamically allocated arrays. (Use with –optMemDbg or tau_exec –memory_debug). |
| TAU_CALLPATH_DEPTH | 2 | Specifies depth of callpath. Setting to 0 generates no callpath or routine information, setting to 1 generates flat profile and context events have just parent information (e.g., Heap Entry: foo) |
| TAU_SAMPLING | 1 | Setting to 1 enables event-based sampling. |
| TAU_TRACK_SIGNALS | 0 | Setting to 1 generate debugging callstack info when a program crashes |
| TAU_COMM_MATRIX | 0 | Setting to 1 generates communication matrix display using context events |
| TAU_THROTTLE | 1 | Setting to 0 turns off throttling. Enabled by default to remove instrumentation in lightweight routines that are called frequently |
| TAU_THROTTLE_NUMCALLS | 100000 | Specifies the number of calls before testing for throttling |
| TAU_THROTTLE_PERCALL | 10 | Specifies value in microseconds. Throttle a routine if it is called over 100000 times and takes less than 10 usec of inclusive time per call |
| TAU_COMPENSATE | 0 | Setting to 1 enables runtime compensation of instrumentation overhead |
| TAU_PROFILE_FORMAT | Profile | Setting to "merged" generates a single file. "snapshot" generates xml format |
| TAU_METRICS | TIME | Setting to a comma separated list generates other metrics. (e.g., TIME,P_VIRTUAL_TIME,PAPI_FP_INS,PAPI_NATIVE_<event>:<subevent>) |

| Environment Variable | Default | Description |
| --- | --- | --- |
| TAU_TRACK_MEMORY_LEAKS | 0 | Tracks allocates that were not de-allocated (needs –optMemDbg or tau_exec –memory) |
| TAU_EBS_SOURCE | TIME | Allows using PAPI hardware counters for periodic interrupts for EBS (e.g., TAU_EBS_SOURCE=PAPI_TOT_INS when TAU_SAMPLING=1) |
| TAU_EBS_PERIOD | 100000 | Specifies the overflow count for interrupts |
| TAU_MEMDBG_ALLOC_MIN/MAX | 0 | Byte size minimum and maximum subject to bounds checking (used with TAU_MEMDBG_PROTECT_*) |
| TAU_MEMDBG_OVERHEAD | 0 | Specifies the number of bytes for TAU's memory overhead for memory debugging. |
| TAU_MEMDBG_PROTECT_BELOW/ABOVE | 0 | Setting to 1 enables tracking runtime bounds checking below or above the array bounds (requires –optMemDbg while building or tau_exec –memory) |
| TAU_MEMDBG_ZERO_MALLOC | 0 | Setting to 1 enables tracking zero byte allocations as invalid memory allocations. |
| TAU_MEMDBG_PROTECT_FREE | 0 | Setting to 1 detects invalid accesses to deallocated memory that should not be referenced until it is reallocated (requires –optMemDbg or tau_exec –memory) |
| TAU_MEMDBG_ATTEMPT_CONTINUE | 0 | Setting to 1 allows TAU to record and continue execution when a memory error occurs at runtime. |
| TAU_MEMDBG_FILL_GAP | Undefined | Initial value for gap bytes |
| TAU_MEMDBG_ALINGMENT | Sizeof(int) | Byte alignment for memory allocations |
| TAU_EVENT_THRESHOLD | 0.5 | Define a threshold value (e.g., .25 is 25%) to trigger marker events for min/max |

# Simplifying TAU's usage (tau_exec)

- Uninstrumented execution
  - % mpirun -np 4  ./a.out
- Track MPI performance
  - % mpirun -np 4   tau_exec  ./a.out
- Track POSIX I/O and MPI performance (MPI enabled by default)
  - % mpirun -np 4  tau_exec –T mpi,pdt  –io  ./a.out
- Track memory operations
  - % export TAU_TRACK_MEMORY_LEAKS=1
  - % mpirun –np 8 tau_exec –memory_debug ./a.out (bounds check)
- Use event based sampling (compile with –g)
  - % mpirun –np 8 tau_exec –ebs ./a.out
  - Also –ebs_source=<PAPI_COUNTER> -ebs_period=<overflow_count>
- Load wrapper interposition library
  - % mpirun –np 8 tau_exec –loadlib=<path/libwrapper.so> ./a.out
- Track GPGPU operations
  - % mpirun –np 8 tau_exec –cupti ./a.out
  - % mpirun –np 8 tau_exec –opencl ./a.out

- Support for both static and dynamic executables
- Specify a list of routines to instrument
- Specify the TAU measurement library to be injected
- MAQAO:

```
% tau_rewrite -T [tags] a.out -o a.inst
```

- Dyninst:

```
% tau_run -T [tags] a.out -o a.inst
```

- Pebil:

```
% tau_pebil_rewrite -T [tags] a.out \
  -o a.inst
```

- Execute the application to get measurement data:

```
% mpirun -np 256 ./a.inst
```

# ParaProf Profile Analysis Framework

% export TAU_COMM_MATRIX=1

# Event Based Sampling in TAU



```
% export TAU_MAKEFILE=$TAU/Makefile.tau-icpc-papi-mpi-pdt
% make CC=tau_cc.sh CXX=tau_cxx.sh
% export TAU_SAMPLING=1
% mpirun –np 256 ./a.out
% paraprof
```

TAU: ParaProf: Statistics for: node 0, thread 0 – /usr/global/tools/tau/training/tau-2.23/examples/mm

File   Options   Windows   Help

| Name △ | Exclu... | Inclu... | C... | C... |
|---|---|---|---|---|
| int main(int, char **) C [{matmult.pomp.c} {224,1}–{294,1}] | 0 | 3.134 | 1 | 3 |
| MPI_Finalize() | 0.191 | 0.191 | 1 | 5 |
| MPI_Init_thread() | 1.148 | 1.151 | 1 | 45 |
| double do_work(void) C [{matmult.pomp.c} {185,1}–{216,1}] | 0.001 | 1.792 | 1 | 8 |
| double **allocateMatrix(int, int) C [{matmult.pomp.c} {38,1}–{45,1}] | 0.003 | 0.003 | 3 | 0 |
| void compute(double **, double **, double **, int, int, int) C [{matmult.pomp.c} {111,1}–{146,1}] | 0 | 0.97 | 1 | 1 |
| parallel fork/join [OpenMP] | 0 | 0.97 | 1 | 1 |
| parallel (parallel fork/join) [OpenMP location: file:/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c <80, 96>] | 0 | 0.97 | 1 | 1 |
| parallel begin/end [OpenMP] | 0 | 0.97 | 1 | 1 |
| parallel (parallel begin/end) [OpenMP location: file:/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c <80, 96>] | 0 | 0.97 | 1 | 2 |
| barrier enter/exit [OpenMP] | 0 | 0.003 | 1 | 1 |
| parallel (barrier enter/exit) [OpenMP location: file:/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c <80, 96>] | 0.003 | 0.003 | 1 | 0 |
| for enter/exit [OpenMP] | 0 | 0.967 | 1 | 1 |
| for (loop body) [OpenMP location: file:/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c <84, 95>] | 0.967 | 0.967 | 1 | 0 |
| [CONTEXT] for (loop body) [OpenMP location: file:/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c <84, 95>] | 0 | 0.93 | 58 | 0 |
| [SAMPLE] L_compute_118__par_region0_2_204 [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.pomp.c} {127}] | 0.039 | 0.039 | 3 | 0 |
| [SAMPLE] L_compute_118__par_region0_2_204 [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.pomp.c} {131}] | 0.891 | 0.891 | 55 | 0 |
| void compute_interchange(double **, double **, double **, int, int, int) C [{matmult.pomp.c} {148,1}–{183,1}] | 0 | 0.812 | 1 | 1 |
| parallel fork/join [OpenMP] | 0 | 0.812 | 1 | 1 |
| parallel (parallel fork/join) [OpenMP location: file:/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c <101, 117>] | 0 | 0.811 | 1 | 1 |
| parallel begin/end [OpenMP] | 0 | 0.811 | 1 | 1 |
| parallel (parallel begin/end) [OpenMP location: file:/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c <101, 117>] | 0 | 0.811 | 1 | 2 |
| void initialize(double **, int, int) C [{matmult_initialize.pomp.c} {5,1}–{33,1}] | 0 | 0.007 | 3 | ? |

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-icpc-papi-mpi-pdt-opari-openmp
% make CC=tau_cc.sh CXX=tau_cxx.sh
% export TAU_SAMPLING=1; export OMP_NUM_THREADS=16
% mpirun –np 256 ./a.out
% paraprof
```

# TAU's Support for Intel OMPT



TAU: ParaProf: Statistics for: node 0, thread 0 – /usr/global/tools/tau/training/tau-2.23/examples/mm

| Name △ | Exclusi... | Inclusi... | Calls | Child Calls |
|---|---|---|---|---|
| .TAU application | 0 | 2.943 | 1 | 1 |
| int main(int, char **) C [{matmult.c} {159,1}–{229,1}] | 0 | 2.942 | 1 | 3 |
| MPI_Finalize() | 0.034 | 0.034 | 1 | 5 |
| MPI_Init_thread() | 1.148 | 1.151 | 1 | 45 |
| double do_work(void) C [{matmult.c} {120,1}–{151,1}] | 0 | 1.757 | 1 | 8 |
| double **allocateMatrix(int, int) C [{matmult.c} {36,1}–{43,1}] | 0.003 | 0.003 | 3 | 0 |
| void compute(double **, double **, double **, int, int, int) C [{matmult.c} {78,1}–{97,1}] | 0 | 0.951 | 1 | 1 |
| OpenMP_PARALLEL_REGION: [OPENMP] compute [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c} {80}] | 0 | 0.951 | 1 | 2 |
| OpenMP_BARRIER: [OPENMP] compute [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c} {80}] | 0.004 | 0.004 | 1 | 0 |
| [CONTEXT] OpenMP_BARRIER: [OPENMP] compute [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c} {80}] | 0 | 0.007 | 1 | 0 |
| OpenMP_LOOP: [OPENMP] compute [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c} {80}] | 0.946 | 0.946 | 1 | 0 |
| [CONTEXT] OpenMP_LOOP: [OPENMP] compute [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c} {80}] | 0 | 0.944 | 62 | 0 |
| [SAMPLE] L_compute_80__par_region0_2_150 [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c} {87}] | 0.047 | 0.047 | 2 | 0 |
| [SAMPLE] L_compute_80__par_region0_2_150 [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult.c} {91}] | 0.897 | 0.897 | 60 | 0 |
| void compute_interchange(double **, double **, double **, int, int, int) C [{matmult.c} {99,1}–{118,1}] | 0 | 0.787 | 1 | 1 |
| void initialize(double **, int, int) C [{matmult_initialize.c} {3,1}–{16,1}] | 0.005 | 0.016 | 3 | 3 |
| OpenMP_PARALLEL_REGION: [OPENMP] initialize [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult_initialize.c} {5}] | 0 | 0.011 | 3 | 6 |
| OpenMP_BARRIER: [OPENMP] initialize [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult_initialize.c} {5}] | 0.01 | 0.01 | 3 | 0 |
| [CONTEXT] OpenMP_BARRIER: [OPENMP] initialize [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult_initializ | 0 | 0.03 | 2 | 0 |
| OpenMP_LOOP: [OPENMP] initialize [{/usr/global/tools/tau/training/tau-2.23/examples/mm/matmult_initialize.c} {5}] | 0.001 | 0.001 | 3 | 0 |

```
% export TAU_MAKEFILE=$TAU/Makefile.tau-icpc-papi-mpi-pdt-ompt-openmp
% make CC=tau_cc.sh CXX=tau_cxx.sh
% export TAU_SAMPLING=1; export OMP_NUM_THREADS=16
% mpirun –np 256 tau_exec –T ompt –loadlib=$TAU/libiomp5.so  ./a.out
% paraprof
```

% export TAU_MAKEFILE=$TAUROOT/mic_linux/lib/Makefile.tau-icpc-papi-mpi-pdt
% export TAU_METRICS=TIME,PAPI_NATIVE_VPU_ELEMENTS_ACTIVE,
                     PAPI_NATIVE_VPU_INSTRUCTIONS_EXECUTED

TAU: ParaProf: Context Events for: node 2 – rapl_marker_16p.ppk

| Name △ | MaxValue | MinValue | NumSamples | MeanValue | Std. Dev. |
|---|---|---|---|---|---|
| ▼ int main(int, char **) C [{matmult.c} {165,1}–{237,1}] | | | | | |
| ▼ double do_work(void) C [{matmult.c} {126,1}–{157,1}] | | | | | |
| ▼ void compute(double **, double **, double **, int, int, int) C [{matmult.c} {84,1}–{103,1}] | | | | | |
| [GROUP=MAX_MARKER] rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts) | 17.585 | 17.469 | 5 | 17.521 | 0.037 |
| [GROUP=MAX_MARKER] rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts) | 15.261 | 15.218 | 4 | 15.237 | 0.016 |
| [GROUP=MAX_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE0 (Power in Watts) | 118.903 | 114.923 | 22 | 116.98 | 1.201 |
| [GROUP=MAX_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE1 (Power in Watts) | 113.466 | 110.207 | 22 | 111.778 | 0.996 |
| [GROUP=MAX_MARKER] rapl:::PP0_ENERGY:PACKAGE0 (Power in Watts) | 100.138 | 96.266 | 24 | 98.206 | 1.13 |
| [GROUP=MAX_MARKER] rapl:::PP0_ENERGY:PACKAGE1 (Power in Watts) | 95.846 | 92.758 | 24 | 94.319 | 0.937 |
| [GROUP=MIN_MARKER] rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts) | 17.397 | 17.303 | 4 | 17.358 | 0.035 |
| [GROUP=MIN_MARKER] rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts) | 15.048 | 15.042 | 2 | 15.045 | 0.003 |
| ▼ int mysleep(int) C [{matmult.c} {46,1}–{49,1}] | | | | | |
| [GROUP=MIN_MARKER] rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts) | 15.84 | 15.84 | 1 | 15.84 | 0 |
| [GROUP=MIN_MARKER] rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts) | 14.275 | 14.275 | 1 | 14.275 | 0 |
| [GROUP=MIN_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE1 (Power in Watts) | 96.853 | 96.853 | 1 | 96.853 | 0 |
| [GROUP=MIN_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE0 (Power in Watts) | 93.125 | 93.125 | 1 | 93.125 | 0 |
| [GROUP=MIN_MARKER] rapl:::PP0_ENERGY:PACKAGE0 (Power in Watts) | 75.096 | 75.096 | 1 | 75.096 | 0 |
| [GROUP=MIN_MARKER] rapl:::PP0_ENERGY:PACKAGE1 (Power in Watts) | 79.646 | 79.646 | 1 | 79.646 | 0 |
| ▼ void compute_interchange(double **, double **, double **, int, int, int) C [{matmult.c} {105,1}–{124,1}] | | | | | |
| [GROUP=MAX_MARKER] rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts) | 26.064 | 25.711 | 2 | 25.887 | 0.176 |
| [GROUP=MAX_MARKER] rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts) | 24.373 | 23.965 | 4 | 24.232 | 0.159 |
| [GROUP=MAX_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE0 (Power in Watts) | 126.872 | 125.182 | 6 | 125.732 | 0.557 |
| [GROUP=MAX_MARKER] rapl:::PACKAGE_ENERGY:PACKAGE1 (Power in Watts) | 124.377 | 116.689 | 5 | 122.428 | 2.885 |
| [GROUP=MAX_MARKER] rapl:::PP0_ENERGY:PACKAGE0 (Power in Watts) | 103.981 | 102.21 | 6 | 102.769 | 0.584 |
| [GROUP=MAX_MARKER] rapl:::PP0_ENERGY:PACKAGE1 (Power in Watts) | 102.615 | 101.693 | 4 | 102.115 | 0.33 |
| rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts) | 26.064 | 15.84 | 36 | 19.053 | 3.39 |
| rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts) | 24.373 | 14.275 | 36 | 16.435 | 3.155 |
| rapl:::PACKAGE_ENERGY:PACKAGE0 (Power in Watts) | 126.872 | 93.125 | 36 | 117.729 | 5.403 |
| rapl:::PACKAGE_ENERGY:PACKAGE1 (Power in Watts) | 124.377 | 96.853 | 36 | 112.961 | 4.776 |
| rapl:::PP0_ENERGY:PACKAGE0 (Power in Watts) | 103.981 | 75.096 | 36 | 98.208 | 4.466 |
| rapl:::PP0_ENERGY:PACKAGE1 (Power in Watts) | 102.615 | 79.646 | 36 | 94.872 | 3.662 |

% export TAU_EVENT_THRESHOLD 0.5

| TAU: ParaProf: Context Events for: node 0 – power_rapl_16p.ppk | | | | | | |
|---|---|---|---|---|---|---|
| Name △ | .. | NumSamples | MaxValue | MinValue | MeanValue | Std. Dev. |
| rapl:::DRAM_ENERGY:PACKAGE0 (Power in Watts) | … | 39 | 26.338 | 17.673 | 20.514 | 3.031 |
| rapl:::DRAM_ENERGY:PACKAGE1 (Power in Watts) | … | 39 | 21.29 | 11.462 | 14.182 | 2.275 |
| rapl:::PACKAGE_ENERGY:PACKAGE0 (Power in Watts) | … | 39 | 120.72 | 80.665 | 114.013 | 7.222 |
| rapl:::PACKAGE_ENERGY:PACKAGE1 (Power in Watts) | … | 39 | 117.461 | 77.758 | 106.762 | 6.365 |
| rapl:::PP0_ENERGY:PACKAGE0 (Power in Watts) | … | 39 | 97.753 | 61.266 | 93.768 | 7.06 |
| rapl:::PP0_ENERGY:PACKAGE1 (Power in Watts) | … | 39 | 97.987 | 61.895 | 89.887 | 6.12 |

```
#include <TAU.h>

TAU_TRACK_POWER();   // In Fortran: call TAU_TRACK_POWER()

% sudo chmod –R go+r /dev/cpu/*/msr
% sudo /sbin/setcap cap_sys_rawio=ep ./a.out
% unset LD_LIBRARY_PATH
% ldd ./a.out
should have no "not found" entries, Use –Wl,-rpath,/path while linking
% ./a.out
% paraprof
```

- The Tutorial contains Score-P experiments of BT-MZ
  - class "B", 4 processes with 4 OpenMP threads each
  - collected on a dedicated node of the SuperMUC HPC system at Leibniz Rechenzentrum (LRZ), Munich, Germany

```
% cd
% ls
periscope-1.5                              scorep_bt-mz_B_4x4_sum
README                                     scorep_bt-mz_B_4x4_sum+mets
run.out                                    scorep_bt-mz_B_4x4_trace
scorep-20120913_1740_557443655223384
```

- Start TAU's paraprof GUI with default profile report

```
% paraprof scorep-20120913_1740_557443655223384/profile.cubex
OR
% paraprof scorep_bt-mz_B_4x4_trace/scout.cubex
```

Each color represents an event executing on one or more threads

# ParaProf: Thread Statistics Table



TAU: ParaProf: Statistics for: node 0, thread 0 – scout.cubex

File   Options   Windows   Help

Time ▼

| Name | Exclusive Time ▽ | Inclusive Time | Calls | Child Calls |
|---|---|---|---|---|
| !$omp do @y_solve.f:52 | 5.81? | 5.817 | 3,216 | 0 |
| !$omp do @z_solve.f:52 | 5.657 | 5.657 | 3,216 | 0 |
| !$omp do @x_solve.f:54 | 5.609 | 5.609 | 3,216 | 0 |
| !$omp do @rhs.f:191 | 0.609 | 0.609 | 3,232 | 0 |
| !$omp do @rhs.f:80 | 0.583 | ?83 | 3,232 | 0 |
| MPI_Waitall | 0.402 | | | 0 |
| !$omp implicit barrier | 0.402 | | | |
| !$omp do @rhs.f:301 | 0.36 | | | |
| !$omp implicit barrier | 0.026 | | | |
| !$omp implicit barrier | 0 | | | |
| !$omp do @rhs.f:37 | 0.343 | | | |
| !$omp do @rhs.f:62 | 0.225 | 0.228 | 3,232 | 3,232 |
| !$omp implicit barrier | 0.004 | 0.004 | 3,216 | 0 |
| !$omp implicit barrier | 0 | 0 | 16 | 0 |
| MPI_Init_thread | 0.218 | 0.218 | 1 | 0 |
| !$omp do @rhs.f:384 | 0.199 | 0.199 | 3,232 | 0 |
| !$omp parallel do @add.f:22 | 0.099 | 0.111 | 3,216 | 3,216 |
| !$omp do @rhs.f:428 | 0.069 | 0.069 | 3,232 | 0 |
| MPI_Isend | 0.043 | 0.043 | 603 | 0 |
| !$omp do @initialize.f:50 | 0.04 | 0.04 | 32 | 0 |
| !$omp parallel @rhs.f:28 | 0.03 | 2.536 | 3,232 | 51,712 |
| !$omp parallel do @exch_qbc.f:215 | 0.021 | 0.029 | 6,432 | 6,432 |
| !$omp parallel do @exch_qbc.f:255 | 0.02 | 0.033 | 6,432 | 6,432 |
| !$omp parallel @exch_qbc.f:255 | 0.02 | 0.053 | 6,432 | 6,432 |
| !$omp parallel @exch_qbc.f:244 | | | | |

Click to sort by a given metric, drag and move to rearrange columns

FinderScreenSnapz003.png

# Example: Score-P with TAU (NPB LU)

# ParaProf: Thread Callgraph Window



Click on options to choose a different color or to resize the box based on metrics

## TAU: ParaProf: Call Path Data n,c,t, 0,0,0 — scout.cubex

File   Options   Windows   Help

Metric Name: Time
Sorted By: Exclusive
Units: seconds

```
        0.04          0.04       32/32           !$omp parallel @initialize.f:28
 -->    0.04          0.04       32              !$omp do @initialize.f:50


        0.03          2.536      3232/3232       compute_rhs_
 -->    0.03          2.536      3232            !$omp parallel @rhs.f:28
        9.8E-4        9.8E-4     3232/3232       !$omp master @rhs.f:424
        0.225         0.228      3232/3232       !$omp do @rhs.f:62
        0.002         0.002      3232/3232       !$omp master @rhs.f:74
        0.002         0.002      3232/3232       !$omp master @rhs.f:293
        0.199         0.199      3232/3232       !$omp do @rhs.f:384
        0.002         0.002      3232/3232       !$omp master @rhs.f:183
        0.343         0.343      3232/3232       !$omp do @rhs.f:37
        0.016         0.016      3232/3232       !$omp do @rhs.f:372
        0.014         0.027      3232/3232       !$omp do @rhs.f:413
        0.609         0.609      3232/3232       !$omp do @rhs.f:191
        0.36          0.386      3232/3232       !$omp do @rhs.f:301
        0.583         0.583      3232/3232       !$omp do @rhs.f:80
        0.019         0.019      3232/3232       !$omp do @rhs.f:400
        0.006         0.006      3232/51680      !$omp implicit barrier
        0.069         0.069      3232/3232       !$omp do @rhs.f:428
        0.015         0.015      3232/3232       !$omp do @rhs.f:359


        0.021         0.029      6432/6432       !$omp parallel @exch_qbc.f:215
 -->    0.021         0.029      6432            !$omp parallel do @exch_qbc.f:215
        0.007         0.007      6432/51680      !$omp implicit barrier


        0.02          0.033      6432/6432       !$omp parallel @exch_qbc.f:255
 -->    0.02          0.033      6432            !$omp parallel do @exch_qbc.f:255
        0.013         0.013      6432/51680      !$omp implicit barrier
```

# ParaProf: Topology 3D View

# ParaProf: Node View



TAU: ParaProf: node 0, thread 0 – profile.cubex

File  Options  Windows  Help

Metric: Time
Value: Exclusive
Units: seconds

| Value | Call path |
|-------|-----------|
| 3.71 | MAIN__ => adi_ => y_solve_ => !$omp parallel @y_solve.f:43 => !$omp do @y_solve.f:52 |
| 3.71 | !$omp do @y_solve.f:52 |
| 3.593 | MAIN__ => adi_ => z_solve_ => !$omp parallel @z_solve.f:43 => !$omp do @z_solve.f:52 |
| 3.593 | !$omp do @z_solve.f:52 |
| 3.55 | MAIN__ => adi_ => x_solve_ => !$omp parallel @x_solve.f:46 => !$omp do @x_solve.f:54 |
| 3.55 | !$omp do @x_solve.f:54 |
| 0.4 | !$omp do @rhs.f:191 |
| 0.398 | MAIN__ => adi_ => compute_rhs_ => !$omp parallel @rhs.f:28 => !$omp do @rhs.f:191 |
| 0.383 | !$omp do @rhs.f:80 |
| 0.381 | MAIN__ => adi_ => compute_rhs_ => !$omp parallel @rhs.f:28 => !$omp do @rhs.f:80 |
| 0.352 | !$omp implicit barrier |
| 0.299 | !$omp parallel @rhs.f:28 |
| 0.298 | MAIN__ => adi_ => compute_rhs_ => !$omp parallel @rhs.f:28 |
| 0.28 | !$omp do @rhs.f:37 |
| 0.279 | MAIN__ => adi_ => compute_rhs_ => !$omp parallel @rhs.f:28 => !$omp do @rhs.f:37 |
| 0.261 | !$omp do @rhs.f:301 |
| 0.259 | MAIN__ => adi_ => compute_rhs_ => !$omp parallel @rhs.f:28 => !$omp do @rhs.f:301 |
| 0.228 | !$omp do @rhs.f:62 |
| 0.227 | MAIN__ => adi_ => compute_rhs_ => !$omp parallel @rhs.f:28 => !$omp do @rhs.f:62 |
| 0.214 | MAIN__ => mpi_setup_ => MPI_Init_thread |
| 0.214 | MPI_Init_thread |
| 0.161 | MAIN__ => exch_qbc_ => copy_x_face_ |
| 0.161 | copy_x_face_ |
| 0.16 | MAIN__ => exch_qbc_ => copy_y_face_ |
| 0.16 | copy_y_face_ |
| 0.15 | MAIN__ => exch_qbc_ |
| 0.15 | exch_qbc_ |
| 0.141 | !$omp do @rhs.f:384 |
| 0.14 | MAIN__ => adi_ => compute_rhs_ => !$omp parallel @rhs.f:28 => !$omp do @rhs.f:384 |
| 0.127 | MAIN__ => exch_qbc_ => MPI_Waitall |
| 0.127 | MPI_Waitall |
| 0.103 | MAIN__ => adi_ |
| 0.103 | adi_ |
| 0.094 | MAIN__ => adi_ => add_ => !$omp parallel @add.f:22 => !$omp parallel do @add.f:22 |
| 0.094 | !$omp parallel do @add.f:22 |

# ParaProf: Add Thread to Comparison Window

# ParaProf: Group Changer Window

**http://tau.uoregon.edu**

**http://www.hpclinux.com [LiveDVD, OVA]**

**Free download, open source, BSD license**

# Automatic trace analysis
# with Scalasca

Markus Geimer
Jülich Supercomputing Centre

- ## 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
  - Follow-up to pioneering KOJAK project (started 1998)
- 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 v2.1 (August 2014)

- Open source, BSD 3-clause license
- Fairly portable
  - IBM Blue Gene, IBM SP & blade clusters, Cray XT/XE/XK/XC, SGI Altix, Solaris & Linux clusters, Fujitsu FX10 & K computer, ...
- Uses Score-P instrumenter & measurement libraries
  - Scalasca 2.1 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

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

- 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

- ## One command for (almost) everything…

```
% scalasca
Scalasca 2.1
Toolset for scalable performance analysis of large-scale 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

  -c, --show-config  show configuration and exit
  -h, --help         show this help and exit
  -n, --dry-run      show actions without taking them
      --quickref     show quick reference guide and exit
  -v, --verbose      enable verbose commentary
  -V, --version      show version information and exit
```

- # Scalasca application instrumenter

```
% skin
Scalasca 2.1: application instrumenter using scorep
usage: skin [-v] [-comp] [-pdt] [-pomp] [-user] <compile-or-link-cmd>
   -comp={all|none|...}: routines to be instrumented by compiler
           (... custom instrumentation specification for compiler)
   -pdt:  process source files with PDT 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 Score-P instrumenter
```

- Deprecated command
  - Provides compatibility with Scalasca 1.x
  - Prints corresponding Score-P instrumenter command
  - Helps in transitioning existing configurations
- Recommended: use Score-P instrumenter directly

- ## Scalasca measurement collection & analysis nexus

```
% scan
Scalasca 2.1: 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.
```

- ## Scalasca analysis report explorer

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

  -c <none|quick|full>: Level of sanity checks for newly created reports
  -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_W_4x4_sum
% OMP_NUM_THREADS=4 scan mpiexec -np 4 ./bt-mz_W.4
S=C=A=N: Scalasca 2.1 runtime summarization
S=C=A=N: ./scorep_bt-mz_W_4x4_sum experiment archive
S=C=A=N: Thu Jun 12 18:05:17 2014: Collect start
mpiexec -np 4 ./bt-mz_W.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

 [... More application output ...]

S=C=A=N: Thu Jun 12 18:05:39 2014: Collect done (status=0) 22s
S=C=A=N: ./scorep_bt-mz_W_4x4_sum complete.
```

- Creates experiment directory ./scorep_bt-mz_W_4x4_sum

- ## Score summary analysis report

```
% square -s  scorep_bt-mz_W_4x4_sum
INFO: Post-processing runtime summarization result...
INFO: Score report written to ./scorep_bt-mz_W_4x4_sum/scorep.score
```

- ## Post-processing and interactive exploration with CUBE

```
% square  scorep_bt-mz_W_4x4_sum
INFO: Displaying ./scorep_bt-mz_W_4x4_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

1. Reference preparation for validation
2. Program instrumentation
3. Summary measurement collection
4. Summary analysis report examination
5. Summary experiment scoring
6. Summary measurement collection with filtering
7. Filtered summary analysis report examination
8. Event trace collection
9. Event trace examination & analysis

- Re-run the application using Scalasca nexus with **"-t"** flag

```
% export SCOREP_EXPERIMENT_DIRECTORY=scorep_bt-mz_W_4x4_trace
% OMP_NUM_THREADS=4 scan -t mpiexec –np 4 ./bt-mz_W.4
S=C=A=N: Scalasca 2.1 trace collection and analysis
S=C=A=N: ./scorep_bt-mz_W_4x4_trace experiment archive
S=C=A=N: Thu Jun 12 18:05:39 2014: Collect start
mpiexec –np 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

 [... More application output ...]

S=C=A=N: Thu Jun 12 18:05:58 2014: Collect done (status=0) 19s
 [... continued ...]
```

- ## Continues with automatic (parallel) analysis of trace files

```
S=C=A=N: Thu Jun 12 18:05:58 2014: Analyze start
mpiexec -np 4 scout.hyb ./scorep_bt-mz_W_4x4_trace/traces.otf2
SCOUT    Copyright (c) 1998-2012 Forschungszentrum Juelich GmbH
         Copyright (c) 2009-2012 German Research School for Simulation
                                 Sciences GmbH

Analyzing experiment archive ./scorep_bt-mz_W_4x4_trace/traces.otf2

Opening experiment archive ... done (0.002s).
Reading definition data    ... done (0.004s).
Reading event trace data   ... done (0.130s).
Preprocessing              ... done (0.259s).
Analyzing trace data       ...
  Wait-state detection (fwd)      (1/4) ... done (0.575s).
  Wait-state detection (bwd)      (2/4) ... done (0.138s).
  Synchpoint exchange             (3/4) ... done (0.358s).
  Critical-path analysis          (4/4) ... done (0.288s).
done (1.360s).
Writing analysis report    ... done (0.121s).

Total processing time      : 1.924s
S=C=A=N: Thu Jun 12 18:06:00 2014: Analyze done (status=0) 2s
```

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

```
% square  scorep_bt-mz_W_4x4_trace
INFO: Post-processing runtime summarization result...
INFO: Post-processing trace analysis report...
INFO: Displaying ./scorep_bt-mz_W_4x4_trace/trace.cubex...

              [GUI showing trace analysis report]
```

Additional trace-based metrics in metric hierarchy

Access online metric description via context menu

**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 imbalance highlights inefficient parallelism

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

…and select trace file from the experiment directory

# Show most severe pattern instances



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

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

**Website:** www.scalasca.org
**User support:** scalasca@fz-juelich.de

# Performance Analysis with Vampir

Bert Wesarg, Andreas Knüpfer

ZIH, Technische Universität Dresden

- Visualization of dynamics of complex parallel processes

- Full details for arbitrary temporal and spatial levels

- Supplement to automatic analysis

**Typical questions that Vampir helps to answer:**

- What happens in my application execution during a given time in a given process or thread?

- How do the communication patterns of my application execute on a real system?

- Are there any imbalances in computation, I/O or memory usage and how do they affect the parallel execution of my application?

- ## Directly on front end or local machine

```
% vampir
```



Small/Medium sized trace

Thread parallel analysis

- ## On local machine with remote VampirServer

- Vampir & VampirServer
  - Interactive trace visualization and analysis
  - Intuitive browsing and zooming
  - Scalable to large trace data sizes (20 TByte)
  - Scalable to high parallelism (200000 processes)

- Vampir is available for Linux, Windows and Mac OS X

- ## **Timeline Charts:**
  - Master Timeline
  - Process Timeline
  - Counter Data Timeline
  - Performance Radar

- ## **Summary Charts:**
  - Function Summary
  - Message Summary
  - Process Summary
  - Communication Matrix View



Show application activities and communication along a time axis



Provide quantitative results for the currently selected time interval

# Vampir hands-on

Visualizing and analyzing NPB-MZ-MPI / BT

1. Reference preparation for validation
2. Program instrumentation
3. Summary measurement collection
4. Summary analysis report examination
5. Summary experiment scoring
6. Summary measurement collection with filtering
7. Filtered summary analysis report examination
8. Event trace collection
9. Event trace examination & analysis

# Vampir: Visualization of the NPB-MZ-MPI / BT trace

## Master Timeline



Detailed information about functions, communication and synchronization events for collection of processes.

# Process Timeline



> Detailed information about different levels of function calls in a stacked bar chart for an individual process.

# Typical program phases



Initialization Phase

Computation Phase

## Counter Data Timeline



Detailed counter information over time for an individual process.

# Vampir: Visualization of the NPB-MZ-MPI / BT trace

## Performance Radar



Detailed counter information over time for a collection of processes.

# Zoom in: Initialization Phase



Context View:
Detailed information about function "initialize_".

# Feature: Find Function



> Execution of function "initialize_" results in higher page fault rates.

# Computation Phase



Computation phase results in higher floating point operations.

## Zoom in: Computation Phase



MPI communication results in lower floating point operations.

## Zoom in: Finalization Phase

## Process Summary



**Function Summary**: Overview of the accumulated information across all functions and for a collection of processes.

**Process Summary**: Overview of the accumulated information across all functions and for every process independently.

## Process Summary



Find groups of similar processes and threads by using summarized function information.

Vampir is available at http://www.vampir.eu,
Get support via vampirsupport@zih.tu-dresden.de

# Hardware performance/soft counter measurements hands-on

VI-HPS Team

- If Score-P has been built with performance metric support it is capable of recording performance counter information

- Requested counters will be recorded with every enter/exit event

- Supported metric sources
  - PAPI
  - Resource usage statistics
  - Custom written metric plug-ins

> Note: Additional memory is needed to store metric values. Therefore, you may have to adjust SCOREP_TOTAL_MEMORY, for example as reported using "scorep-score -c"

- # Recording hardware counters via PAPI

```
% export SCOREP_METRIC_PAPI=PAPI_TOT_INS,PAPI_FP_INS
% OMP_NUM_THREADS=4 mpiexec –n 4 ./bt-mz_W.4

 NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark

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

- # Also possible to record them only per rank

```
% export SCOREP_METRIC_PAPI_PER_PROCESS=PAPI_L3_DCM
% OMP_NUM_THREADS=4 mpiexec –n 4 ./bt-mz_W.4

 NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark

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

- ## Available PAPI metrics
  - Preset events: common set of events deemed relevant and useful for application performance tuning
    - Abstraction from specific hardware performance counters, mapping onto available events done by PAPI internally

```
% papi_avail
```

  - Native events: set of all events that are available on the CPU (**platform dependent**)

```
% papi_native_avail
```

> Note:
> Due to hardware restrictions
> - number of concurrently measured events is limited
> - there may be unsupported combinations of concurrent events
> - Use `papi_event_chooser` tool to test event combinations

- ## Recording operating system resource usage

```
% export SCOREP_METRIC_RUSAGE=ru_stime
% OMP_NUM_THREADS=4 mpiexec –n 4 ./bt-mz_W.4

 NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark

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

- ## Also possible to record them only per rank

```
% export SCOREP_METRIC_RUSAGE_PER_PROCESS=ru_maxrss
% OMP_NUM_THREADS=4 mpiexec –n 4 ./bt-mz_W.4

 NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark

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

- ## Available resource usage metrics

> Note:
> (1) Not all fields are maintained on each platform.
> (2) Check scope of metrics (per process vs. per thread)

```
% man getrusage
 [... Output ...]

struct rusage {
    struct timeval ru_utime; /* user CPU time used */
    struct timeval ru_stime; /* system CPU time used */
    long   ru_maxrss;        /* maximum resident set size */
    long   ru_ixrss;         /* integral shared memory size */
    long   ru_idrss;         /* integral unshared data size */
    long   ru_isrss;         /* integral unshared stack size */
    long   ru_minflt;        /* page reclaims (soft page faults) */
    long   ru_majflt;        /* page faults (hard page faults) */
    long   ru_nswap;         /* swaps */
    long   ru_inblock;       /* block input operations */
    long   ru_oublock;       /* block output operations */
    long   ru_msgsnd;        /* IPC messages sent */
    long   ru_msgrcv;        /* IPC messages received */
    long   ru_nsignals;      /* signals received */
    long   ru_nvcsw;         /* voluntary context switches */
    long   ru_nivcsw;        /* involuntary context switches */
};

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

# Score-P Hands-On
## CUDA: Jacobi example

- # Jacobi Example
  - ## Iterative solver for system of equations
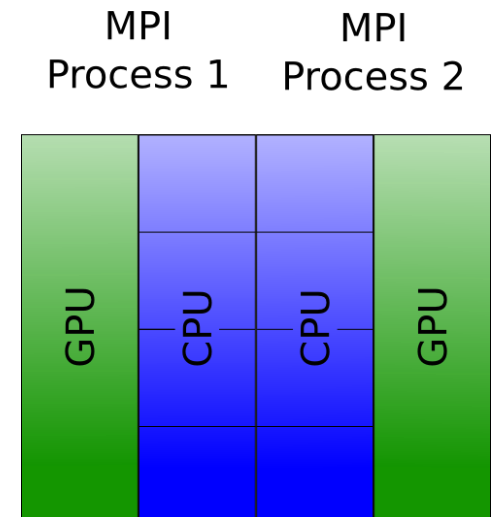
    $$U_{old} = U$$

    $$u_{i,j} = bu_{old,i,j} + a_x(u_{old,i-1,j} + u_{old,i+1,j}) + a_y(u_{old,i,j-1} + u_{old,i,j+1}) - rHs/b$$

  - ## Code uses OpenMP, CUDA and MPI for parallelization

- # Domain decomposition
  - ## Halo exchange at boundaries:
    - ### Via MPI between processes
    - ### Via CUDA between hosts and accelerators

# Jacobi Without Instrumentation

```
# Compile host code
%       mpicc -O3 -fopenmp -DUSE_MPI -I<path_to_cuda_header>
         -c jacobi_cuda.c -o jacobi_mpi+cuda.o

# Compile CUDA kernel
%       nvcc -O3 -c jacobi_cuda_kernel.cu
         -o jacobi_cuda_kernel.o

# Link executable
%       mpicc -fopenmp -lm -L<path_tocuda_libs> -lcudart
         jacobi_mpi+cuda.o jacobi_cuda_kernel.o -o ./jacobi_mpi+cuda
```

# Instrumentation with Score-P

```
# Compile host code
% scorep mpicc -O3 -fopenmp -DUSE_MPI –I<path_to_cuda_header>
        -c jacobi_cuda.c -o jacobi_mpi+cuda.o

# Compile CUDA kernel
% scorep nvcc -O3 -c jacobi_cuda_kernel.cu
        -o jacobi_cuda_kernel.o

# Link executable
% scorep mpicc -fopenmp -lm –L<path_tocuda_libs> -lcudart
        jacobi_mpi+cuda.o jacobi_cuda_kernel.o -o ./jacobi_mpi+cuda
```

- Enable recording of CUDA events with the CUPTI interface via environment variable **`SCOREP_CUDA_ENABLE`**

- Provide a list of recording types, e.g.

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

- Start with using the default configuration

```
% export SCOREP_CUDA_ENABLE=yes
```

- Adjust CUPTI buffer size (in bytes) as needed

```
% export SCOREP_CUDA_BUFFER=100000
```

# SCOREP_CUDA_ENABLE: Recording Types

| Recording type | Remark |
| --- | --- |
| **yes/DEFAULT/1** | "runtime, kernel, memcpy" |
| **no** | Disable CUDA measurement (same as unset SCOREP_CUDA_ENABLE) |
| **runtime** | CUDA runtime API |
| **driver** | CUDA driver API |
| **kernel** | CUDA kernels |
| **kernel_counter** | Fixed CUDA kernel metrics |
| **idle** | GPU compute idle time |
| **pure_idle** | GPU idle time (memory copies are not idle) |
| **memcpy** | CUDA memory copies |
| **sync** | Record implicit and explicit CUDA synchronization |
| **gpumemusage** | Record CUDA memory (de)allocations as a counter |

# Measurement (Profiling)
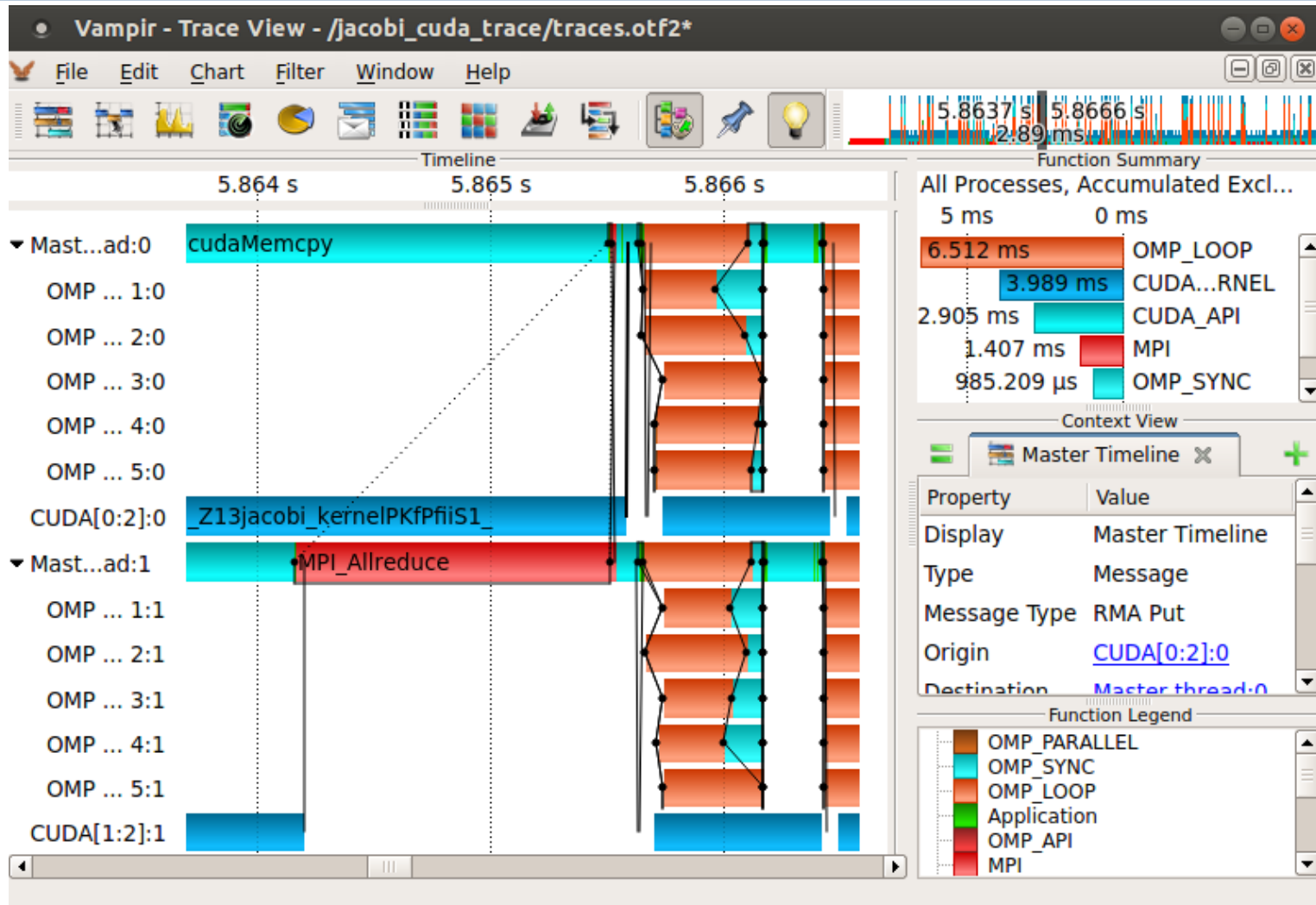
```
% export OMP_NUM_THREADS=6
% export SCOREP_CUDA_ENABLE=yes
% export SCOREP_CUDA_BUFFER=500000
% export SCOREP_EXPERIMENT_DIRECTORY=jacobi_cuda_profile

% mpirun -n 2 ./jacobi_mpi+cuda 4096 4096 0.15


Jacobi relaxation Calculation: 4096 x 4096 mesh with
 2 processes and 6 threads + one Tesla T10 Processor for each process.
 307 of 2049 local rows are calculated on the CPU to balance the load
 between the CPU and the GPU.
     0, 0.113429
  … … … … … …
  900, 0.000101
 total: 12.83581
```
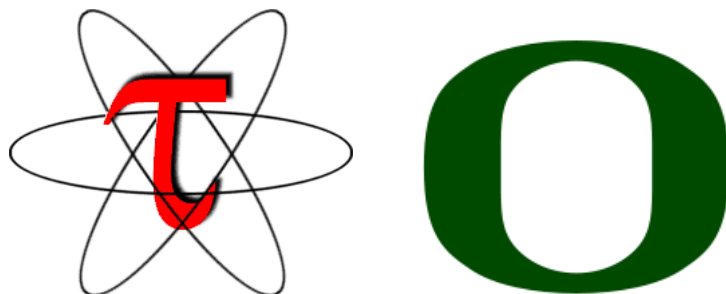
Problem size
(x dimension)

Problem size
(y dimension)

Load balancing factor
(in this example 15% of the
computations are calculated
on the CPU)

# CUBE4 Analysis

```
% cube jacobi_cuda_profile/profile.cubex
```

- Do we need to filter? (Overhead and memory footprint)

```
% scorep-score jacobi_cuda_profile/profile.cubex
Estimated aggregate size of event trace (total_tbc):      3.875.472 bytes
Estimated requirements for largest trace buffer (max_tbc): 1.937.936 bytes
(hint: When tracing set SCOREP_TOTAL_MEMORY > max_tbc to avoid
       intermediate flushes or reduce requirements using file listing
       names of USR regions to be filtered.)

flt type          max_tbc          time      % region
    ALL           1937936         24.97  100.0 ALL
    OMP           1154110         18.78   75.2 OMP
    USR            667480          5.95   23.8 USR
    MPI            116192          0.14    0.5 MPI
    COM               154          0.10    0.4 COM
```

☞ Very small example => no filtering

# Measurement (Tracing)

```
% export OMP_NUM_THREADS=6
% export SCOREP_CUDA_ENABLE=yes
% export SCOREP_CUDA_BUFFER=500000
% export SCOREP_EXPERIMENT_DIRECTORY=jacobi_cuda_trace
% export SCOREP_ENABLE_PROFILING=false
% export SCOREP_ENABLE_TRACING=true

% mpirun -n 2 ./jacobi_mpi+cuda 4096 4096 0.15

Jacobi relaxation Calculation: 4096 x 4096 mesh with
 2 processes and 6 threads + one Tesla T10 Processor for each process.
 307 of 2049 local rows are calculated on the CPU to balance the load
 between the CPU and the GPU.
    0, 0.113429
  … … … … … …
  900, 0.000101
 total: 12.875220 s
```

# Vampir Analysis

```
% vampir jacobi_cuda_trace/traces.otf2
```

# Performance Data Management with TAU PerfExplorer

Sameer Shende

Performance Research Lab, University of Oregon

http://TAU.uoregon.edu

- **Configure TAUdb (Done by each user)**

  % taudb_configure --create-default

  - Choose derby, PostgreSQL, MySQL, Oracle or DB2
  - Hostname
  - Username
  - Password
  - Say yes to downloading required drivers (we are not allowed to distribute these)
  - Stores parameters in your ~/.ParaProf/taudb.cfg file

- **Configure PerfExplorer (Done by each user)**

  % perfexplorer_configure

- **Execute PerfExplorer**

  % perfexplorer

```
% wget http://tau.uoregon.edu/data.tgz   (Contains CUBE profiles from Score-P)
% taudb_configure --create-default
(Chooses derby, blank user/passwd, yes to save passwd, defaults)
% perfexplorer_configure
(Yes to load schema, defaults)
% paraprof
(load each trial: DB -> Add Trial -> Type (Paraprof Packed Profile) -> OK) OR use
    taudb_loadtrial –a "app" –x "experiment" –n "name" file.ppk
Then,
% tar zxf $TAU/data.tgz; cd data/tau;
% taudb_loadtrial –a BT_MZ –x "Class_B" bt-mz_B.*.ppk
% perfexplorer
(Select experiment, Menu: Charts -> Speedup)
```

- Performance knowledge discovery framework
  - Data mining analysis applied to parallel performance data
    - comparative, clustering, correlation, dimension reduction, …
  - Use the existing TAU infrastructure
    - TAU performance profiles, taudb
  - Client-server based system architecture
- Technology integration
  - Java API and toolkit for portability
  - taudb
  - R-project/Omegahat, Octave/Matlab statistical analysis
  - WEKA data mining package
  - JFreeChart for visualization, vector output (EPS, SVG)

- Performance data represented as vectors - each dimension is the cumulative time for an event

- $k$-means: $k$ random centers are selected and instances are grouped with the "closest" (Euclidean) center

- New centers are calculated and the process repeated until stabilization or max iterations

- Dimension reduction necessary for meaningful results

- Virtual topology, summaries constructed

# PerfExplorer - Cluster Analysis (sPPM)

- Describes strength and direction of a linear relationship between two variables (events) in the data

- -0.995 indicates strong, negative relationship
- As CALC_CUT_ BLOCK_CONTRIBUTIO NS() increases in execution time, MPI_Barrier() decreases

- Relative speedup, efficiency
  - total runtime, by event, one event, by phase
- Breakdown of total runtime
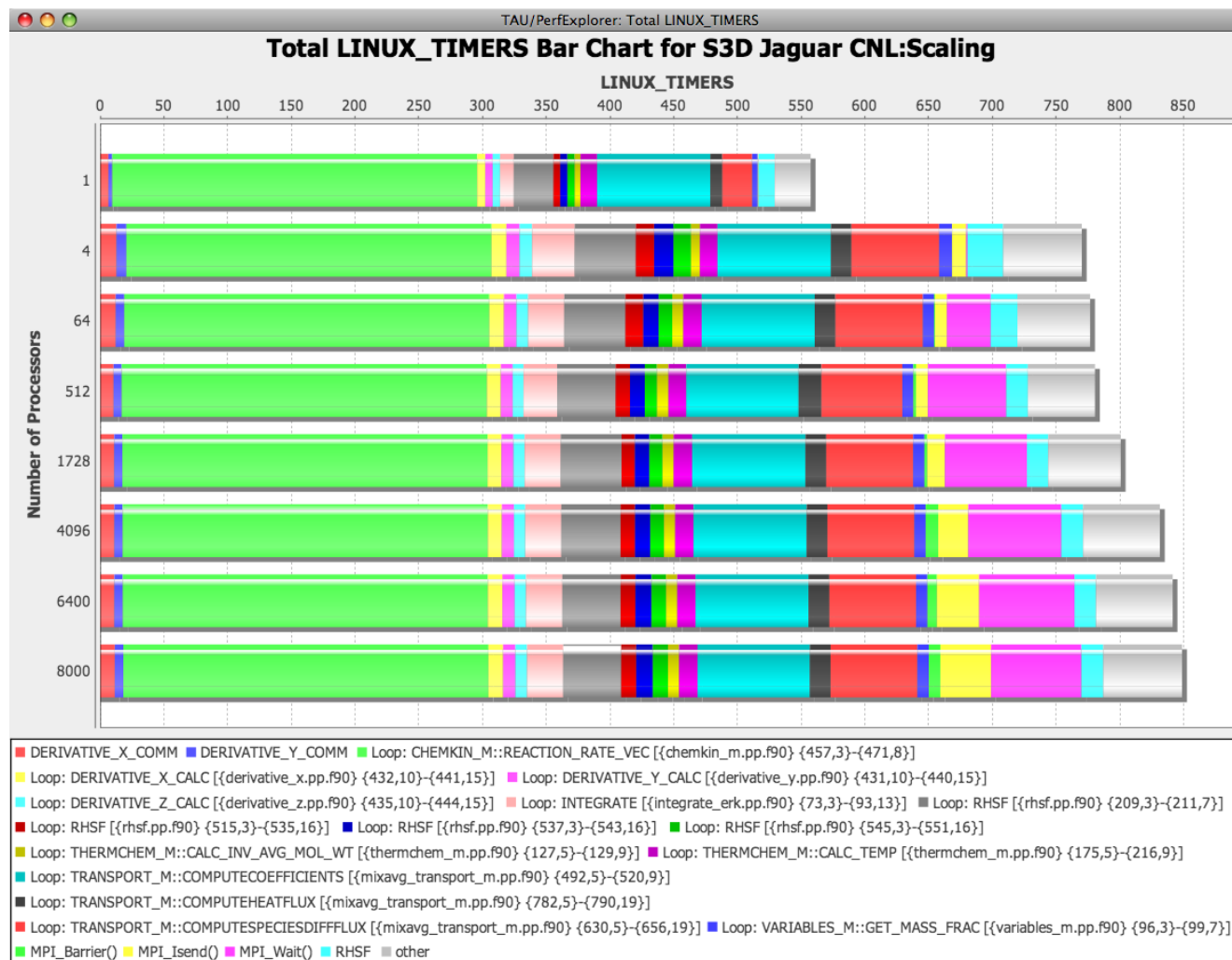- Group fraction of total runtime
- Correlating events to total runtime
- Timesteps per second

# PerfExplorer - Interface

- Goal: How does my application scale? What bottlenecks occur at what core counts?
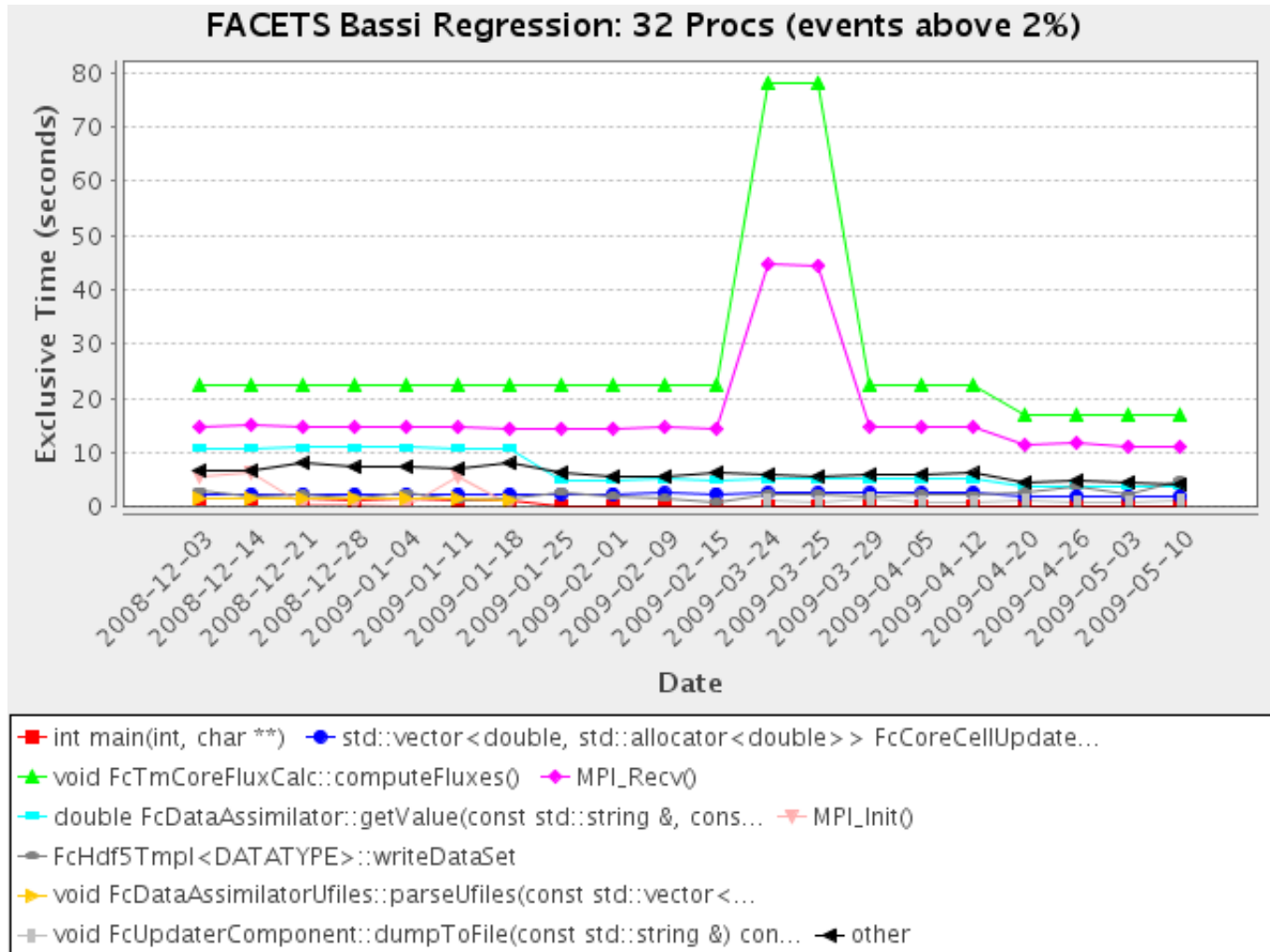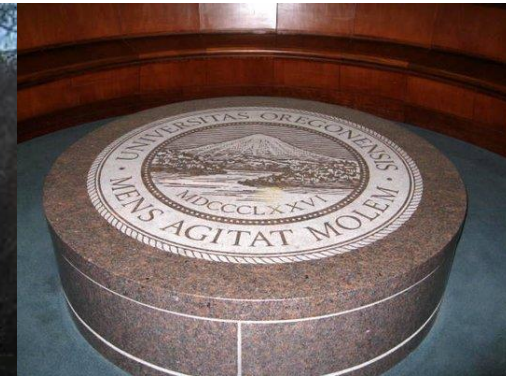- Load profiles in taudb database and examine with PerfExplorer

Total LINUX_TIMERS Bar Chart for S3D Jaguar CNL:Scaling

TAU/PerfExplorer: Total TIME

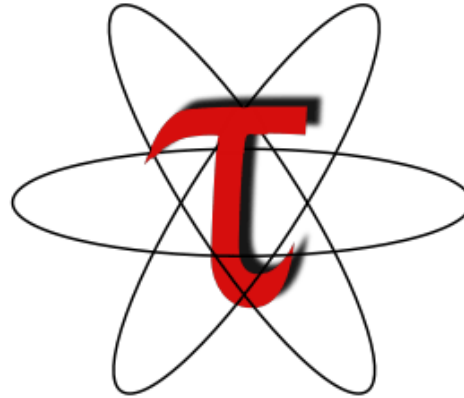**Total TIME Bar Chart for IRMHD:Scaling_BGP**

Total TIME Breakdown for IRMHD:Scaling_BGP

www.uoregon.edu

- U.S. Department of Energy (DOE)
    - Office of Science
    - PNNL, LBL, ORNL
    - ASC/NNSA, Tri-labs (LLNL,LANL, SNL)
- U.S. Department of Defense (DoD)
    - HPC Modernization Office (HPCMO)
- NSF Software Development for Cyberinfrastructure (SDCI)
- Juelich Supercomputing Center, NIC
- Argonne National Laboratory
- T.U. Dresden
- ParaTools, Inc.

**http://tau.uoregon.edu**

**http://www.hpclinux.com [LiveDVD, OVA]**

**Free download, open source, BSD license**

# Typical performance bottlenecks and how they can be identified

Bert Wesarg

ZIH, Technische Universität Dresden

- **Case I:**
  - **Load imbalances in OpenMP codes**


- **Case II:**
  - **Communication and computation overlapping in MPI codes**

- **Note:** We won't do the complete performance engineering cycle here.

$$\begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} = \begin{pmatrix} a_{11} & \cdots & a_{n1} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

- A sparse matrix is a matrix populated primarily with zeros

- Only non-zero elements of $a_{ij}$ are saved efficiently in memory

- Algorithm

```
foreach row r in A
  y[r.x] = 0
  foreach non-zero element e in row
    y[r.x] += e.value * x[e.y]
```

- Naive OpenMP Algorithm

```
#pragma omp parallel for
foreach row r in A
  y[r.x] = 0
  foreach non-zero element e in row
    y[r.x] += e.value * x[e.y]
```

- Distributes the rows of A evenly across the threads in the parallel region

- The distribution of the non-zero elements may influence the load balance in the parallel application

# Case I: Load imbalances in OpenMP codes
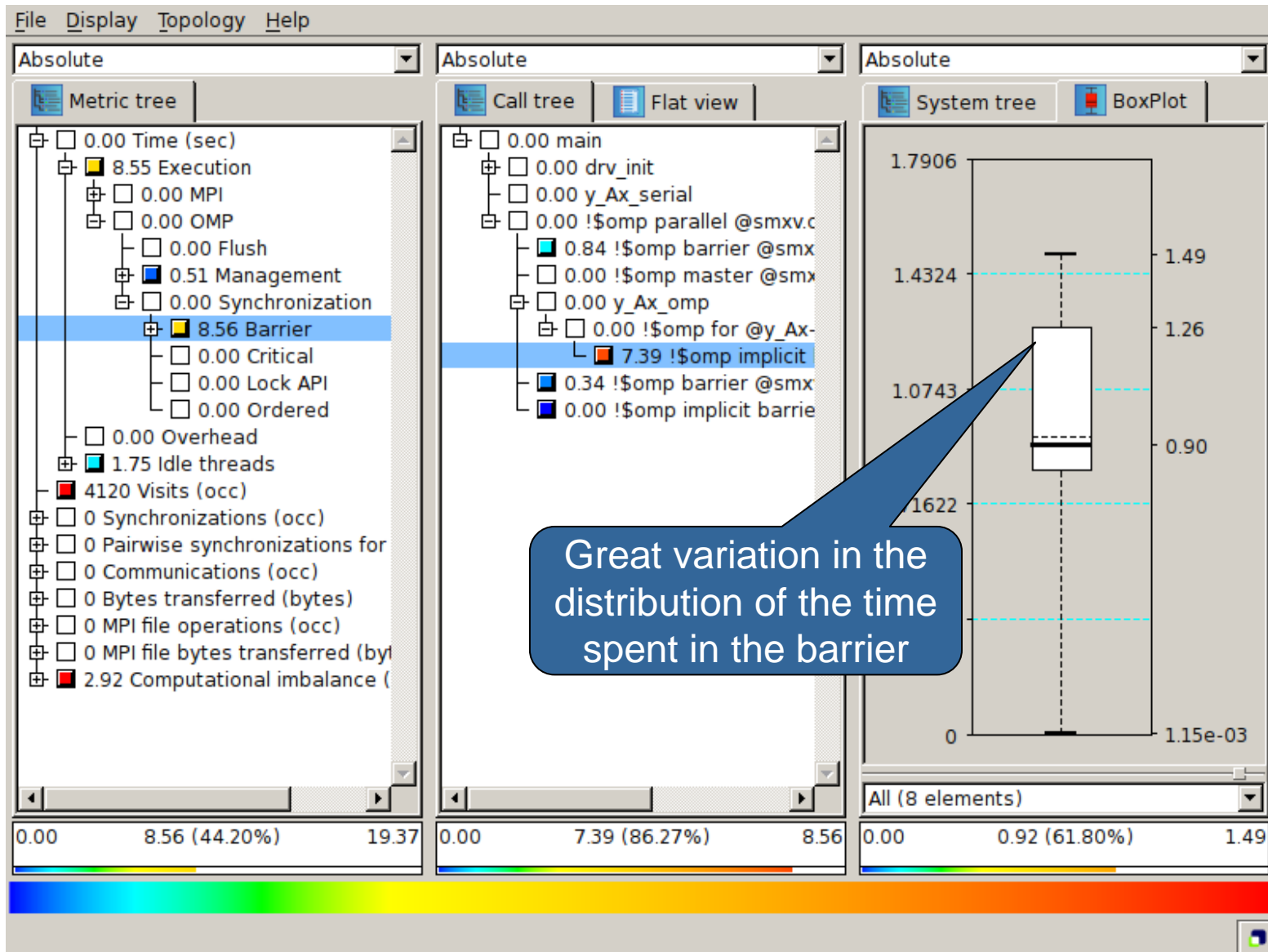
- ## Measuring the static OpenMP application

```
% cd ~/Bottlenecks/smxv
% make PREP=scorep
scorep   gcc -fopenmp -DLITTLE_ENDIAN \
    -DFUNCTION_INC='"y_Ax-omp.inc.c"' -DFUNCTION=y_Ax_omp \
    -o smxv-omp smxv.c -lm
scorep   gcc -fopenmp -DLITTLE_ENDIAN \
    -DFUNCTION_INC='"y_Ax-omp-dynamic.inc.c"' \
    -DFUNCTION=y_Ax_omp_dynamic -o smxv-omp-dynamic smxv.c -lm
% OMP_NUM_THREADS=8 scan -t ./smxv-omp yax_large.bin
```

- Two metrics which indicate load imbalances:
  - Time spent in OpenMP barriers
  - Computational imbalance

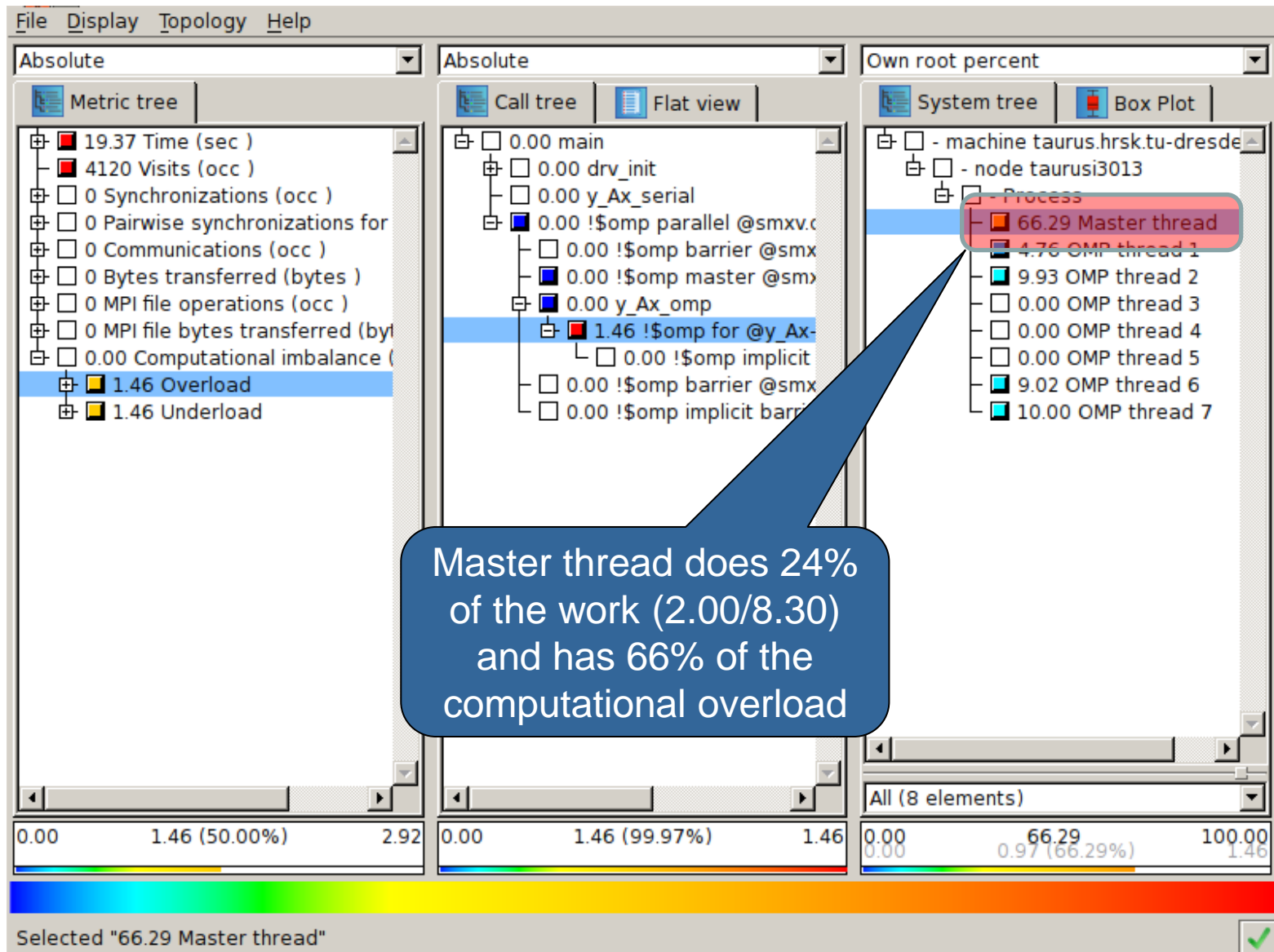- Open prepared measurement on the LiveDVD with Cube

```
% cube ~/Bottlenecks/smxv/scorep_smxv-omp_large/trace.cubex

                [CUBE GUI showing trace analysis report]
```

Master thread does 24% of the work (2.00/8.30) and has 66% of the computational overload

- # Improved OpenMP Algorithm

```
#pragma omp parallel for schedule(dynamic,1000)
foreach row r in A
  y[r.x] = 0
  foreach non-zero element e in row
    y[r.x] += e.value * x[e.y]
```
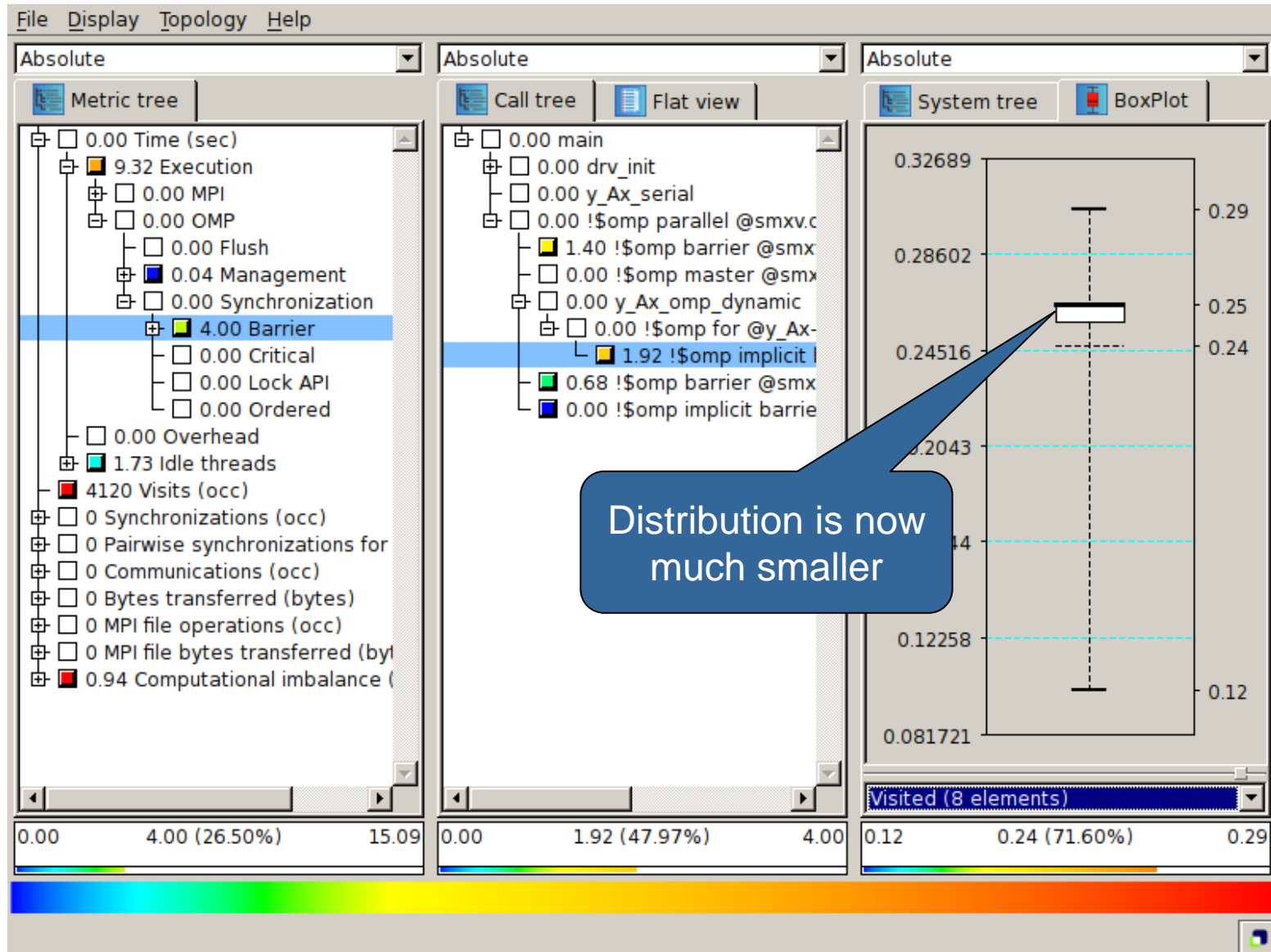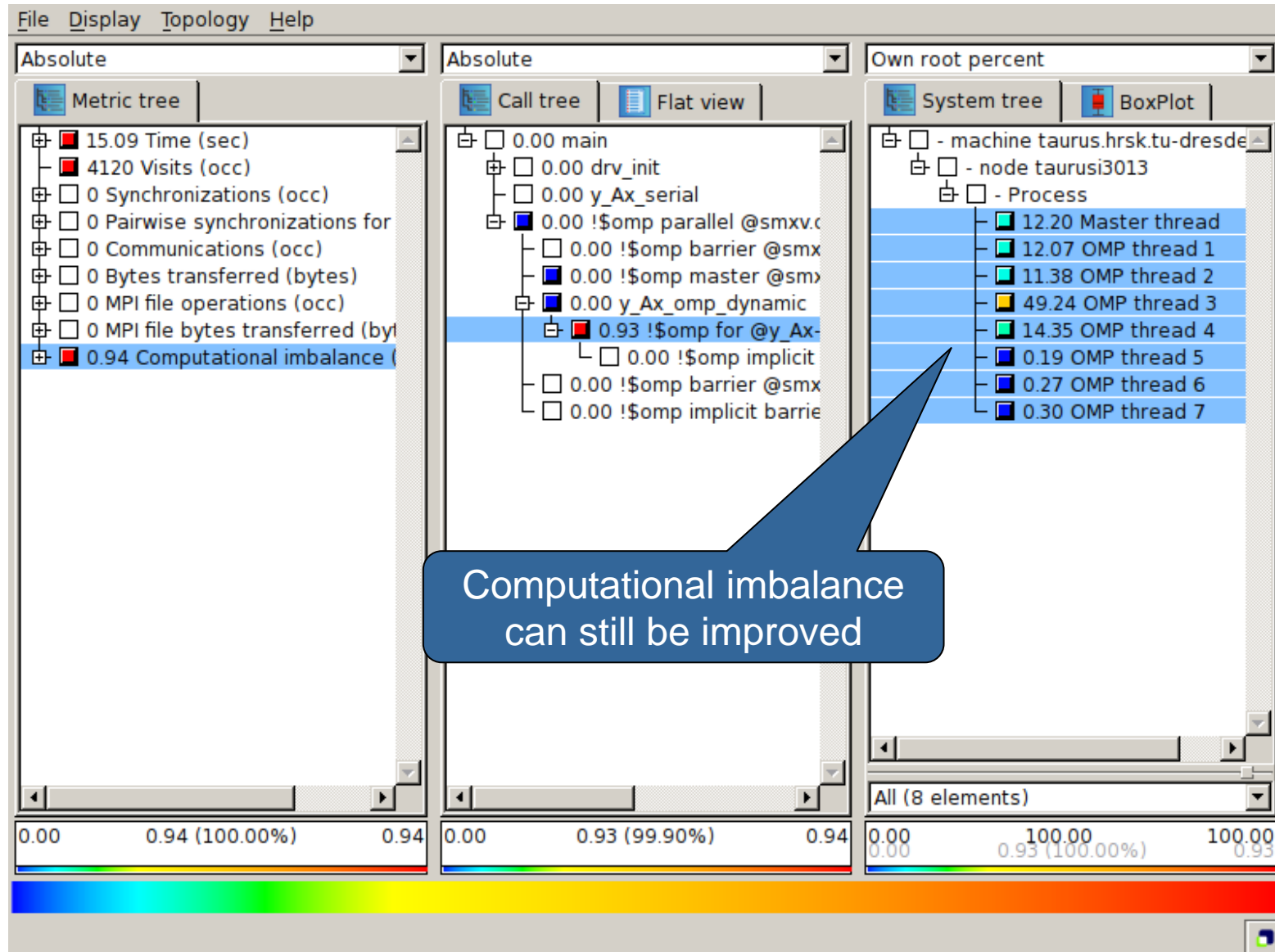
- # Distributes the rows of A *dynamically* across the threads in the parallel region

- # Two metrics which indicate load imbalances
  - ## Time spent in OpenMP barriers
  - ## Computational imbalance

- # Open prepared measurement on the LiveDVD with Cube

```
% cube ~/Bottlenecks/smxv/scorep_smxv-omp-dynamic_large/trace.cubex

                [CUBE GUI showing trace analysis report]
```

# Case I: Time spent in OpenMP barriers

# Case I: Computational imbalance



Computational imbalance
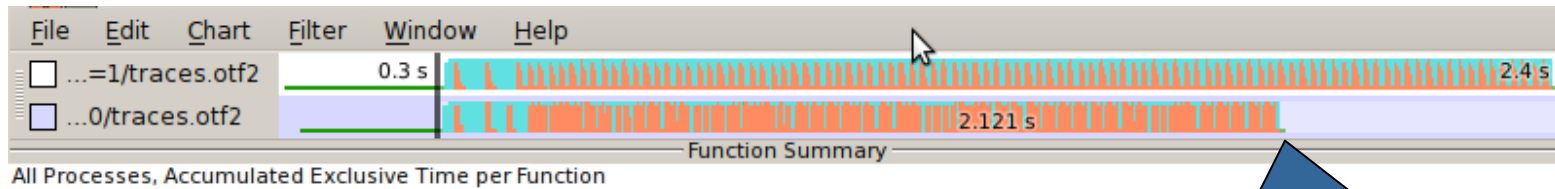can still be improved

- Open prepared measurement on the LiveDVD with Vampir

```
% vampir ~/Bottlenecks/smxv/scorep_smxv-omp_large/traces.otf2 \
        ~/Bottlenecks/smxv/scorep_smxv-omp-dynamic_large/traces.otf2

                        [Vampir GUI showing trace]
```
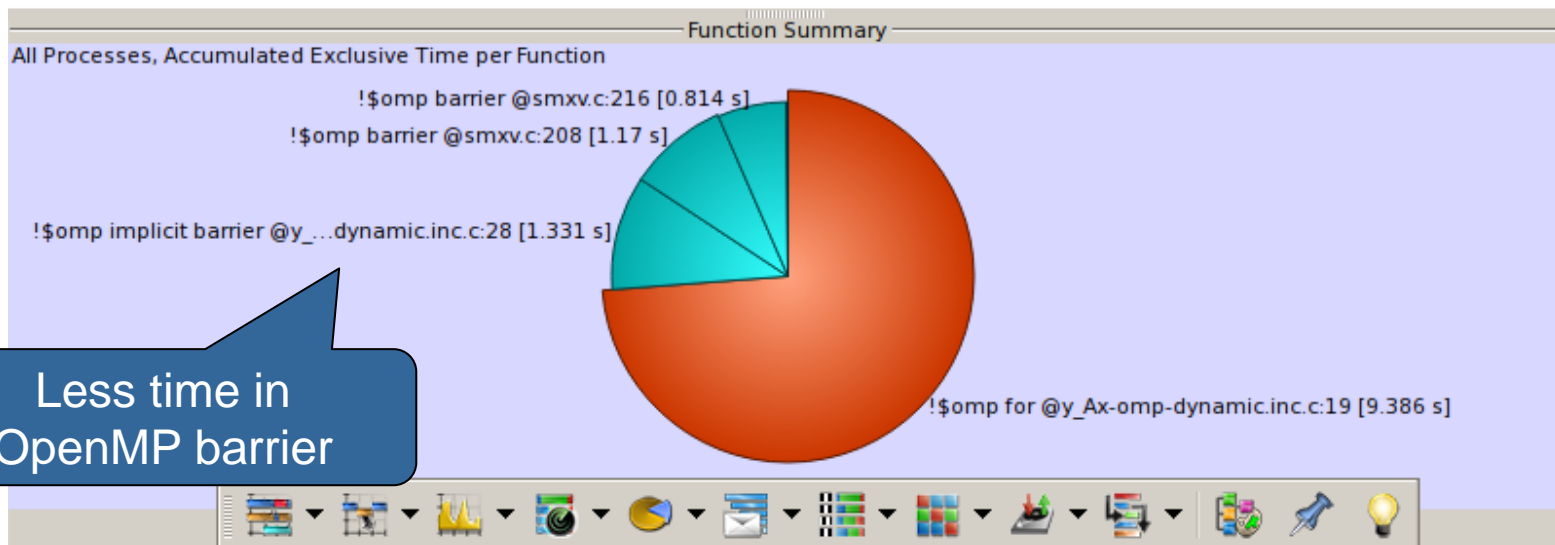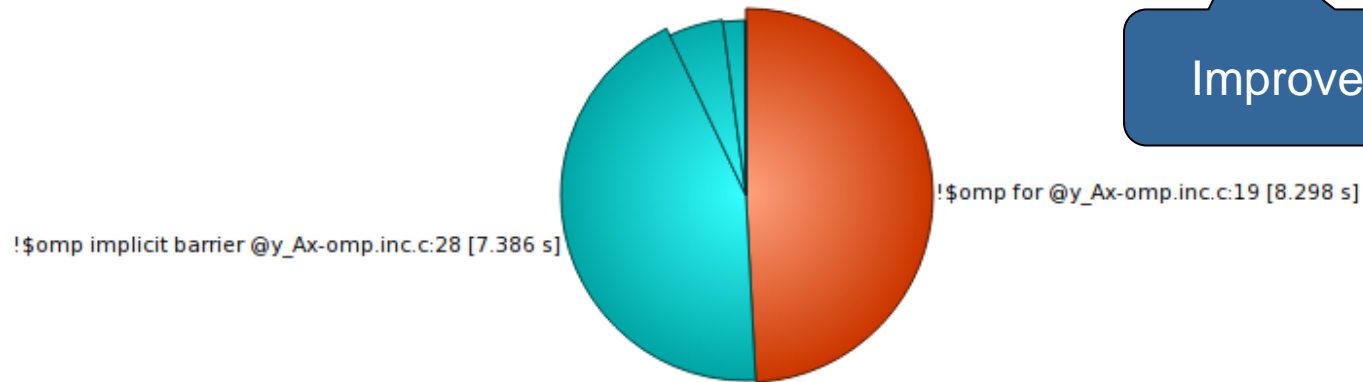
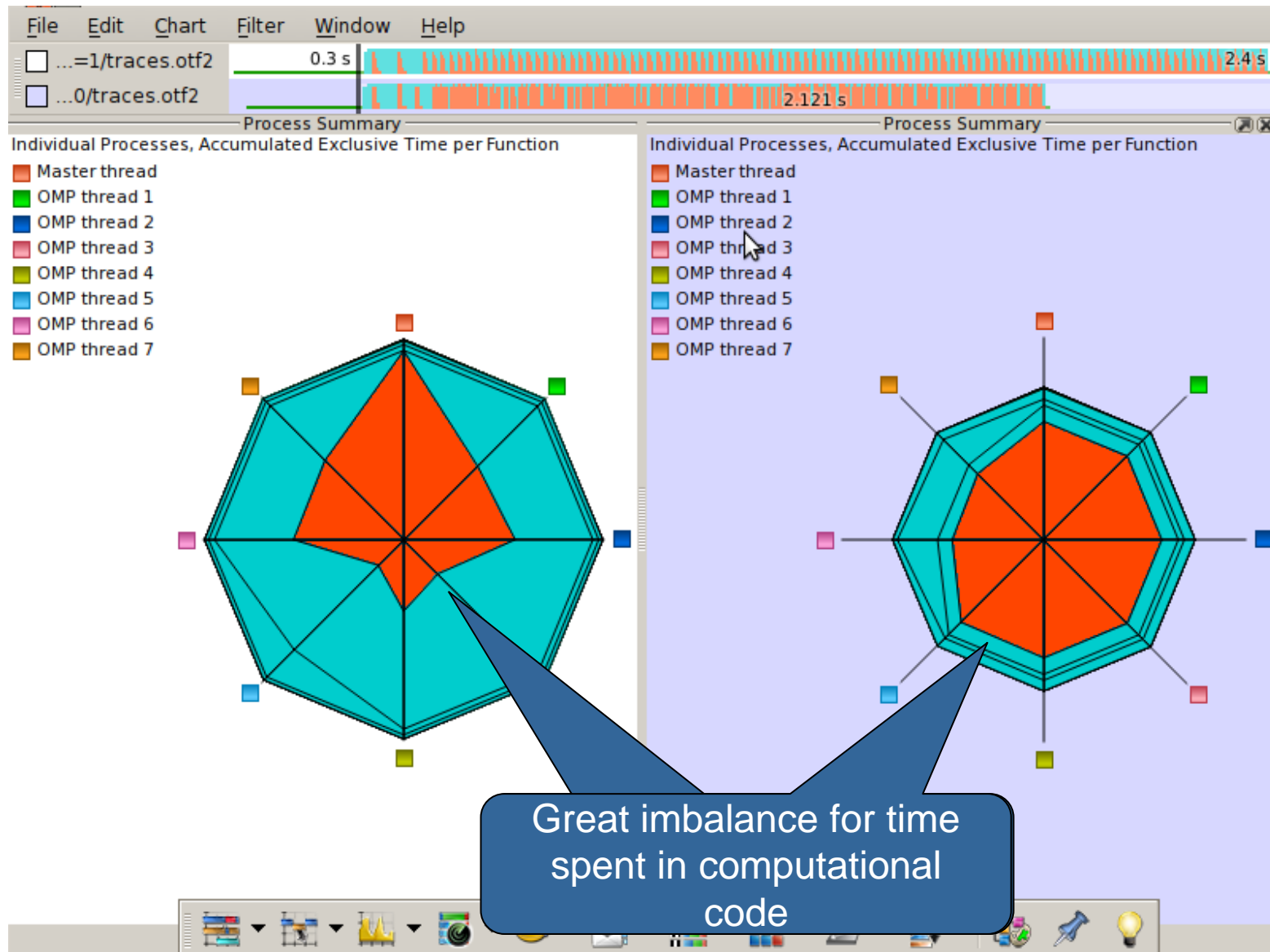# Case I: Time spent in OpenMP barriers

# Case I: Computational imbalance



Great imbalance for time spent in computational code

- **Case I:**
  - Load imbalances in OpenMP codes

- **Case II:**
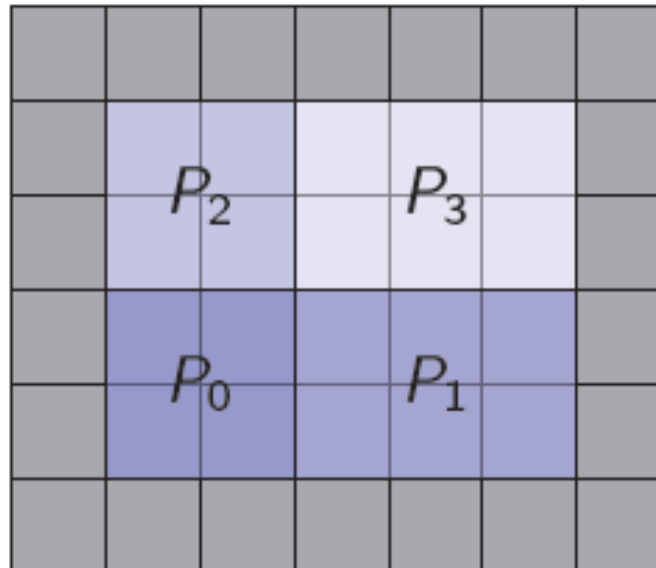  - **Communication and computation overlapping in MPI codes**

- Calculating the heat conduction at each time step
- Discretized formula for space $dx, dy$ and time $dt$

$$\theta_{i,j}^{t+1} = \theta_{i,j}^t + \left( \frac{\theta_{i+1,j}^t - 2\theta_{i,j}^t + 2\theta_{i-1,j}^t}{dx^2} + \frac{\theta_{i,j+1}^t - 2\theta_{i,j}^t + 2\theta_{i,j-1}^t}{dy^2} \right) \cdot k \cdot dt$$
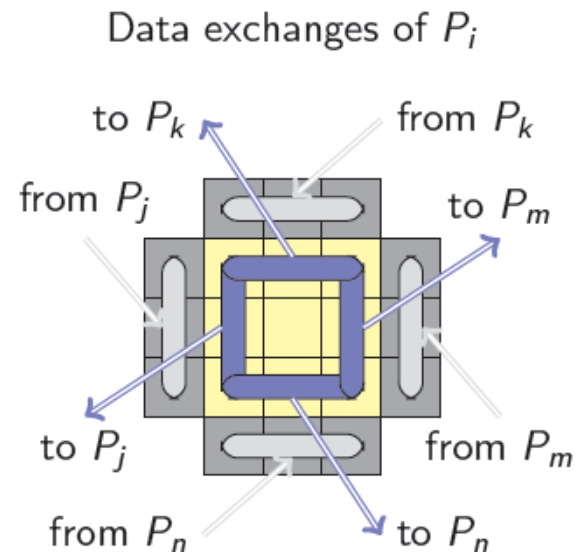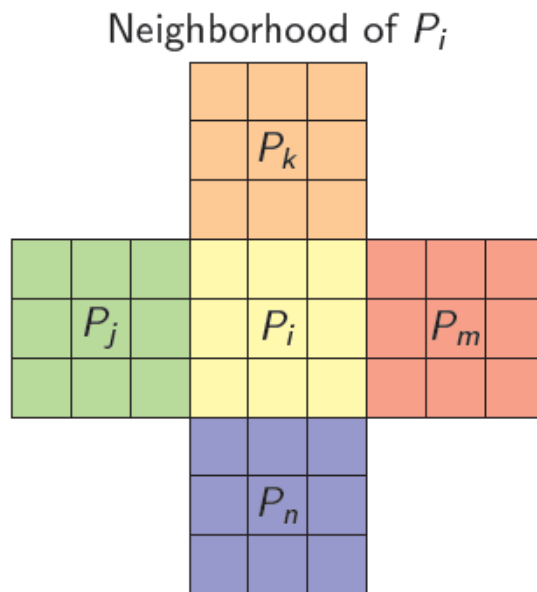
- Application uses MPI for boundary exchange
- Simulation grid is distributed across MPI ranks

- Ranks need to exchange boundaries before next iteration step



Neighborhood of $P_i$

Data exchanges of $P_i$

- ## MPI algorithm

```
foreach step in [1:nsteps]
  exchangeBoundaries
  computeHeatConduction
```

- ## Building and measuring the heat conduction application

```
% cd ~/Bottlenecks/heat
% make PREP='scorep --user'
  [... make output ...]
% scan mpirun –np 16 ./heat-MPI 3072 32
```
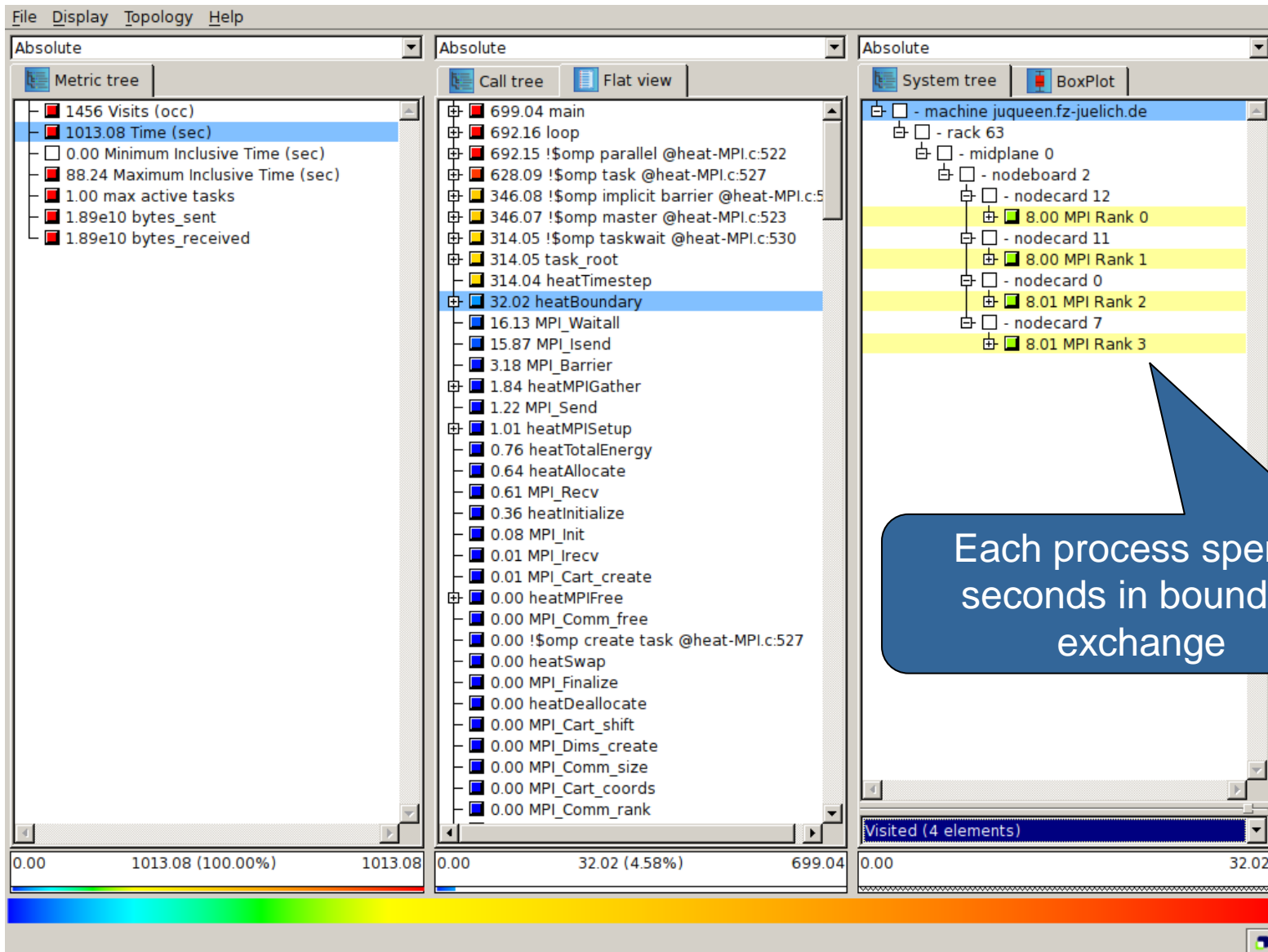
- ## Open prepared measurement on the LiveDVD with Cube

```
% cube ~/Bottlenecks/heat/scorep_heat-MPI_small/profile.cubex

              [CUBE GUI showing trace analysis report]
```

# Case II: Time spent in Boundary Exchange



Each process spent 8 seconds in boundary exchange

# Case II: Time spent in Boundary Exchange



… that's ~10% of the computation

- Step 1: Compute heat in the area which is communicated to your neighbors



Compute heat conduction in the boundaries of $P_i$

- Step 2: Start communicating boundaries with your neighbors

- Step 3: Compute heat in the interior area



Compute heat conduction
in the interior of $P_i$

- ## Improved MPI algorithm

```
foreach step in [1:nsteps]
  computeHeatConductionInBoundaries
  startBoundaryExchange
  computeHeatConductionInInterior
  waitForCompletionOfBoundaryExchange
```

- Note: As not all MPI implementations support overlapping, it is here done with the help of OpenMP tasks.

- Measuring the improved heat conduction application

```
% scan mpirun –np 16 ./heat-MPI-overlap 3072 32
```

- ## Open prepared measurement on the LiveDVD with Cube

```
% cube ~/Bottlenecks/heat/scorep_heat-MPI-overlap_small/profile.cubex

              [CUBE GUI showing trace analysis report]
```
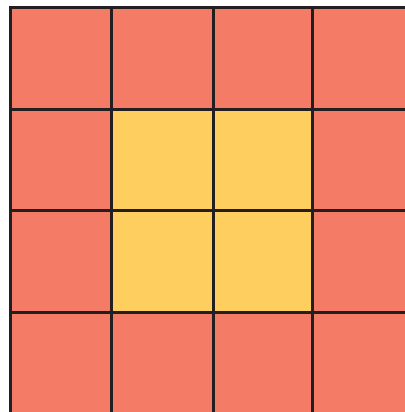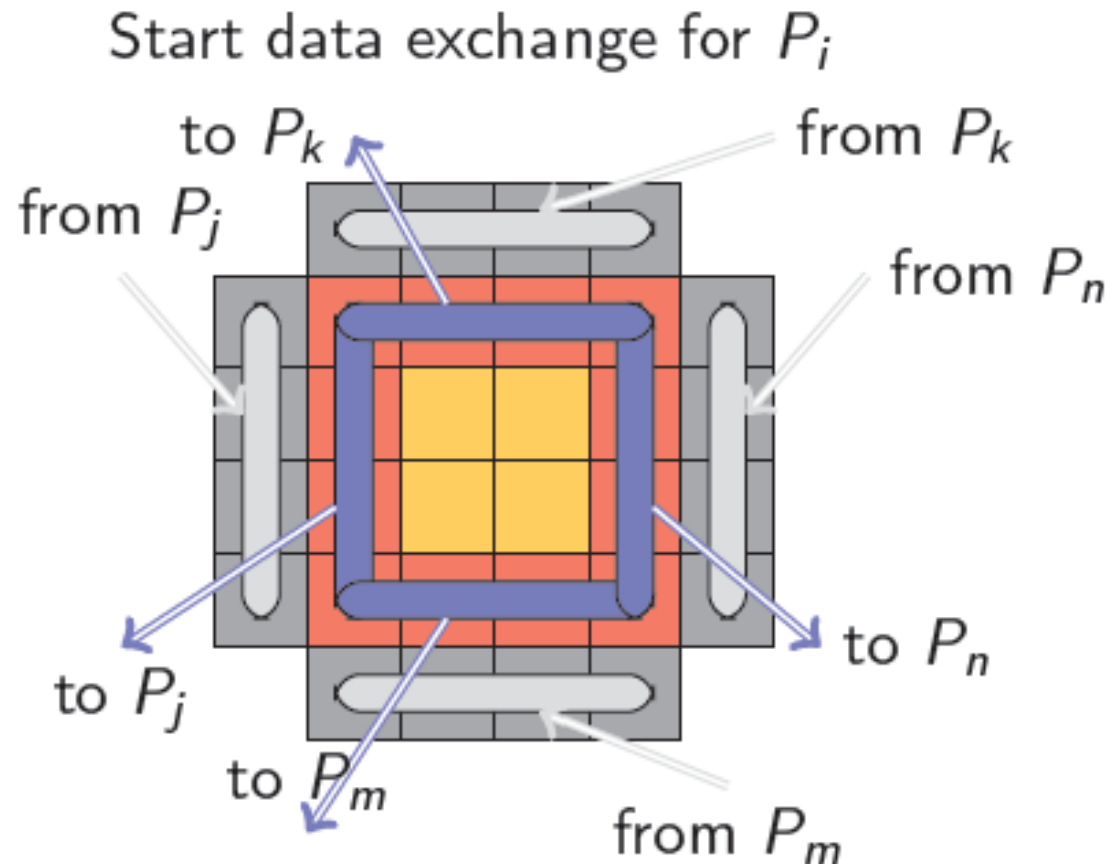
Still ~8 seconds in boundary exchange

- ## Calculate differences between profiles

```
% cube_diff ~/Bottlenecks/heat/scorep_heat-MPI_small/profile.cubex \
            ~/Bottlenecks/heat/scorep_heat-MPI-overlap_small/profile.cubex
```

- ## Open prepared profile diff on the LiveDVD with Cube

```
% cube ~/Bottlenecks/heat/diff.cubex

                [CUBE GUI showing trace analysis report]
```

But all threads spend ~8 seconds less in the main loop

- Open prepared measurement on the LiveDVD with Vampir

```
% vampir ~/Bottlenecks/heat/scorep_heat-MPI_small/traces.otf2 \
        ~/Bottlenecks/heat/scorep_heat-MPI-overlap_small/traces.otf2

                    [Vampir GUI showing trace]
```

# Case II: Trace Comparison



Improved runtime

- Thanks to Dirk Schmidl, RWTH Aachen, for providing the sparse matrix vector multiplication code

# Review

## Brian Wylie
## Jülich Supercomputing Centre

You've been introduced to a variety of tools

- – with hints to apply and use the tools effectively

- Tools provide complementary capabilities
  - – computational kernel & processor analyses
  - – communication/synchronization analyses
  - – load-balance, scheduling, scaling, …
- Tools are designed with various trade-offs
  - – general-purpose versus specialized
  - – platform-specific versus agnostic
  - – simple/basic versus complex/powerful

- Which tools you use and when you use them likely to depend on situation
  - which are available on (or for) your computer system
  - which support your programming paradigms and languages
  - which you are familiar (comfortable) with using
  - which type of issue you suspect
  - which question you want to have answered

- Being aware of (potentially) available tools and their capabilities can help finding the most appropriate tools

- ## First ensure that the parallel application runs correctly
  - no-one will care how quickly you can get invalid answers or produce a directory full of corefiles
  - parallel debuggers help isolate known problems
  - correctness checking tools can help identify other issues
  - (that might not cause problems right now, but will eventually)
    - e.g., race conditions, invalid/non-compliant usage

- ## Generally valuable to start with an overview of execution performance
  - fraction of time spent in computation vs comm/synch vs I/O
  - which sections of the application/library code are most costly

- ## and how it changes with scale or different configurations
  - processes vs threads, mappings, bindings

- Communication/synchronization issues generally apply to every computer system (to different extents) and typically grow with the number of processes/threads
  - *Weak scaling*: fixed computation per thread, and perhaps fixed localities, but increasingly distributed
  - *Strong scaling*: constant total computation, increasingly divided amongst threads, while communication grows
  - Collective communication (particularly of type "all-to-all") result in increasing data movement
  - Synchronizations of larger groups are increasingly costly
  - Load-balancing becomes increasingly challenging, and imbalances increasingly expensive
    - generally manifests as waiting time at following collective ops

- ## Waiting times are difficult to determine in basic profiles
  - Part of the time each process/thread spends in communication & synchronization operations may be wasted waiting time
  - Need to correlate event times between processes/threads
    - *Periscope* uses augmented messages to transfer timestamps and additional on-line analysis processes
    - Post-mortem event trace analysis avoids interference and provides a complete history
    - *Scalasca* automates trace analysis and ensures waiting times are completely quantified
    - *Vampir* allows interactive exploration and detailed examination of reasons for inefficiencies

# Effective computation within processors/cores is also vital

- Optimized libraries may already be available
- Optimizing compilers can also do a lot
  - provided the code is clearly written and not too complex
  - appropriate directives and other hints can also help
- Processor hardware counters can also provide insight
  - although hardware-specific interpretation required
- Tools available from processor and system vendors help navigate and interpret processor-specific performance issues

# Technologies and their integration

KCACHEGRIND

LWM2 / MAP / MPIP / O|SS / MAQAO

TAU

SCORE-P

PAPI

Hardware monitoring

Automatic profile & trace analysis

PERISCOPE

MUST

SCALASCA

DDT

Debugging, error & anomaly detection

Visual trace analysis

VAMPIR / PARAVER

STAT

Execution

Optimization

RUBIK / MAQAO

SYSMON / SPINDLE / SIONLIB / OPENMPI

- **Score-P**
  - community-developed instrumenter & measurement libraries for parallel profiling and event tracing
- **CUBE & ParaProf/PerfExplorer**
  - interactive parallel profile analyses
- **Scalasca**
  - automated event-trace analysis
- **Vampir**
  - interactive event-trace visualizations and analyses
- **TAU/PDT**
  - comprehensive performance system

- Website
  - Introductory information about the VI-HPS portfolio of tools for high-productivity parallel application development
    - VI-HPS Tools Guide
    - links to individual tools sites for details and download
  - Training material
    - tutorial slides
    - latest ISO image of VI-HPS Linux DVD with productivity tools
    - user guides and reference manuals for tools
  - News of upcoming events
    - tutorials and workshops
    - mailing-list sign-up for announcements

## http://www.vi-hps.org