# Trace Based Performance Analysis at Large Scale on Jaguar and Titan

Guido Juckeland (guido.juckeland@tu-dresden.de)

Performance tools will not automatically make you code run faster. They help you understand, what your code does  and where to put in work.

# Agenda

**Introduction**

**The Vampir Workflow**

**VampirTrace**

- Instrumentation & Run-Time Measurement

**The Vampir Displays**

- "Seeing" Performance Bottlenecks

**Vampir on Jaguar and Titan**

- Scalability & GPU Support
- How to use Vampir at OLCF

**Conclusions**

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
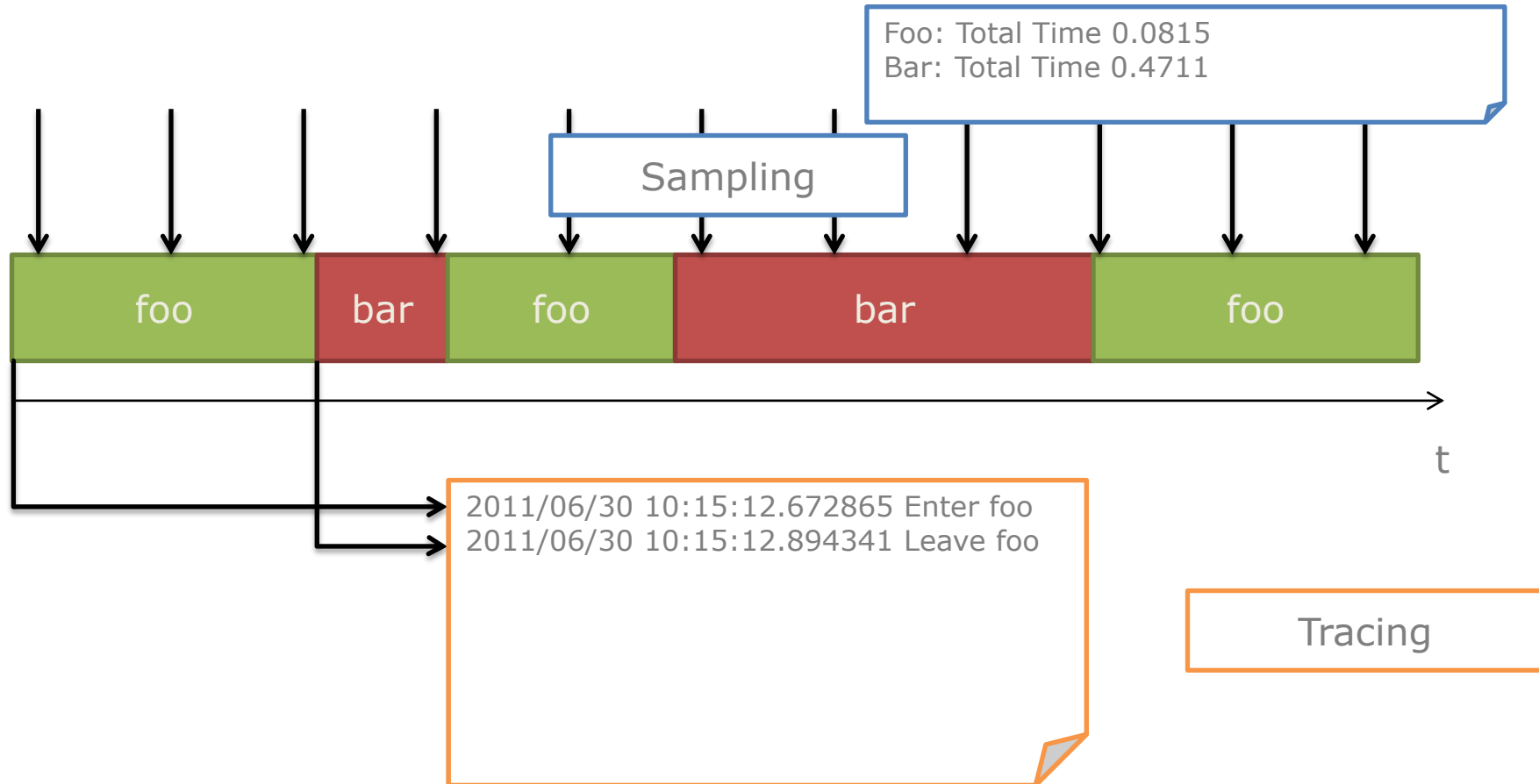High Performance Computing

# Introduction

## Why bother with performance analysis?

- Efficient usage of expensive and limited resources
- Scalability to achieve next bigger simulation
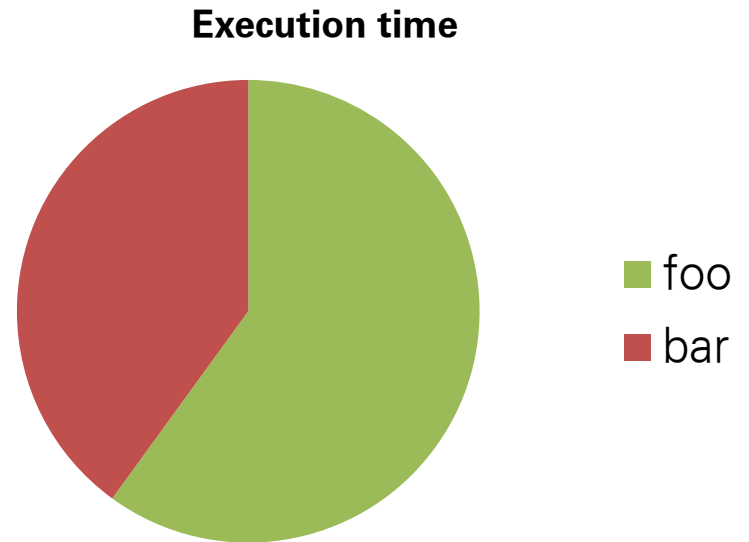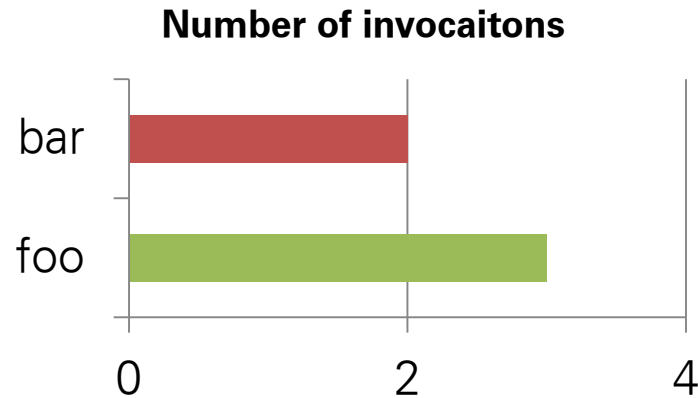
## Profiling and Tracing

- Have an optimization phase
  - just like testing and debugging phase
- Use tools!
- Avoid *do-it-yourself-with-printf* solutions, really!

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Sampling vs. Tracing

Foo: Total Time 0.0815
Bar: Total Time 0.4711

Sampling

| foo | bar | foo | bar | foo |

t

2011/06/30 10:15:12.672865 Enter foo
2011/06/30 10:15:12.894341 Leave foo

Tracing

**TECHNISCHE
UNIVERSITÄT
DRESDEN**

ZIH
Center for Information Services &
High Performance Computing

# Program profiles

Statistics

**Number of invocaitons**

**Execution time**



Timelines

| foo | bar | foo | bar | foo |
|-----|-----|-----|-----|-----|

# Instrumentation & Measurement

## What do you need to do for it?

- Use VampirTrace

Instrumentation (automatic with compiler wrappers)

| | |
|---|---|
| CC=icc | CC=vtcc |
| CXX=icpc | CXX=vtcxx |
| F90=ifc | F90=vtf90 |
| MPICC=mpicc | MPICC=vtcc -vt:cc mpicc |

Re-compile & re-link

Trace Run (run with appropriate test data set)

More details later

# Instrumentation & Measurement

**What does VampirTrace do in the background?**

## Instrumentation:

- Via compiler wrappers
- By underlying compiler with specific options
- MPI instrumentation with replacement lib
- OpenMP instrumentation with Opari
- Also binary instrumentation with Dyninst
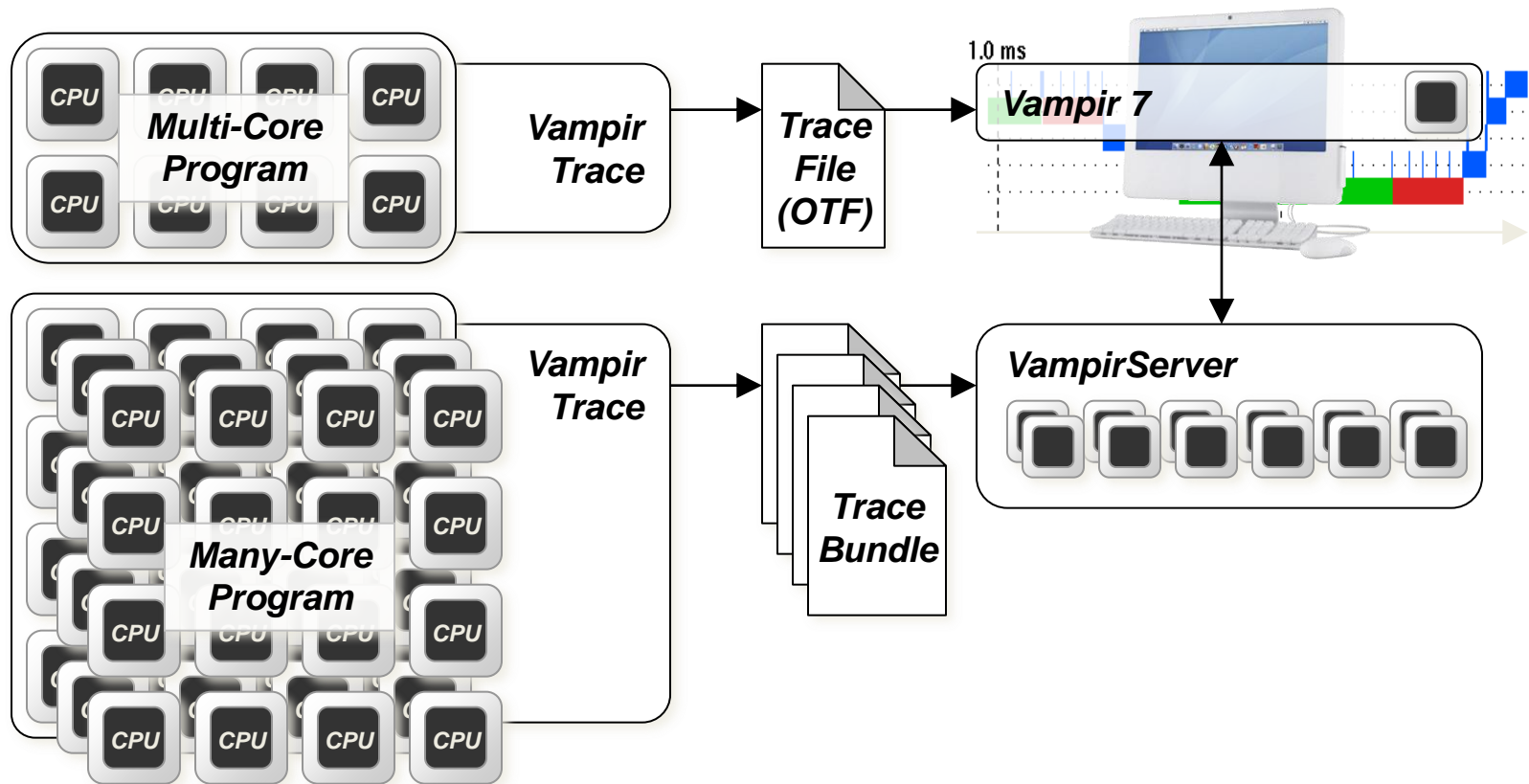- Partial manual instrumentation

# Instrumentation & Measurement

**What does VampirTrace do in the background?**

## Trace Run:

- Event data collection
- Precise time measurement
- Parallel timer synchronization
- Collecting parallel process/thread traces
- Collecting performance counters (from PAPI, memory usage, POSIX I/O calls and fork/system/exec calls, CUDA, and more … )
- Monitor GPU usage
- Filtering and grouping of function calls

# Event Trace Visualization

## Trace Visualization

- Alternative and supplement to automatic analysis
- Show dynamic run-time behavior graphically
- Provide statistics and performance metrics
  - Master timeline for parallel processes/threads
  - Process timeline plus performance counters
  - Statistics summary display
  - Message statistics
  - more
- Interactive browsing, zooming, selecting
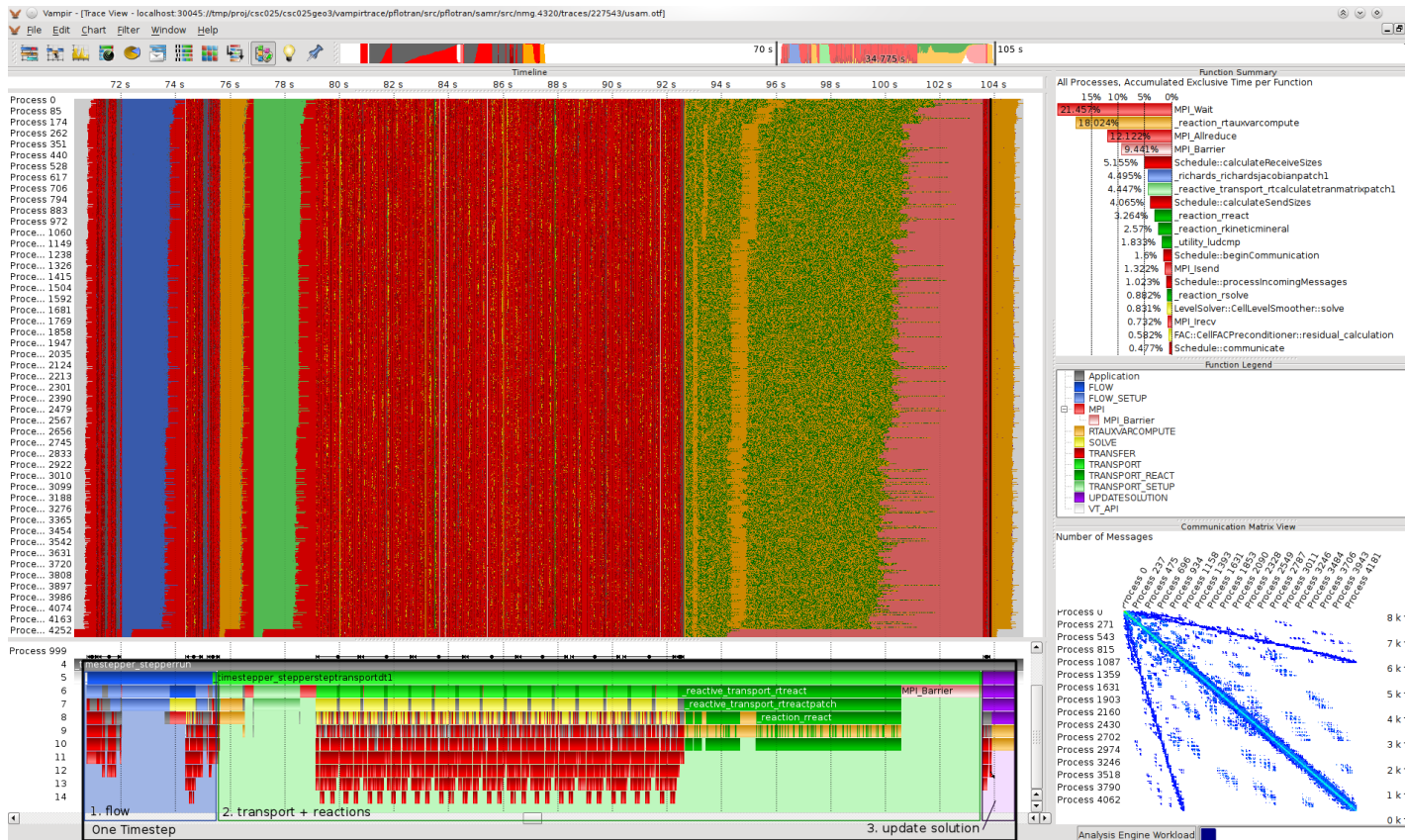  - Adapt statistics to zoom level (time interval)
  - Also for very large and highly parallel traces

TECHNISCHE UNIVERSITÄT DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Vampir Workflow

# Vampir

GUI to analyze trace files (OTF)

Main concept: Timeline + statistics

GUI is QT based – available on Linux, Mac, Windows

# VampirServer

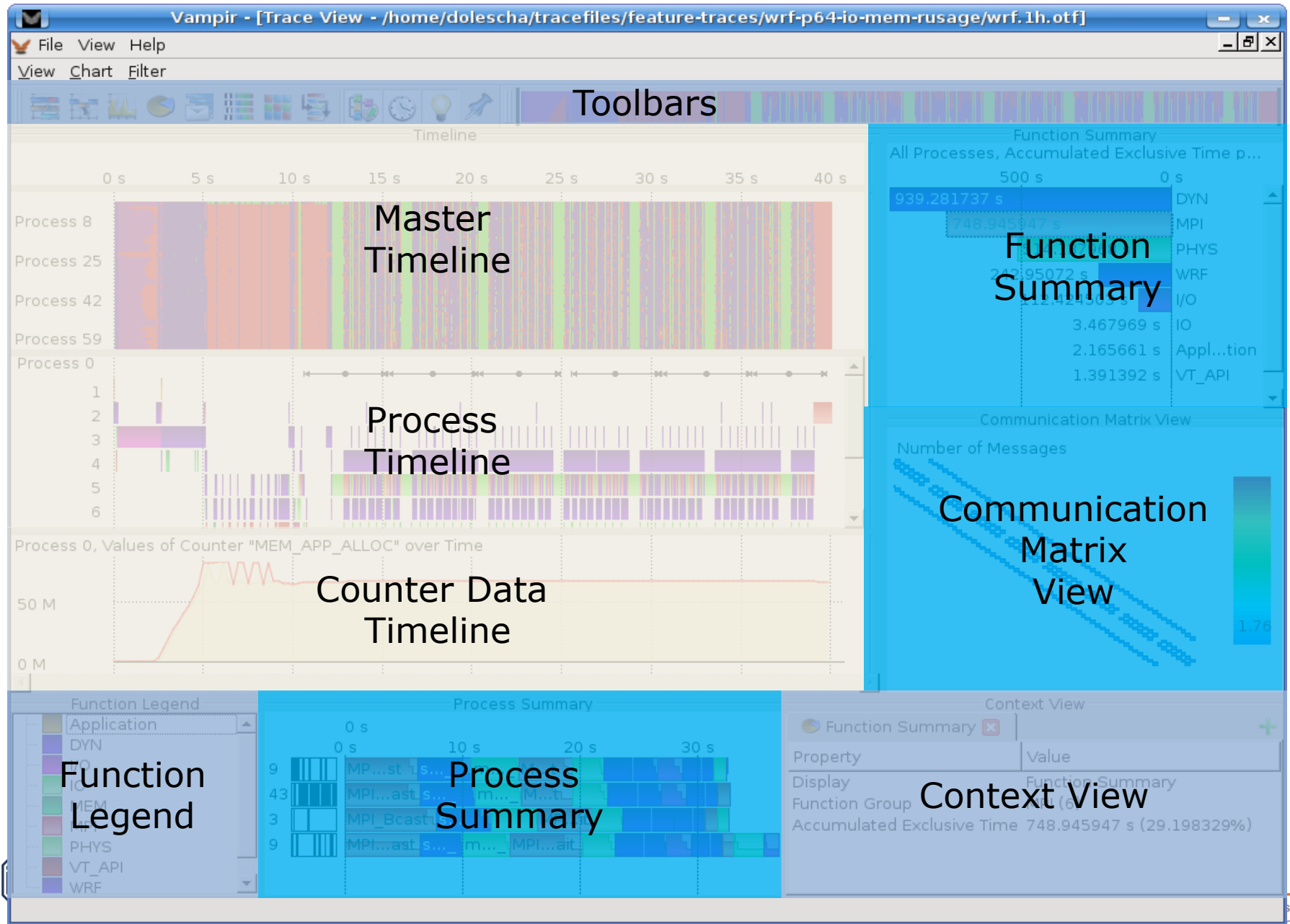**Parallel analysis engine for Vampir**

- MPI
- pthreads

**Scales to > 10,000 analysis processes**

**Loads the entire uncompressed trace into memory**

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Vampir 7: Displays for a WRF trace with 64 processes

# Current Scalability Features in Vampir

## Fit to chart height feature of master timeline and performance radar
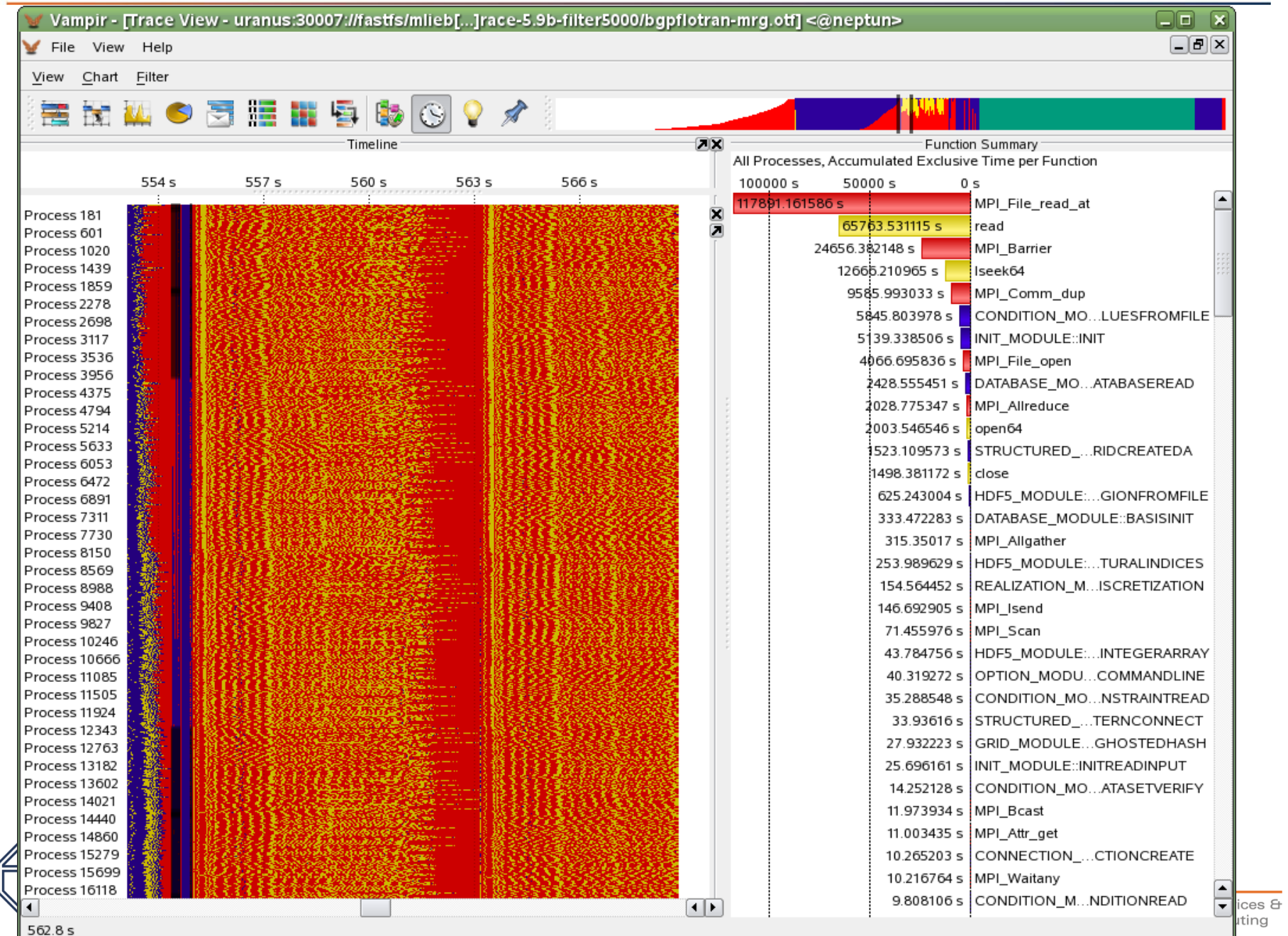
- Allows visualization of more processes than pixels of screen are available

## Clustering

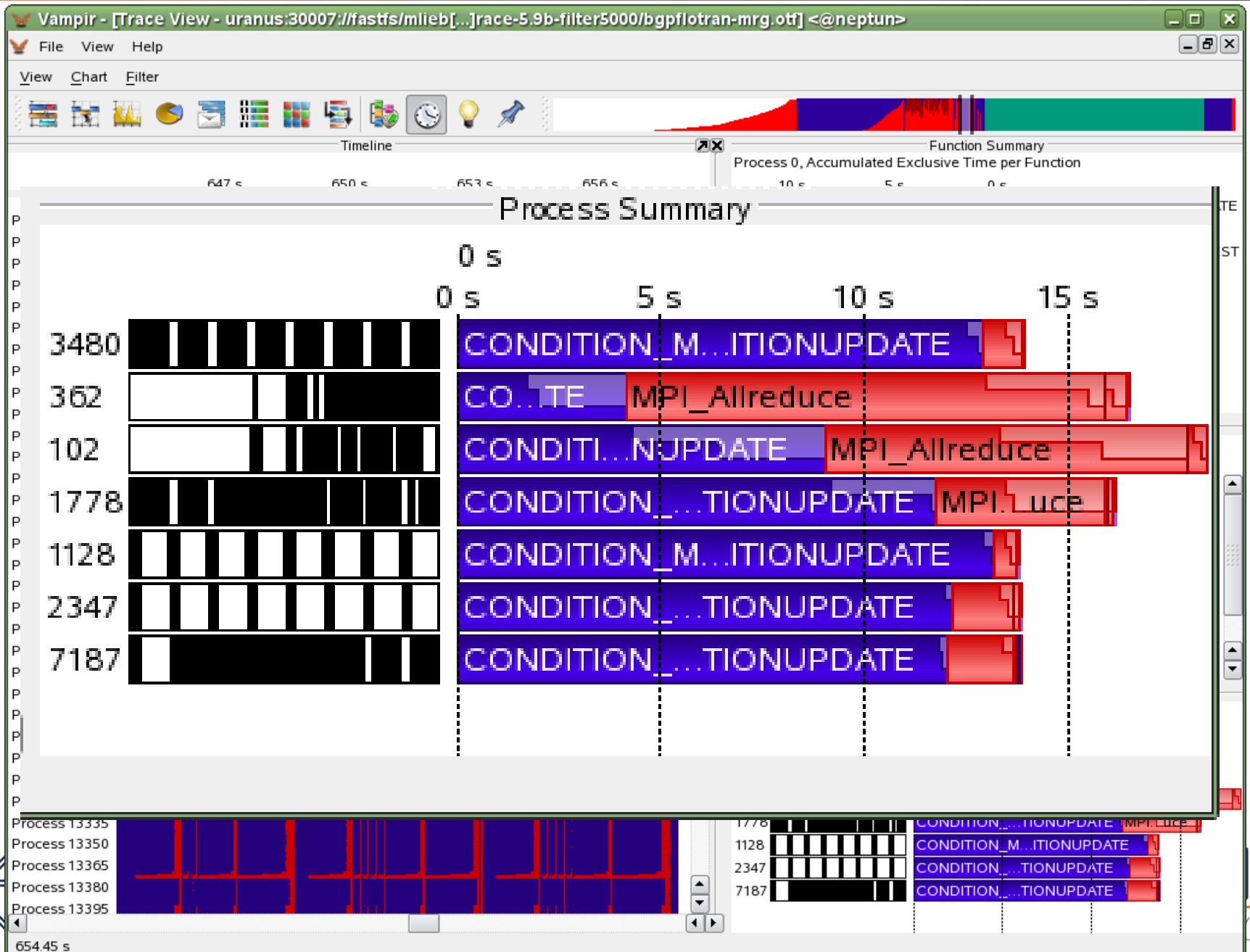- Allows detection of outlier processes and groups with similar behavior

## Performance radar

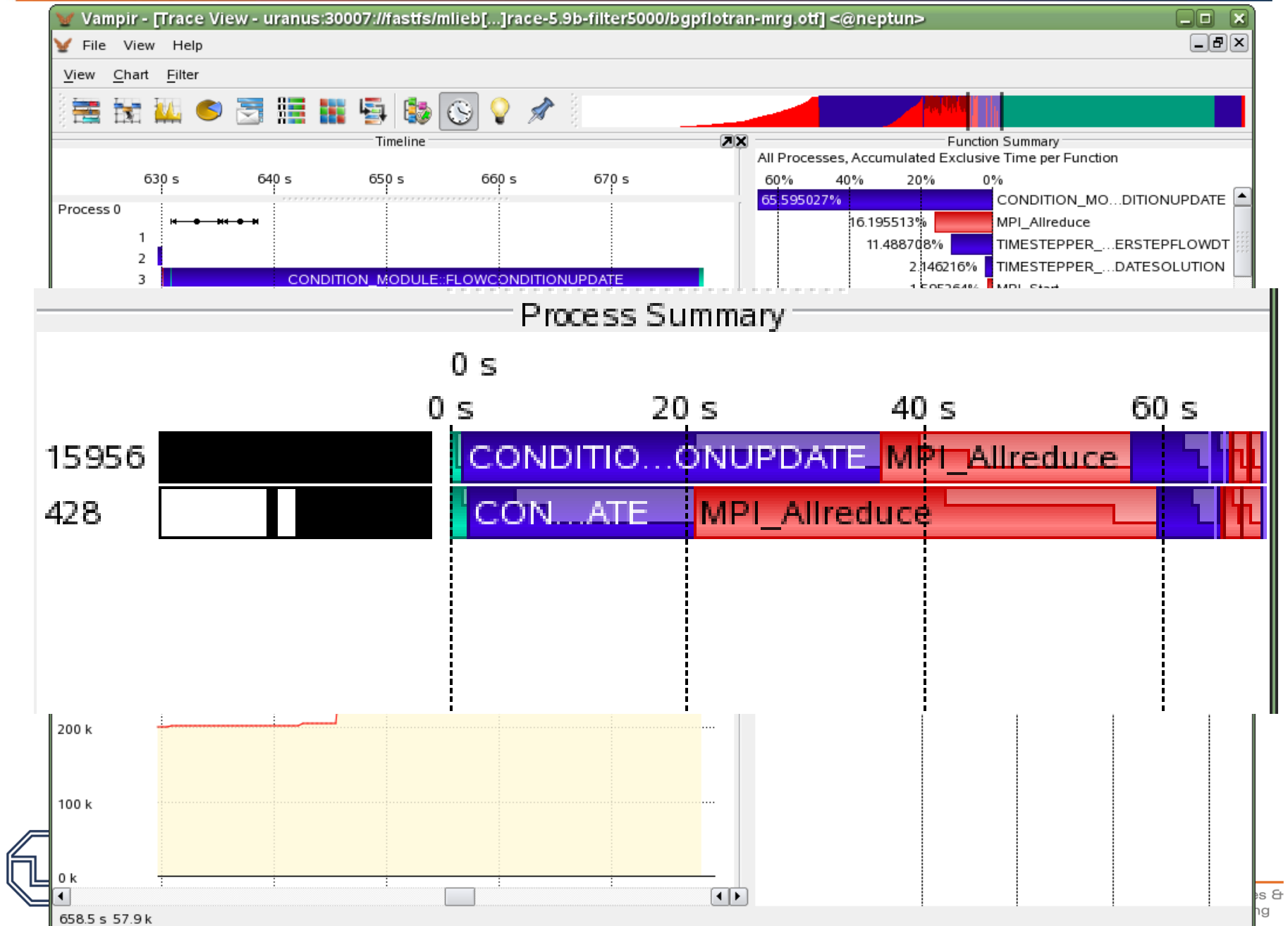- Highlighting performance conditions of your program in a global timeline

**TECHNISCHE UNIVERSITÄT DRESDEN**

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

# Scalability Feature I: Fit to chart height: Pflotran initialization + I/O

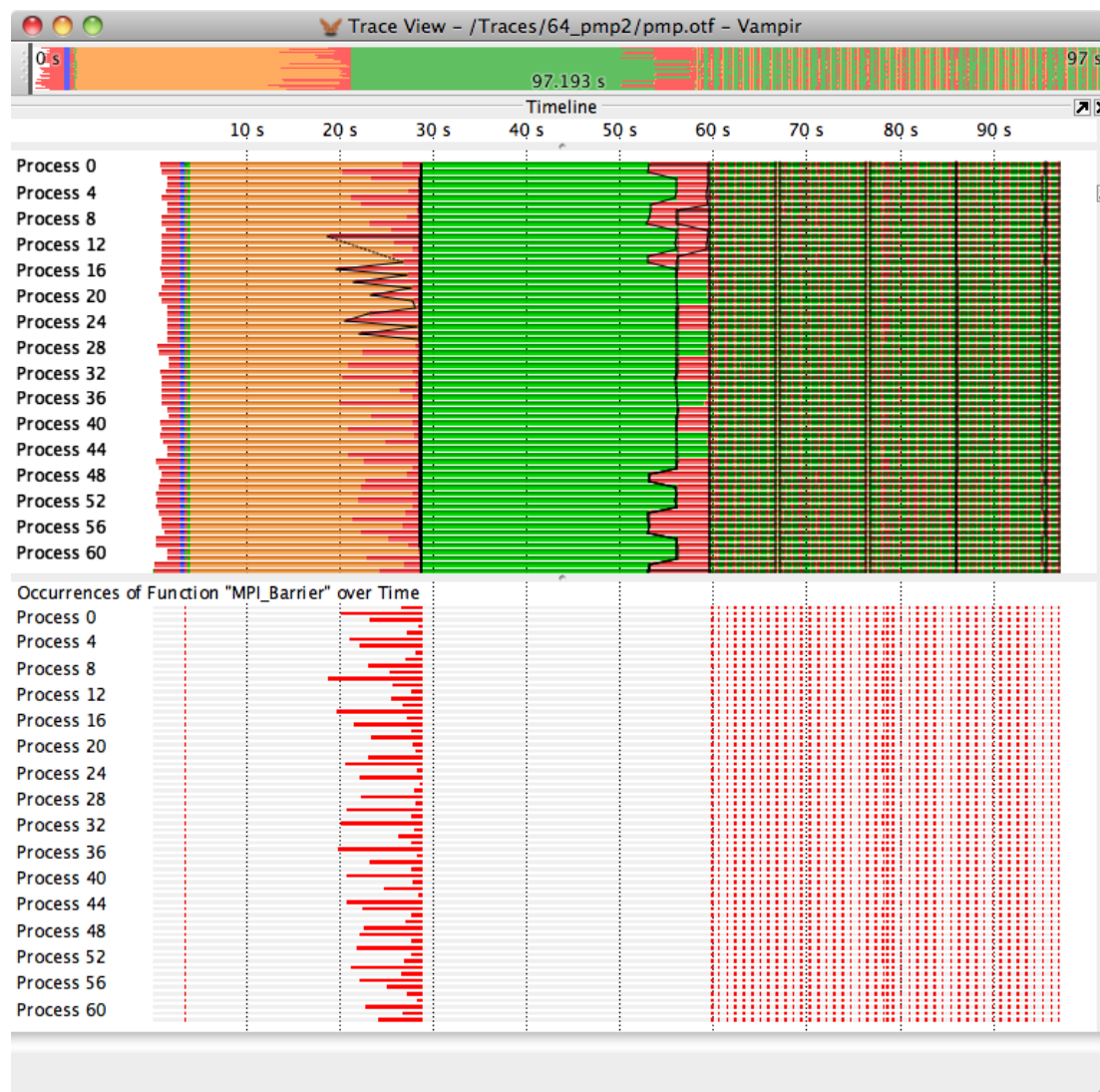# Scalability Feature III: Performance Radar

Display objectives:

- Identification of performance relevant trace parts

- Assistance to users to navigate in trace data and to spot interesting sections

- Performance of basic arithmetics on counter data

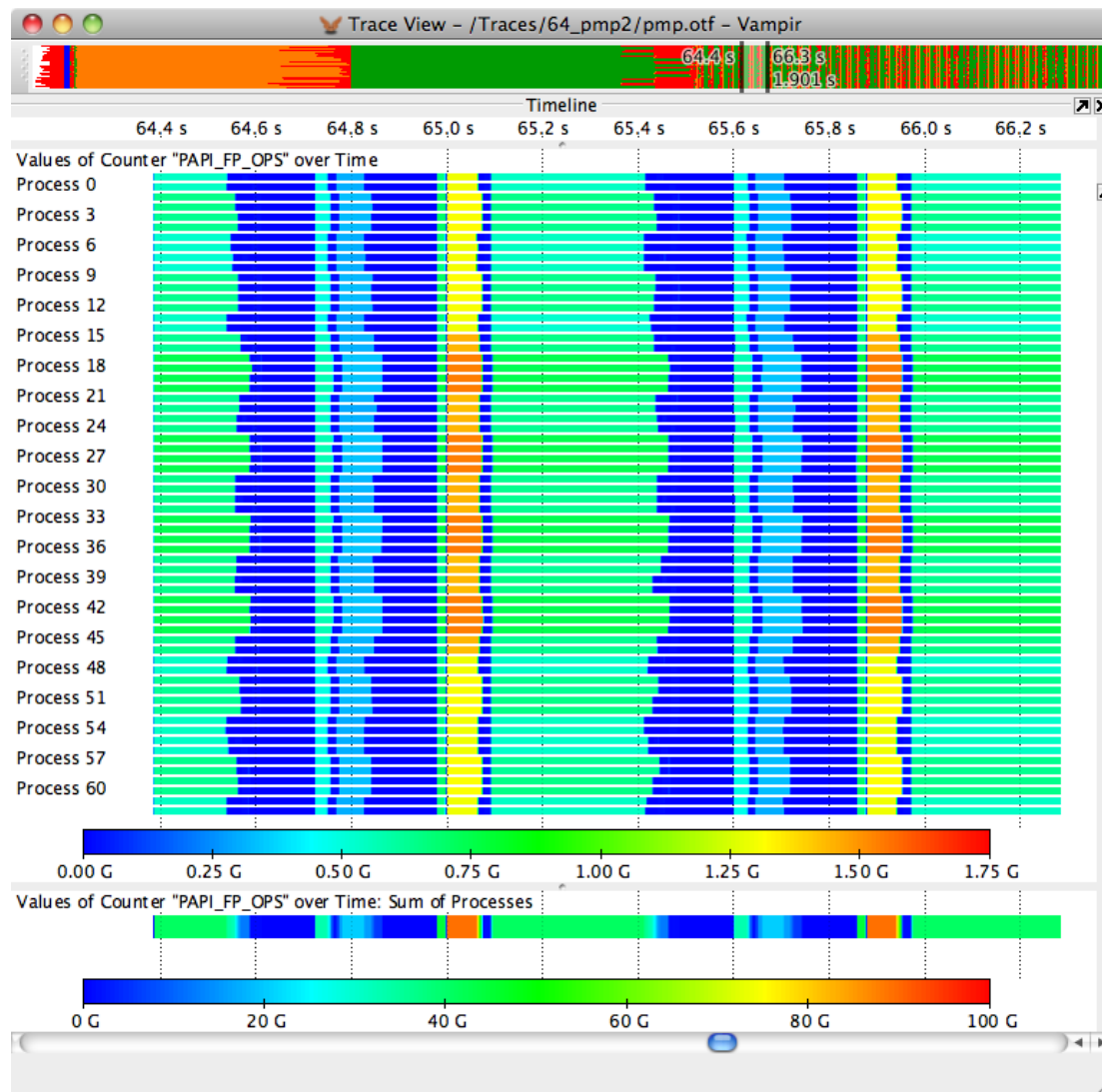# Scalability Feature II: Performance Radar

Features:

- Detection of occurrences of functions/function groups

- Visualization of call density of functions/function groups to help to find performance relevant candidates

- Construction of filter based on function occurrences over time for further usage in calculations
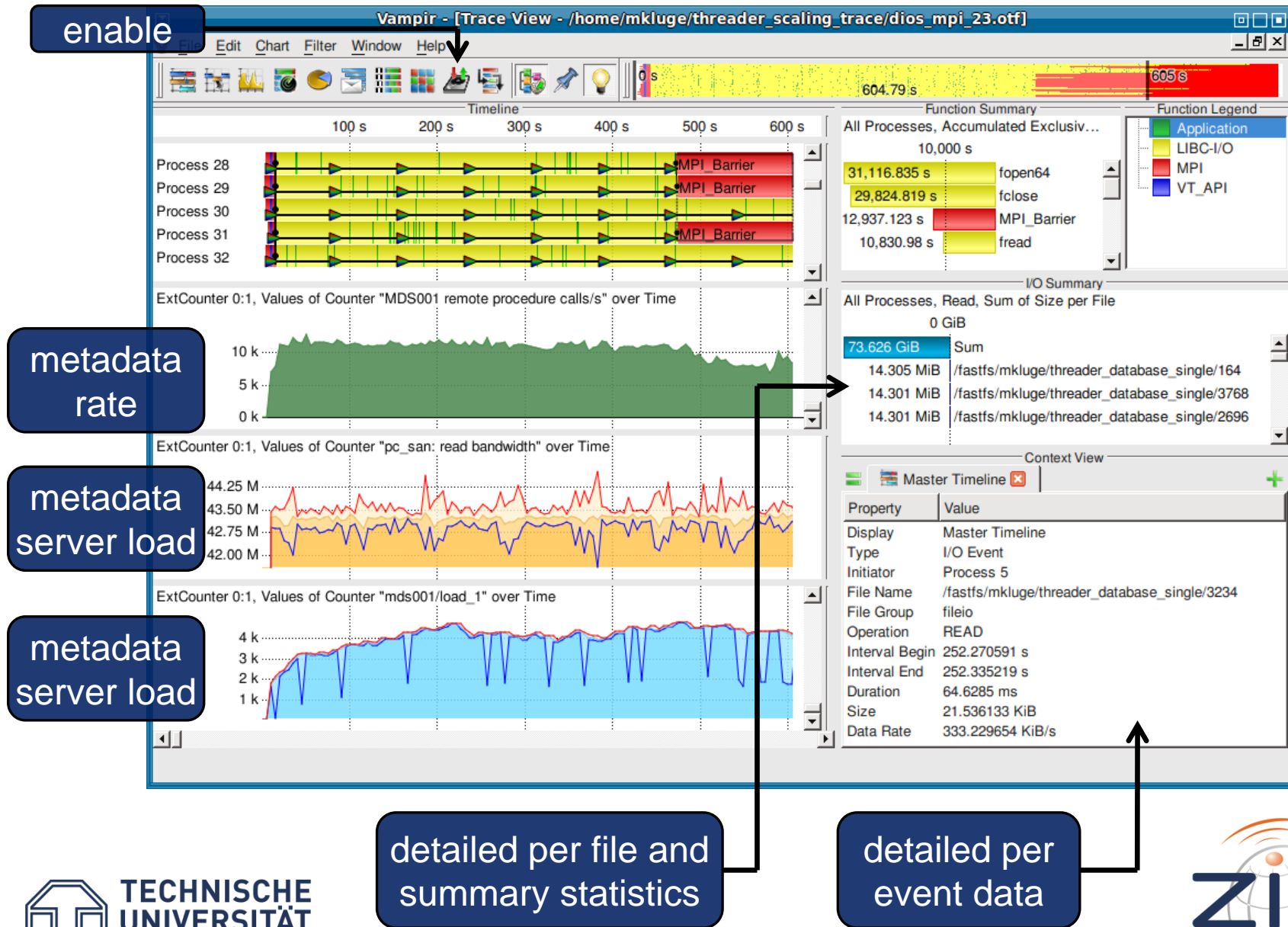
# Scalability Feature II: Performance Radar

Features:

- Performance of arbitrary calculations on counter data, e.g. add up all floating point operations over time, differentiation of performance counter

- Support for concatenation of several calculations

- Utilization of filter in calculations, e.g. only add up FLOPS of function x

TECHNISCHE UNIVERSITÄT DRESDEN

ZIH
Center for Information Services & High Performance Computing

# Vampir I/O Analysis



Titan Summit @ ORNL, Aug 16, 2011 – Guido Juckeland – Slide 22

# Finding Performance Bottlenecks

Inefficient Communication patterns
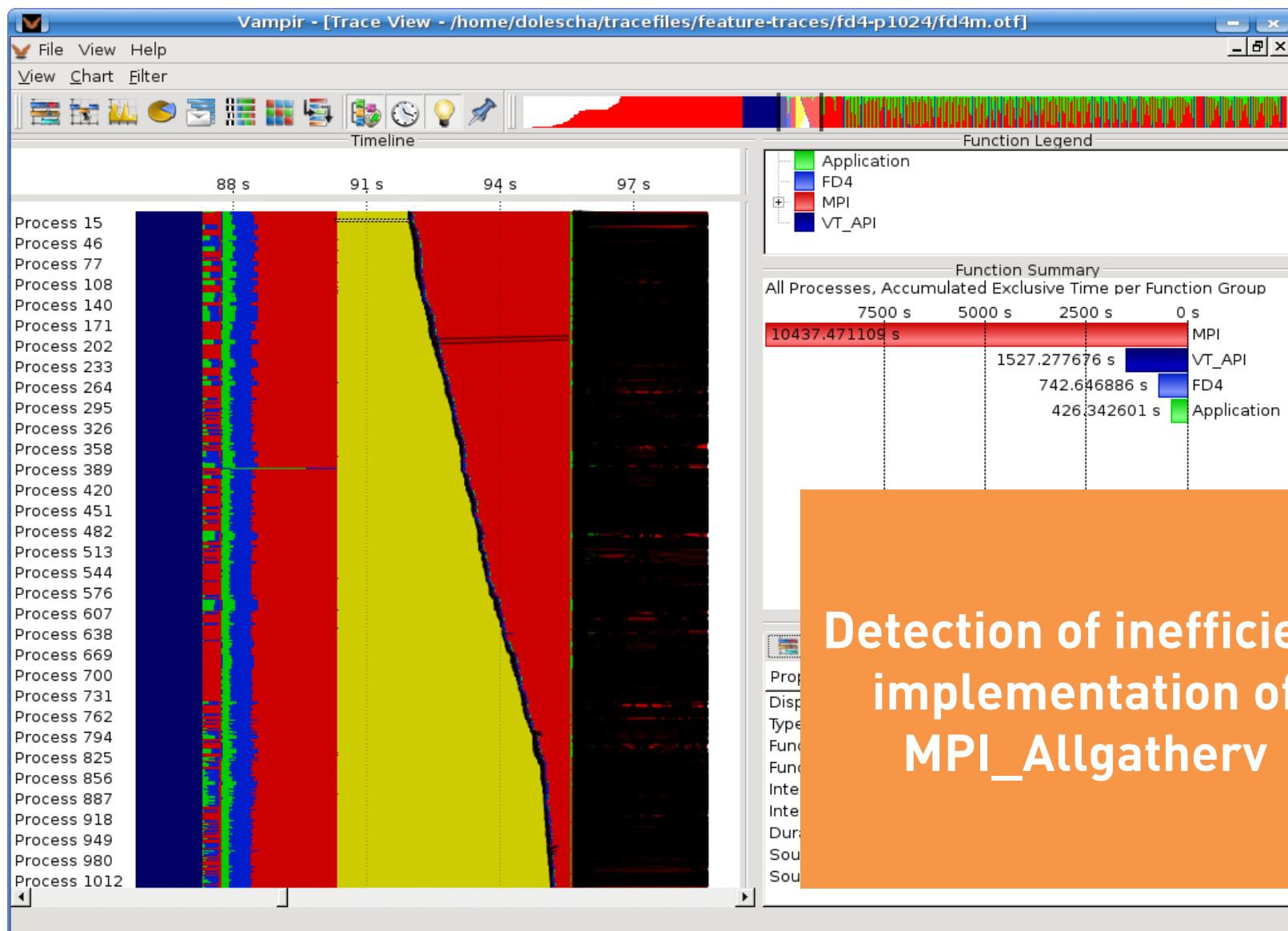
Load imbalance / serial parts of the application

Memory bound computation

- Inefficient cache usage
- TLB misses
- Use HW counters (PAPI) to detect

I/O bottlenecks

Most time consuming function

TECHNISCHE UNIVERSITÄT DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Scalability Feature I: Fit to chart height: Inefficient MPI_Allgatherv

# Challenges on Jaguar and Titan

**Scalability**

- Overcome I/O problems
- Find ways to show that much data

**GPU Support**

- Hitting a moving target
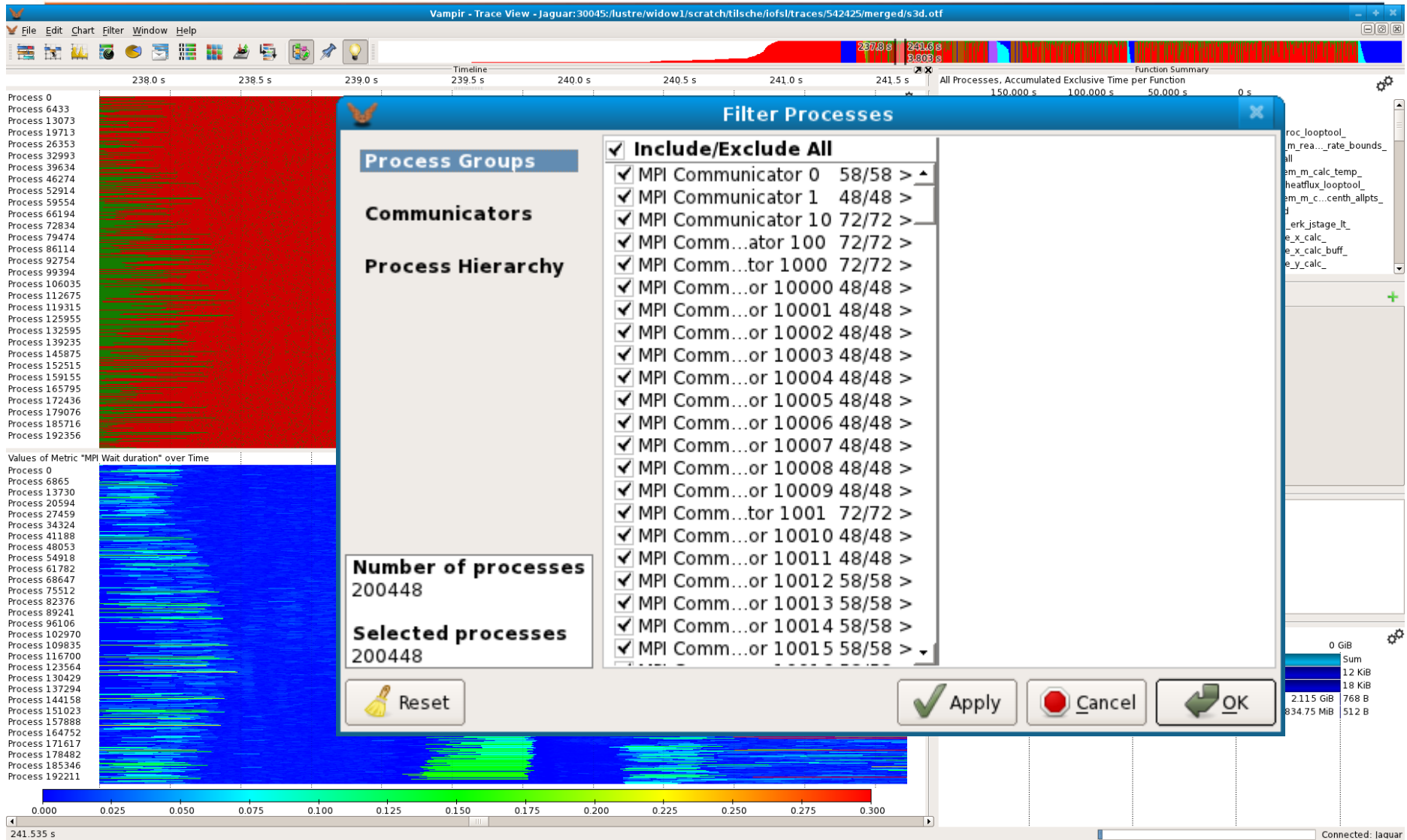- Another layer of heterogenity to display

# Scalability

**Current limits on Jaguar**

- Tracing up to 8000 Processes
  - I/O Problem (too many file creates – one per process)
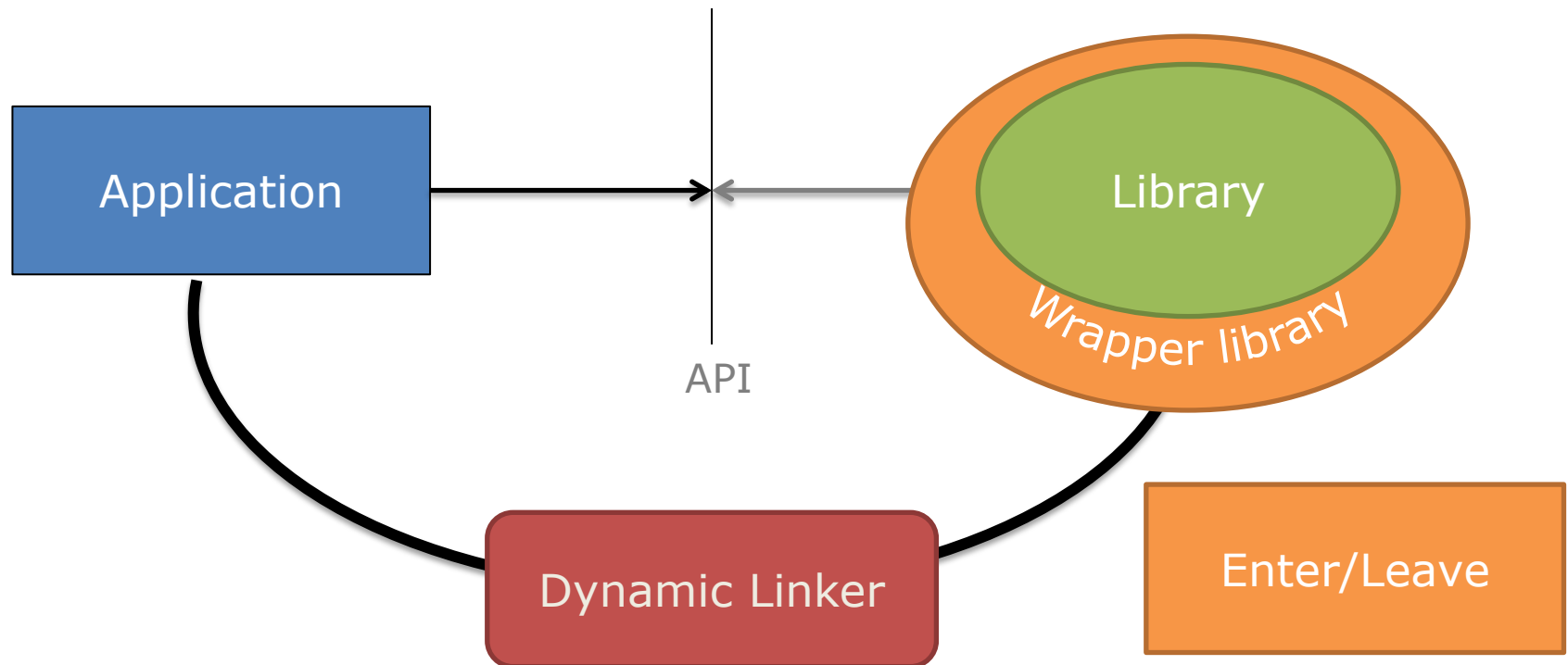
**Prototype already working on Jaguar**

- Tracing 200,000+ processes
- Opened with 20,000 VampirServer processes

TECHNISCHE
UNIVERSITÄT
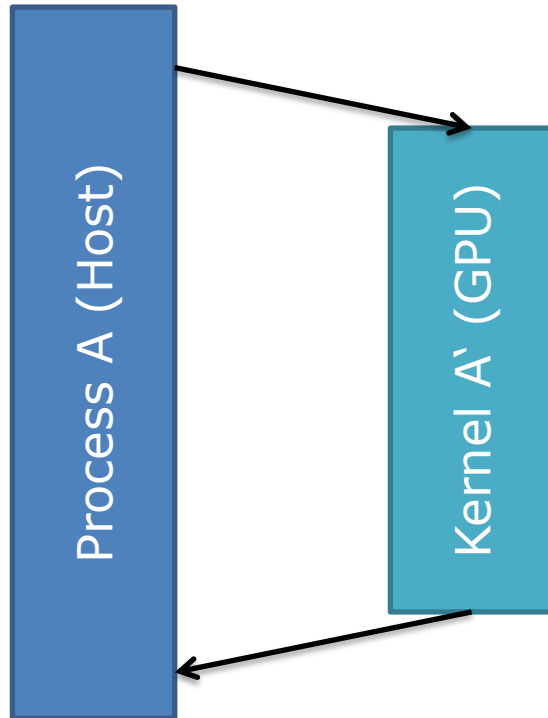DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# 200,000+ Processes in Vampir

## Currently done with Library Wrapping

# Reuse known metrics

Thread = Kernel                    Message = cudaMemCpy

# Asynchronous Events



| foo | launch | foo | Sync_wait | Host |
|-----|--------|-----|-----------|------|

GPU

bar

t

GPU execution queue

bar

# Time stamps for asynchronous events

# PIConGPU – an example for a multi-hybrid application



http://picongpu.fzd.de

# 1 GPU

# MPI + CUDA



Low GPU Utilization

# MPI + pthreads + CUDA

# What are the GPUs doing



Green: Running a kernel          Blue: Idle

# Connection to host processes

# One Process & Thread Group

# Fours steps to get you going at ORNL

**1. Recompile application using VampirTrace**

**2. Run the application**

**3. Start VampirServer**

**4. Connect Vampir to VampirServer**

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing

# Step 1. Recompile application using VampirTrace

First: `module load vampirtrace`

Use the appropriate compiler wrappers

- vtcc, vtCC, vtf77, vtf90

- Pick appropriate library (seq, mpi, mt, hyb)

  – e.g. `vtcc -vt:hyb` (recommended)

- Pick instrumentation type

  – `-vt:inst compinst` (default, compiler instrumentation, all functions)

  – `-vt:inst manual` (MPI, OpenMP, CUDA and manual)

- See what's going on behind `-vt:show` or `-vt:verbose`

- More details `-vt:help`

Tell your build system to use VampirTrace

- `./configure -with-CC="vtcc -vt:hyb" …`

# Step 2. Running the application

Run your application as usual

Make sure the VampirTrace module is loaded

- Load module before and `qsub -V`

- `module load vampirtrace` in the runscript

Control VampirTrace by using environment variables

- `VT_PFORM_GDIR=traces`

  - `mkdir $VT_PFORM_GDIR`

  - `lfs setstripe -c 1 $VT_PFORM_GDIR`

  - place trace files in that directory

  - Use /lustre/widow1 for large runs

**TECHNISCHE UNIVERSITÄT DRESDEN**

ZIH
Center for Information Services & High Performance Computing

# Step 2. Running the application (contd.)

Control VampirTrace by using environment variables

- `VT_FILTER_SPEC=default.filter`
  - `cp $VAMPIRTRACE_DIR/etc/default.filter .`
  - Reduces resulting trace size by filtering frequently called functions
- `VT_MAX_FLUSHES=20`
  - Defaults to 1
  - Allows VT to flush buffer during execution
- `VT_BUFFER_SIZE=100M`
  - Defaults to 32M
  - Increases internal buffer size
- `VT_SYNC_FLUSH=yes`
  - Use that if you have collective ops in your code
  - Avoids all processes waiting for one process to do trace I/O
  - But has overhead itsself (extra allreduce on all collective ops!)
- `VT_METRICS=PAPI_FP_OPS:PAPI_TOT_INS`

**TECHNISCHE UNIVERSITÄT DRESDEN**

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

# Step 3+4. Analyze small traces locally (‹ 100 MB)

Copy the Vampir GUI / Client to your workstation/laptop

- `scp home.ccs.ornl.gov:/sw/sources/vampir/client/`
  `vampir-7.*-i386.tar.gz .`

- `tar -xzf vampir-7.*-i386.tar.gz`

- `./vampir/bin/vampir`

Client is available for:

- Linux {i386, amd64}.tar.gz

- Windows {x86, x64}.exe

- MacOS as .dmg

Copy trace files to your workstation/laptop

- `scp -r jaguarpf.ccs.ornl.gov:/tmp/…/traces/ .`

Start Vampir on your workstation/laptop

- Open trace and enjoy!

# Step 3. Start VampirServer

VampirServer runs on the compute nodes and needs access to the generated trace files (.otf, .z)

Currently available on jaguar, smoky, lens

```
$ module load vampir
    $ qsub -V -I -l size=120,walltime=2:0:0 -A <ACC>
    $ vngd-start.sh
```

Number of processes depends on (size of) the trace

- Try ¼ of the cores used by the application

- Running out of memory → add nodes

- Getting strange MPI errors → reduce process count

**TECHNISCHE UNIVERSITÄT DRESDEN**

ZIH
Center for Information Services &
High Performance Computing

# Step 4. Connect Vampir to VampirServer (Linux/Mac)

Open an SSH tunnel to the compute node

- `jaguarpf$` `vngd-start.sh`
  `Launching VampirServer Version 2.3.0 with  on 12 processes ...`
  `Found license file: /tmp/work/tilsche/.vampir/etc/lic.dat`
  `Running 11 analysis processes... (abort with Ctrl-C or vngd-`
  `shutdown)`
  `Server listens on:` <u>`nid13928:30051`</u>

- `workstation$` `ssh -L 30051:`<u>`nid13928:30051`</u>` jaguarpf.ccs.ornl.gov`

Open Vampir and "remote open" to

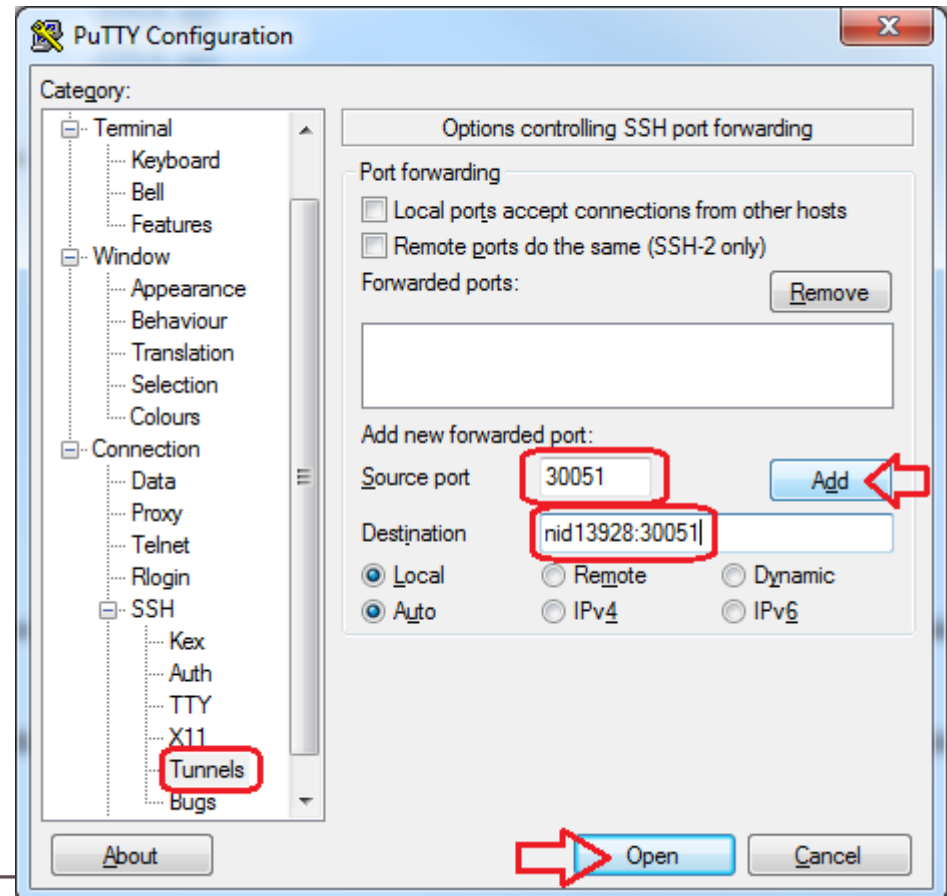- Server: `localhost`

- Port: `30051`

# Step 4. Connect Vampir to VampirServer (Windows)

We need an SSH tunnel to the compute node

- `jaguarpf$ vngd-start.sh`
  ```
  Launching VampirServer Version 2.3.0 with  on 12 processes ...
  Found license file: /tmp/work/tilsche/.vampir/etc/lic.dat
  Running 11 analysis processes... (abort with Ctrl-C or vngd-shutdown)
  Server listens on: nid13928:30051
  ```
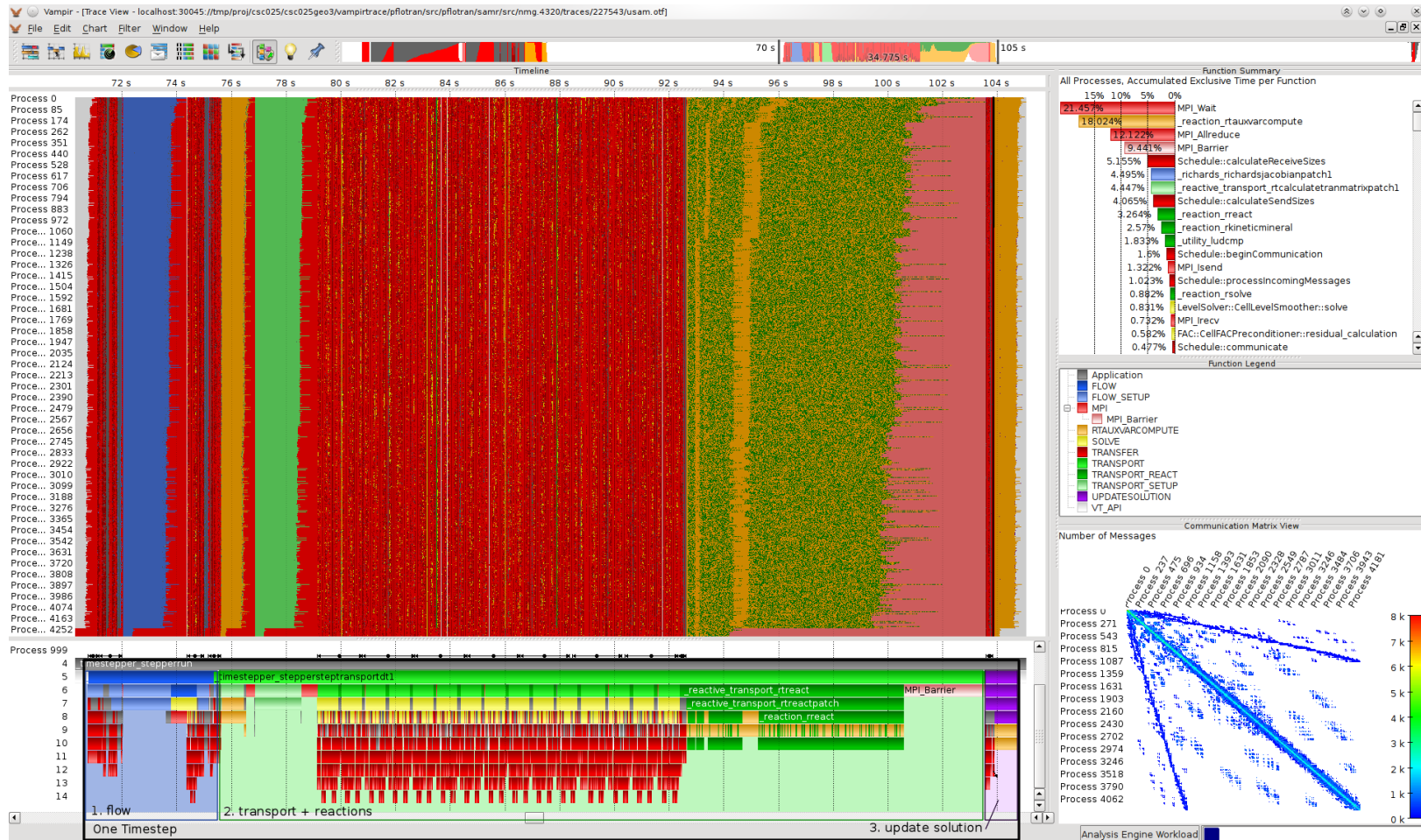
Vampir "remote open" to

- Server: `localhost`

- Port: `30051`

TECHNISCHE UNIVERSITÄT DRESDEN

High Performance Computing

# View your application trace:

"Play around" with Vampir to get a feeling about the features

# Caution: Effects due to Tracing

I/O overhead (flush)

- Visibly marked in the trace

- 'Long' time for I/O

- Ideally only once at the end (invisible) or during barriers

- Avoid by applying runtime filters

Measurement overhead

- Overhead on function calls

- Invisible

- Avoid instrumenting tiny frequently called functions

- Compare total runtime to get an upper bound on overhead

**TECHNISCHE UNIVERSITÄT DRESDEN**

ZIH
Center for Information Services &
High Performance Computing

# Summary

- **Vampir & VampirServer**
  - interactive trace visualization and analysis
  - intuitive browsing and zooming
  - scalable to "quite large" trace data sizes (1,5 TByte)
  - scalable to high parallelism (200,000 processes)
- Vampir available for Windows, Linux/Unix and Mac OS X

- **VampirTrace**
  - convenient instrumentation and measurement infrastructure
  - hides away complicated details
  - provides many options and switches for experts
- VampirTrace is part of Open MPI 1.3 and higher

**TECHNISCHE UNIVERSITÄT DRESDEN**

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

# Acknowledgements

This work would have been impossible without the dedication of:

- Matthias Lieber (Tracing & Analysis)
- Matthias Jurenz (VampirTrace Software & Support)
- Matthias Weber (Vampir Software & Support)

- ## The Vampir Team:

Matthias Jurenz, Andreas Knüpfer, Ronny Brendel, Matthias Lieber, Jens Doleschal, Jens Domke, Holger Mickler, Daniel Hackenberg, Michael Heyde, Thomas Ilsche, Guido Juckeland, Dietrich Robert, Johannes Spazier, Michael Kluge, Matthias Müller, Holger Brunst, Ronald Geisler, Reinhard Neumann, Heide Rohling, Rene Widera, Thomas Ilsche, Matthias Weber, Bert Wesarg, Hartmut Mix, Thomas William, Wolfgang E. Nagel

# Further information about VampirTrace & Vampir

- www.nccs.gov/computing-resources/jaguar/software/?software=vampir

- vampir-7.4.0-OLCF3/doc/Manual.pdf

- www.vampir.eu

- /sw/sources/vampirtrace/5.11ornl/doc/UserManual.pdf

- www.tu-dresden.de/zih/vampirtrace

TECHNISCHE
UNIVERSITÄT
DRESDEN

ZIH
Center for Information Services &
High Performance Computing