**SOFTWARE**

FAST SOLUTIONS

PRODUCTIVITY

# 8th VI-HPS Tuning Workshop
## hosted by GRS in Aachen

5-9 September 2011

Brian Wylie

Jülich Supercomputing Centre

b.wylie@fz-juelich.de

- Presenters
  - Shirley Moore (University of Tennessee ICL)
  - Sameer Shende (University of Oregon PRL)
  - Tobias Hilbrich & Joachim Protze (TU Dresden)
  - Yury Oleynik & Josef Wiedendorfer (TU München)
  - Wolfgang Frings & Brian Wylie (Jülich Supercomputing Centre)
  - Judit Gimenez & Jesus Labarta (Barcelona Supercomp. Center)

- Thanks
  - Local arrangements & facilities
    - Daniel Becker, Marc-André Hermanns (GRS)
    - Systems: JSC & RWTH
  - Sponsor: Bull
  - You
    - Your questions, suggestions & feedback are highly valued

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

- Funded by Helmholtz Association of German Research Centres

HELMHOLTZ | ASSOCIATION

- Activities
  - Development and integration of HPC programming tools
    - ▸ Correctness checking & performance analysis
  - Training workshops
  - Service
    - ▸ Support email lists
    - ▸ Application engagement
  - Academic workshops

www.vi-hps.org

**VI-HPS**

## Forschungszentrum Jülich
- Jülich Supercomputing Centre

## RWTH Aachen University
- Centre for Computing & Communication

## Technical University of Dresden
- Centre for Information Services & HPC

## University of Tennessee (Knoxville)
- Innovative Computing Laboratory

## German Research School
- Laboratory of Parallel Programming

## Technical University of Munich
- Chair for Computer Architecture

## University of Oregon
- Performance Research Laboratory

## University of Stuttgart
- HPC Centre

- KCachegrind
  - Callgraph-based cache analysis
- Marmot/MUST
  - MPI 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/VampirTrace
  - Event tracing and graphical trace visualization & analysis

- 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 & analyse/tune your own code(s)!

- VI-HPS Tuning Workshop series
  - Aachen (3/08), Dresden (10/08), Jülich (2/09), Bremen (9/09), Garching (3/10), Amsterdam (05/10), Stuttgart (03/11)
- VI-HPS Tutorial series
  - SC'08, ICCS'09, SC'09, Cluster'10, SC'10, **SC'11**
- Training with individual tools & platforms (e.g., BlueGene)

- SC'11 tutorial (13 Nov 2011, Seattle, WA, USA)
  - full-day hands-on tutorial using Live DVD
  - "Practical hybrid parallel application performance engineering"

- Further events to be determined
  - (one-day) tutorials
    - ▸ with guided exercises using Live DVD
  - (multi-day) training workshops
    - ▸ with your own applications on real 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 x86-architecture notebook computer
  - VI-HPS tools: KCachegrind, Marmot, PAPI, Periscope, Scalasca, TAU, VT/Vampir*
  - Also: Eclipse/PTP, TotalView*, etc.
    - ▶ * 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

September 2010

**POINT VI-HPS**

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

## Monday 5 Sept.

- 09:00 (early registration & set-up, individual preparation)
- 12:00-13:30 (lunch)
- Welcome & introduction to VI-HPS
- Introduction to parallel performance analysis
- 15:00-15:30 (break)
- Overview of VI-HPS tools
- Lab setup
- 17:30 (adjourn)

- 19:00 Dinner sponsored by Bull, "Im Alten Zollhaus"

Tuesday 6 Sept.

- 09:00-10:30 **Scalasca**
- 11:00-12:30 **Periscope**

Wednesday 7 Sept.

- 09:00-10:30 **TAU**
- 11:00-12:30 **KCachegrind**

Thursday 8 Sept.

- 09:00-10:30 **Vampir**
- 11:00-12:30 **Paraver**

Friday 9 Sept.

- 09:00-10:30 **Marmot / MUST**
- 11:00-12:30 **VI-HPS libraries: PAPI & SIONlib**

- Hands-on exercises part of each tool presentation every morning session

- Hands-on coaching to apply tools to analyse & tune your own codes on workshop HPC systems each afternoon to 17:30

12

- Ensure your application codes build and run to completion with appropriate datasets
  - initial configuration should ideally run in less than 15 minutes with 1-4 compute nodes (up to 96 processes/threads)
    - ▸ to facilitate rapid turnaround and quick experimentation
  - larger/longer scalability configurations are also interesting
    - ▸ turnaround may be limited due to busyness of batch queues
- Compare your application performance on other systems
  - VI-HPS tools already installed on a number of HPC systems
    - ▸ if not, ask your system administrator to install them (or install a personal copy yourself)

# VI-HPS productivity tools

- **KCachegrind**
  - Callgraph-based cache analysis
- **Marmot/MUST**
  - MPI 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/VampirTrace**
  - Event tracing and graphical trace visualization & analysis

Cachegrind: cache analysis by simple cache simulation

- Captures dynamic callgraph
- Based on valgrind dynamic binary instrumentation
- Runs on x86/PowerPC/ARM unmodified binaries
  - ▸ No root access required
- ASCII reports produced

[KQ]Cachegrind GUI

- Visualization of cachegrind output

Developed by TU Munich

- Released as GPL open-source
- http://kcachegrind.sf.net/

Binary → Valgrind

Debug Info

Memory Accesses

Event Counters → Profile

2-level $ Simulator

Event cost tree map

Source code view

Call graph view

Machine code annotation

Tool to check for correct MPI usage at runtime
- Checks conformance to MPI standard
  - ▸ Supports Fortran & C bindings of MPI-1.2
- Checks parameters passed to MPI
- Monitors MPI resource usage

Implementation
- C++ library gets linked to the application
- Does not require source code modifications
- Additional process used as DebugServer
- Results written in a log file (ASCII/HTML/CUBE)

Developed by HLRS & TU Dresden
- Released as open-source
- http://www.hlrs.de/organization/av/amt/projects/marmot

# Marmot logfiles

Next generation MPI runtime error detection tool
- Successor of the Marmot and Umpire tools
- Initial merge of Marmot's many local checks with Umpire's non-local checks
- Improved scalability expected in future

Developed by TU Dresden, LLNL & LANL
- to be released as open-source (BSD license)
- currently in beta-testing for first release in November 2011
- http://tu-dresden.de/.../must

Portable performance counter library & utilities

- Configures and accesses hardware/system counters
- Predefined events derived from available native counters
- Core component for CPU/processor counters
  - ▸ instructions, floating point operations, branches predicted/taken, cache accesses/misses, TLB misses, cycles, stall cycles, …
  - ▸ performs transparent multiplexing when required
- Extensible components for off-processor counters
  - ▸ InfiniBand network, Lustre filesystem, system hardware health, …
- Used by multi-platform performance measurement tools
  - ▸ Periscope, Scalasca, TAU, VampirTrace, ...

Developed by UTK-ICL

- Available as open-source for most modern processors
  http://icl.cs.utk.edu/papi/

# PAPI preset counters (and their definitions)

```
juropa$ papi_avail

Available events and hardware information.
-------------------------------------------------
PAPI Version            : 4.1.0.0
Vendor string and code  : GenuineIntel (1)
Model string and code   : Intel(R) Xeon(R) CPU
                          X5570 @ 2.93GHz (26)
CPU Revision            : 5.000000
CPUID Info              : Family: 6  Model: 26
                          Stepping: 5
CPU Megahertz           : 1600.000000
CPU Clock Megahertz     : 1600
Hdw Threads per core    : 2
Cores per Socket        : 4
NUMA Nodes              : 2
CPU's per Node          : 8
Total CPU's             : 16
Number Hardware Counters : 16
Max Multiplex Counters  : 512
-------------------------------------------------
    Name       Code    Avail Deriv Description

PAPI_L1_DCM  0x80000000  Yes   No
                 Level 1 data cache misses

PAPI_L1_ICM  0x80000001  Yes   No
                 Level 1 instruction cache misses
...
-------------------------------------------------
Of 107 possible events, 35 are available, of
which 9 are derived.
```

```
juropa$ papi_avail -d

...
Symbol       Event Code  Count  |Short Descr.|
 |Long Description|
 |Developer's Notes|
 |Derived|
 |PostFix|
 Native Code[n]: <hex> |name|

PAPI_L1_DCM    0x80000000  1 |L1D cache misses|
 |Level 1 data cache misses|
 ||
 |NOT_DERIVED|
 ||
 Native Code[0]: 0x40002028 |L1D:REPL|

PAPI_L1_ICM    0x80000001  1 |L1I cache misses|
 |Level 1 instruction cache misses|
 ||
 |NOT_DERIVED|
 ||
 Native Code[0]: 0x40001031 |L1I:MISSES|

PAPI_L2_DCM    0x80000002  2 |L2D cache misses|
 |Level 2 data cache misses|
 ||
 |DERIVED_SUB|
 ||
 Native Code[0]: 0x40000437 |L2_RQSTS:MISS|
 Native Code[1]: 0x40002037 |
L2_RQSTS:IFETCH_MISS|

...
```

# PAPI native counters (and qualifiers)

```
juropa$ papi_native_avail

Available native events and hardware information.

...

Event Code   Symbol  | Long Description |
-----------------------------------------------------------------------
0x40000000   UNHALTED_CORE_CYCLES  | count core clock cycles whenever the cloc |
             | k signal on the specific core is running (not halted). Alias to e |
             | vent CPU_CLK_UNHALTED:THREAD                                        |
-----------------------------------------------------------------------
0x40000001   INSTRUCTION_RETIRED  | count the number of instructions at retire |
             | ment. Alias to event INST_RETIRED:ANY_P                            |
-----------------------------------------------------------------------
...
-----------------------------------------------------------------------
0x40000086   UNC_SNP_RESP_TO_REMOTE_HOME  | Remote home snoop response - LLC d |
             | oes not have cache line                                           |
  40000486   :I_STATE  | Remote home snoop response - LLC does not have cache |
             | line                                                             |
  40000886   :S_STATE  | Remote home snoop response - LLC has  cache line in S |
             |  state                                                           |
  40001086   :FWD_S_STATE  | Remote home snoop response - LLC forwarding cache |
             |  line in S state.                                                 |
  40002086   :FWD_I_STATE  | Remote home snoop response - LLC has forwarded a  |
             | modified cache line                                              |
  40004086   :CONFLICT  | Remote home conflict snoop response                 |
  40008086   :WB  | Remote home snoop response - LLC has cache line in the M s |
             | tate                                                             |
  40010086   :HITM  | Remote home snoop response - LLC HITM                     |
-----------------------------------------------------------------------
Total events reported: 135
```

Automated profile-based performance analysis

- Iterative on-line performance analysis
  - ▸ Multiple distributed hierarchical agents
- Automatic search for bottlenecks based on properties formalizing expert knowledge
  - ▸ MPI wait states
  - ▸ Processor utilization hardware counters
- Clustering of processes/threads with similar properties
- Eclipse-based integrated environment

Supports

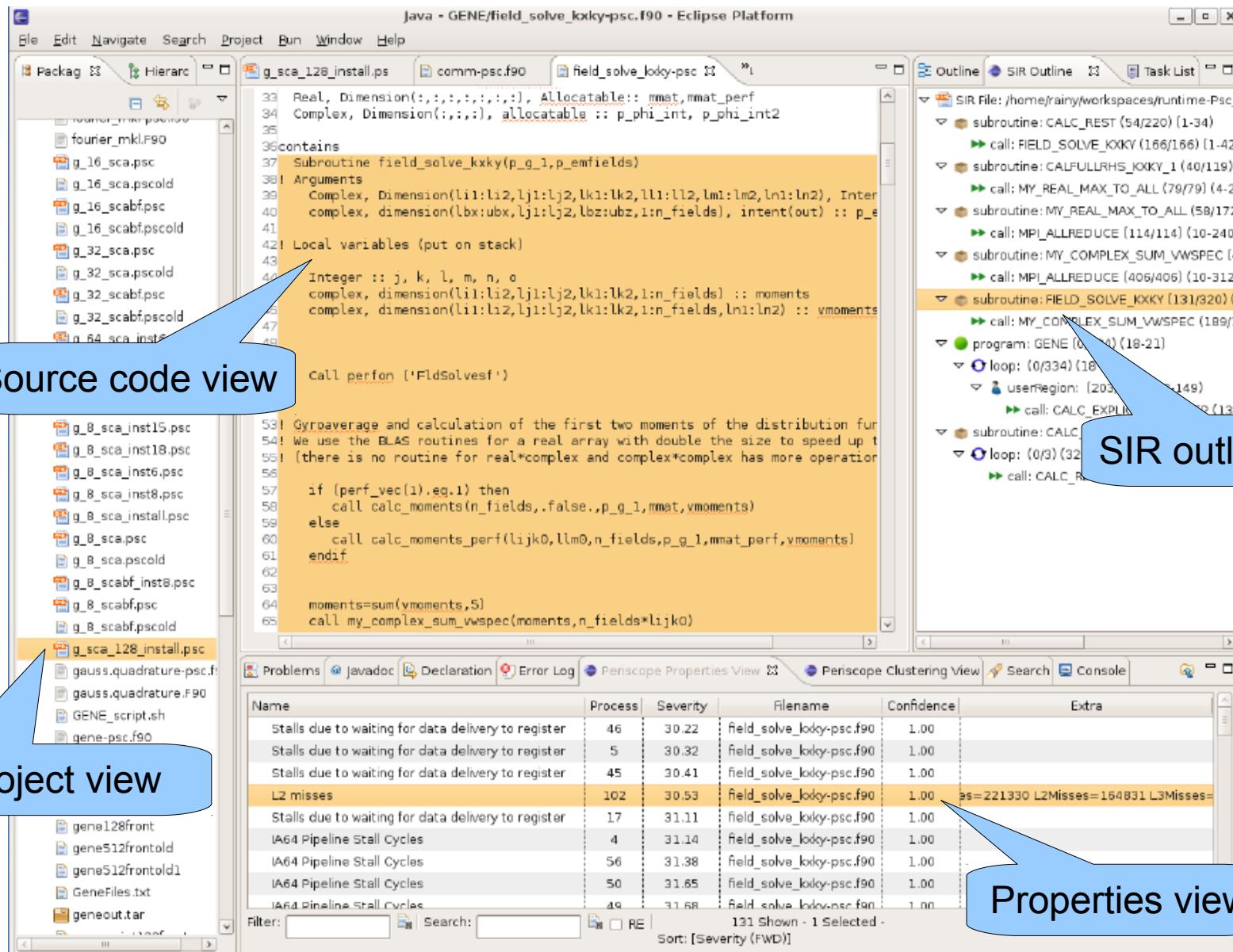- SGI Altix Itanium2, IBM Power and x86-based architectures

Developed by TU Munich

- Released as open-source
- http://www.lrr.in.tum.de/periscope

MPI

- Excessive MPI communication time
- Excessive MPI time due to many small messages
- Excessive MPI time in receive due to late sender
- ...

Hardware performance counters (platform-specific)

- Cycles lost due to cache misses
  - ▸ High L1/L2/L3 demand load miss rate
- Cycles lost due to store instructions
- Cycles lost due to address translation misses
- Cycles lost due to no instruction to dispatch
- ...

# Periscope plug-in to Eclipse environment

Automatic performance analysis toolset

- Scalable performance analysis of large-scale applications
  - ▶ particularly focused on MPI & OpenMP paradigms
  - ▶ analysis of communication & synchronization overheads
- Automatic and manual instrumentation capabilities
- Runtime summarization and/or event trace analyses
- Automatic search of event traces for patterns of inefficiency
  - ▶ Scalable trace analysis based on parallel replay
- Interactive exploration GUI and algebra utilities for XML callpath profile analysis reports

Developed by JSC & GRS

- Released as open-source
- http://www.scalasca.org/

# Scalasca hybrid analysis report

Integrated performance toolkit

- Instrumentation, measurement, analysis & visualization
  - ▸ Highly customizable installation, API, envvars & GUI
  - ▸ Supports multiple profiling & tracing capabilities
- Performance data management & data mining
- Targets all parallel programming/execution paradigms
  - ▸ Ported to a wide range of computer systems
- Performance problem solving framework for HPC
- Extensive bridges to/from other performance tools
  - ▸ PerfSuite, Scalasca, Vampir, ...

Developed by U. Oregon/PRL

- Broadly deployed open-source software
- http://tau.uoregon.edu/

## TAU Architecture

## Program Analysis

## Performance Data Mining

### PerfExplorer

## Parallel Profile Analysis

### PerfDMF

### ParaProf

## Performance Monitoring

### TAUoverSupermon

Interactive event trace analysis

- Alternative & supplement to automatic trace analysis
- Visual presentation of dynamic runtime behaviour
  - ▶ event timeline chart for states & interactions of processes/threads
  - ▶ communication statistics, summaries & more
- Interactive browsing, zooming, selecting
  - ▶ linked displays & statistics adapt to selected time interval (zoom)
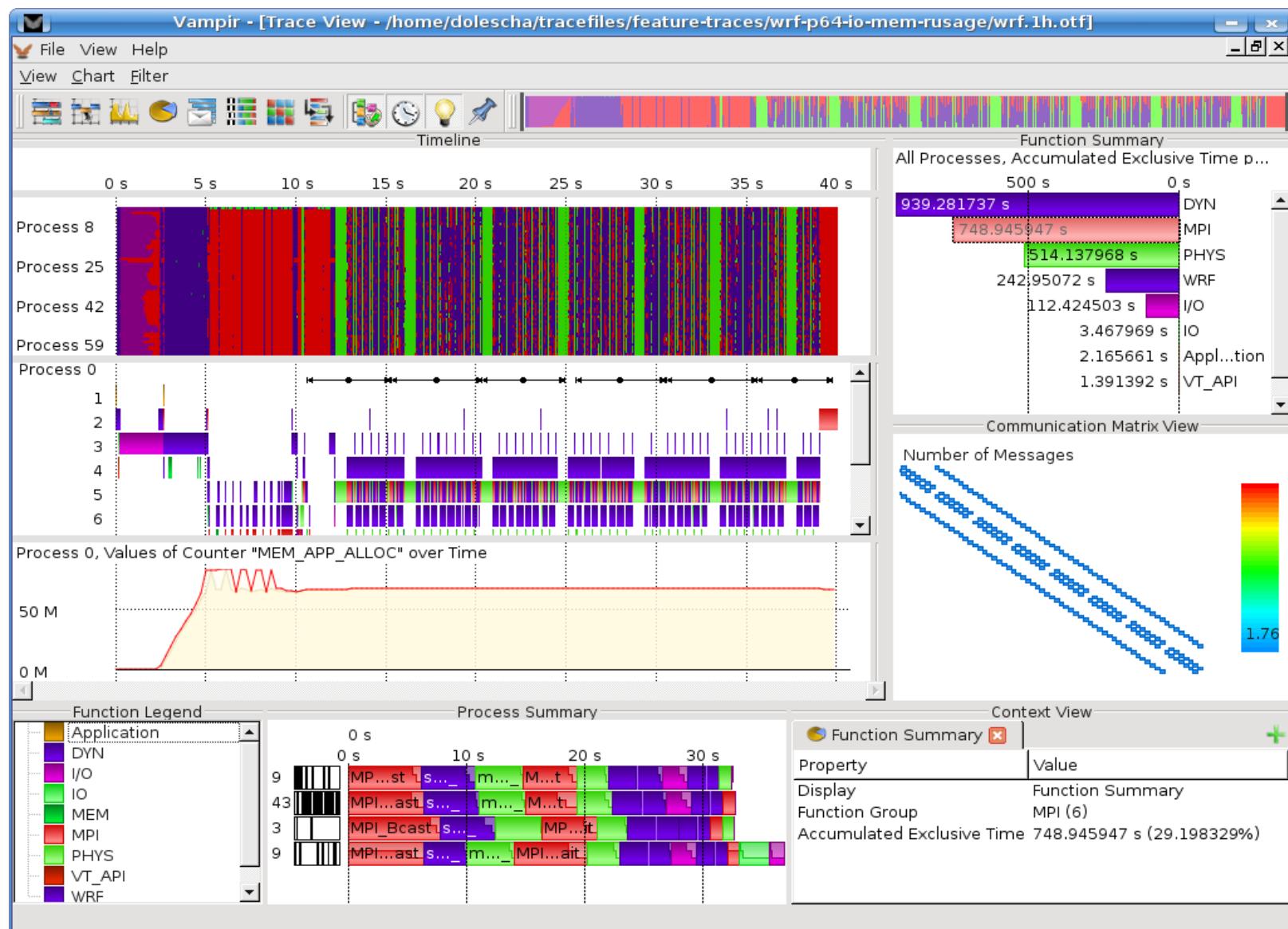  - ▶ scalable server runs in parallel to handle larger traces

Developed by TU Dresden ZIH

- Open-source VampirTrace library bundled with OpenMPI 1.3
- http://www.tu-dresden.de/zih/vampirtrace/
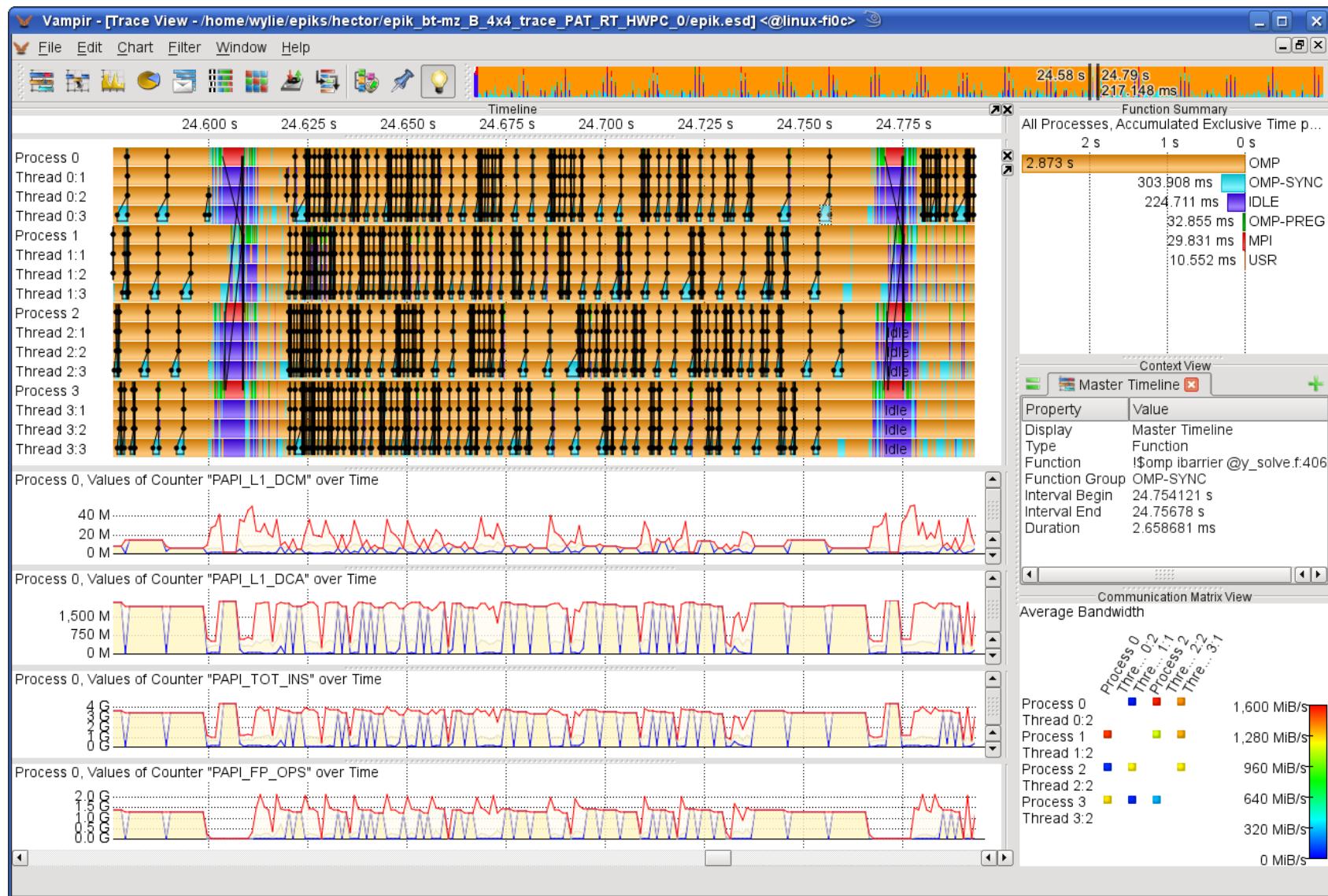- Vampir Server & GUI have a commercial license
- http://www.vampir.eu/

# Vampir interactive trace analysis GUI
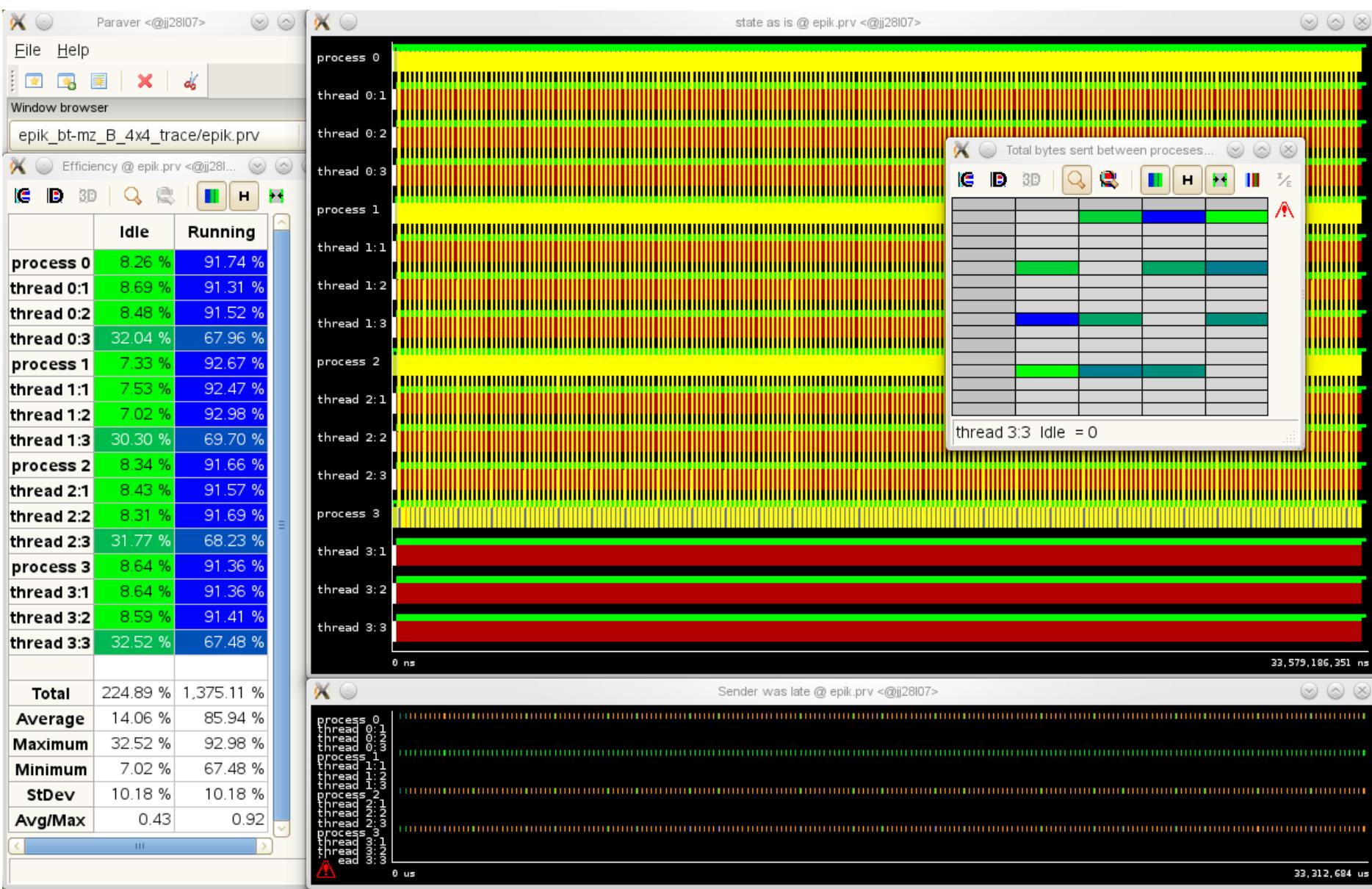
# Vampir interactive trace analysis GUI (zoom)

- Interactive event trace analysis
  - Visual presentation of dynamic runtime behaviour
    - ▸ event timeline chart for states & interactions of processes
    - ▸ Interactive browsing, zooming, selecting
  - Large variety of highly configurable analyses & displays
- Developed by Barcelona Supercomputing Center
  - Paraver trace analyser and Extrae measurement library
  - Open source available from http://www.bsc.es/paraver/

Key tool components also provided as open-source

- Program/library instrumentation
  - ▶ COBI, OPARI, PDToolkit
- MPI library/tool integration
  - ▶ UniMCI
- Scalable I/O
  - ▶ SIONlib
- Libraries & tools for handling (and converting) traces
  - ▶ EPILOG, PEARL, OTF
- Analysis algebra & hierarchical/topological presentation
  - ▶ CUBE

# Portable native parallel I/O library & utilities

- Scalable massively-parallel I/O to task-local files
- Manages single or multiple physical files on disk
  - ▶ optimizes bandwidth available from I/O servers by matching blocksizes/alignment, reduces metadata-server contention
- POSIX-I/O-compatible sequential & parallel API
  - ▶ adoption requires minimal source-code changes
- Tuned for common parallel filesystems
  - ▶ GPFS (BlueGene), Lustre (Cray), ...
- Convenient for application I/O, checkpointing,
  - ▶ Used by Scalasca tracing (when configured)

# Developed by JSC

- Available as open-source from
  http://www.fz-juelich.de/jsc/sionlib/