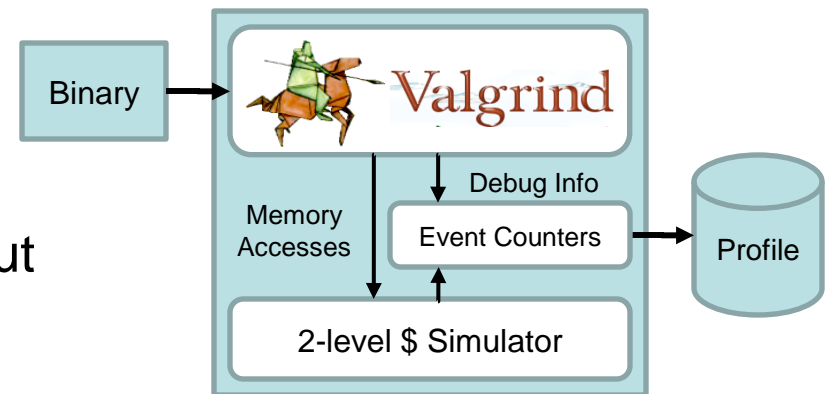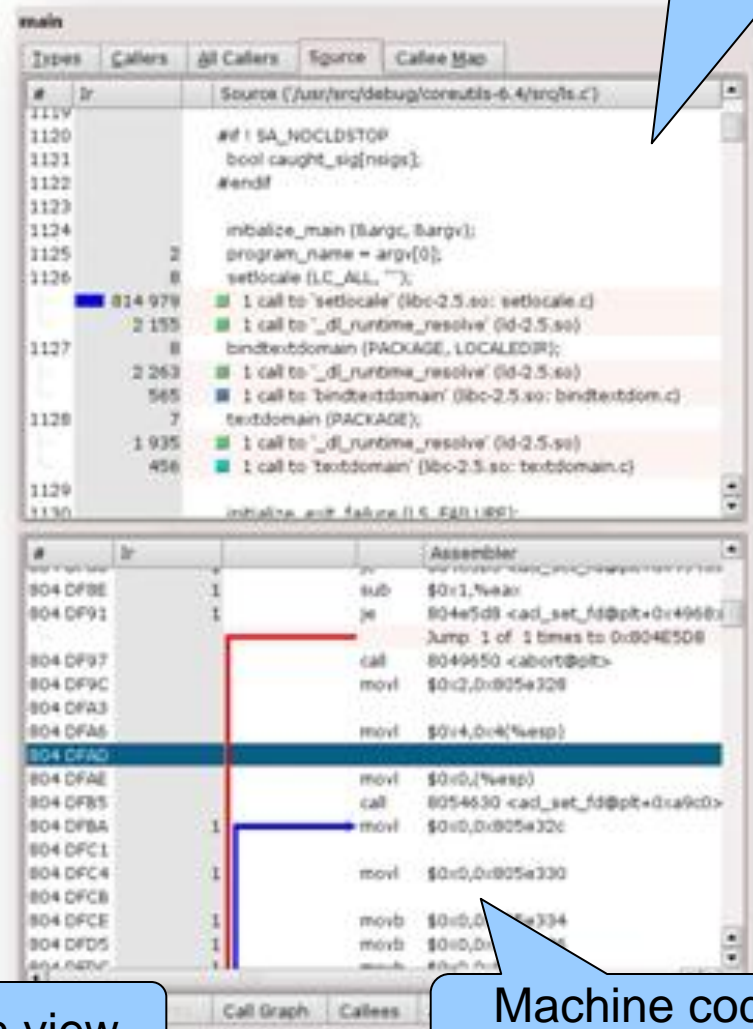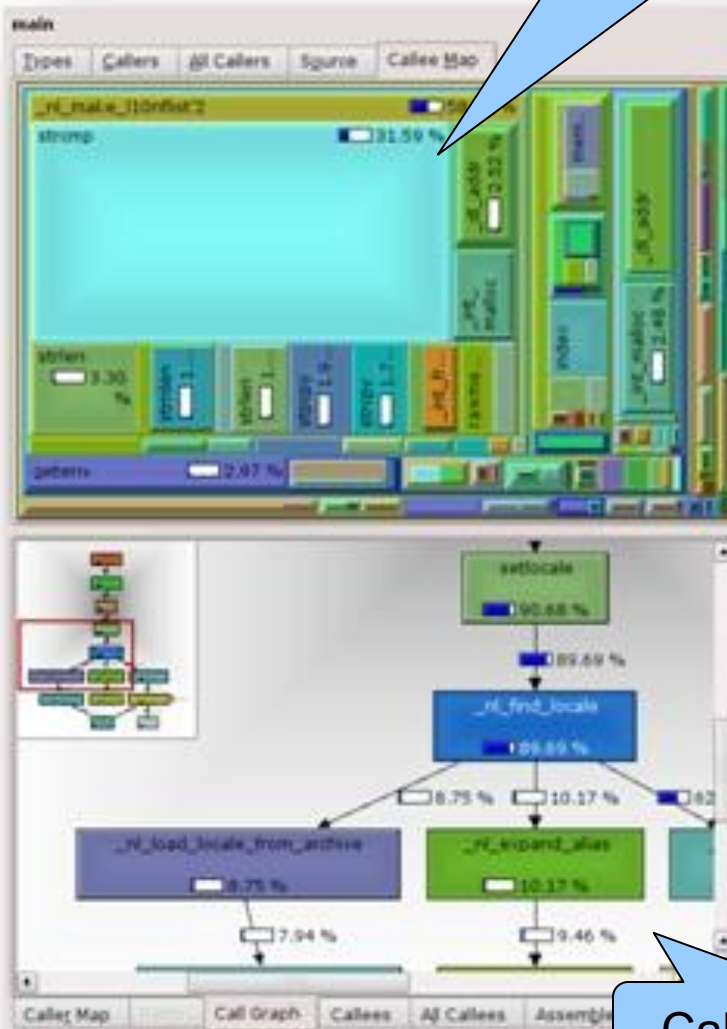# VI-HPS productivity tools suite

Markus Geimer
Jülich Supercomputing Centre

- 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/

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
    - Exploits CMake, GTI & PnMPI infrastructure

- Developed by TU Dresden, LLNL & LANL
    - BSD license open-source initial release in November 2011
    - http://tu-dresden.de/zih/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, OpenMP overheads and imbalances
    - 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
  - …

- OpenMP
  - Load imbalance in parallel region/section
  - Sequential computation in master/single/ordered region
  - ...

- Hardware performance counters (platform-specific)
  - Cycles lost due to cache misses
    - High L1/L2/L3 demand load miss rate
  - Cycles lost due to no instruction to dispatch
  - ...

# Periscope plug-in to Eclipse environment



Source code view

SIR outline view

Project view

Properties view

- 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 automatic trace analysis report

# Scalasca hybrid analysis report

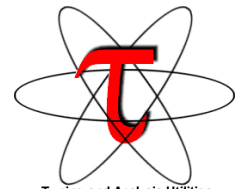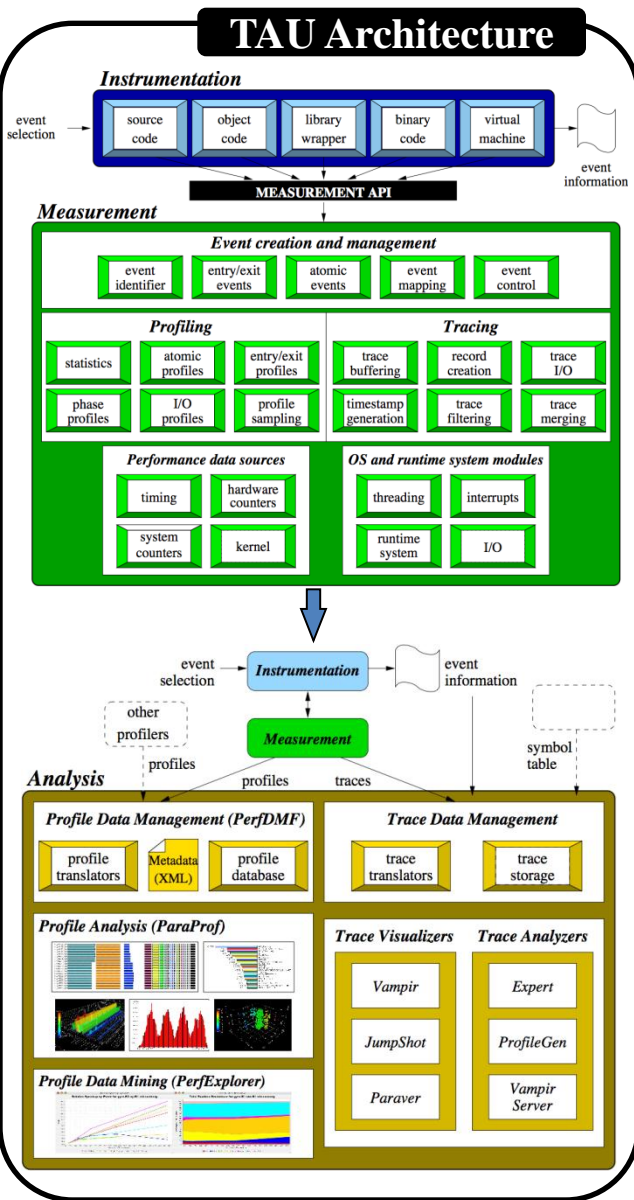# Scalasca automatic trace analysis report

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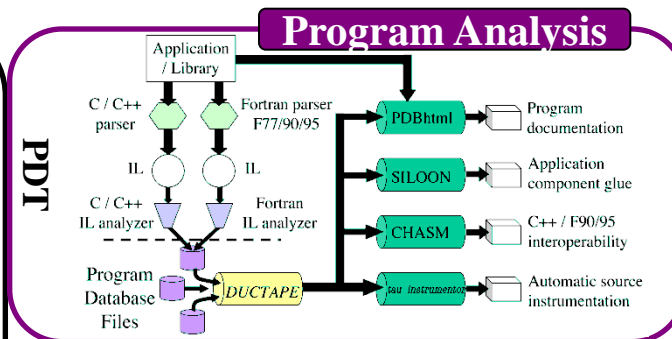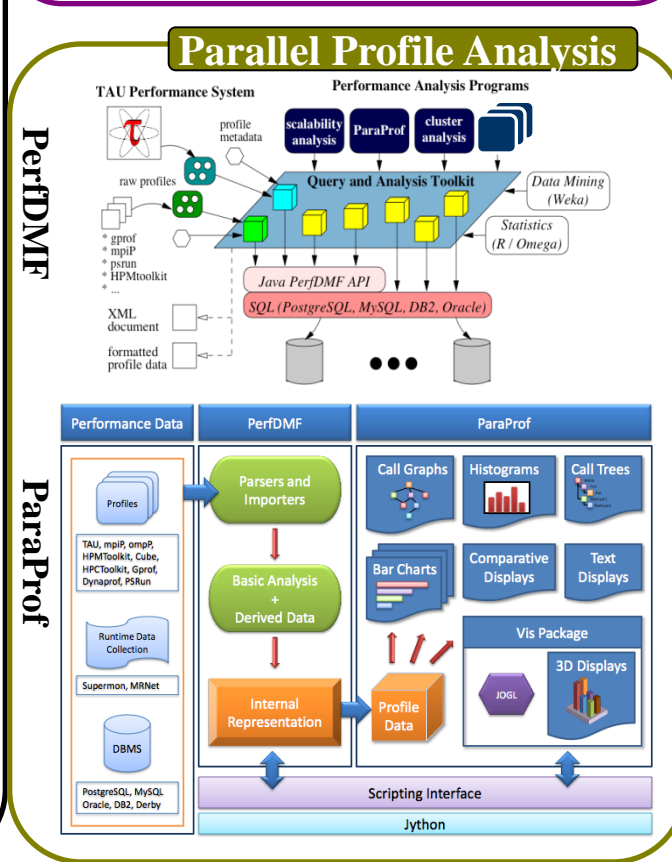  - Broadly deployed open-source software
  - http://tau.uoregon.edu/

- ## 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

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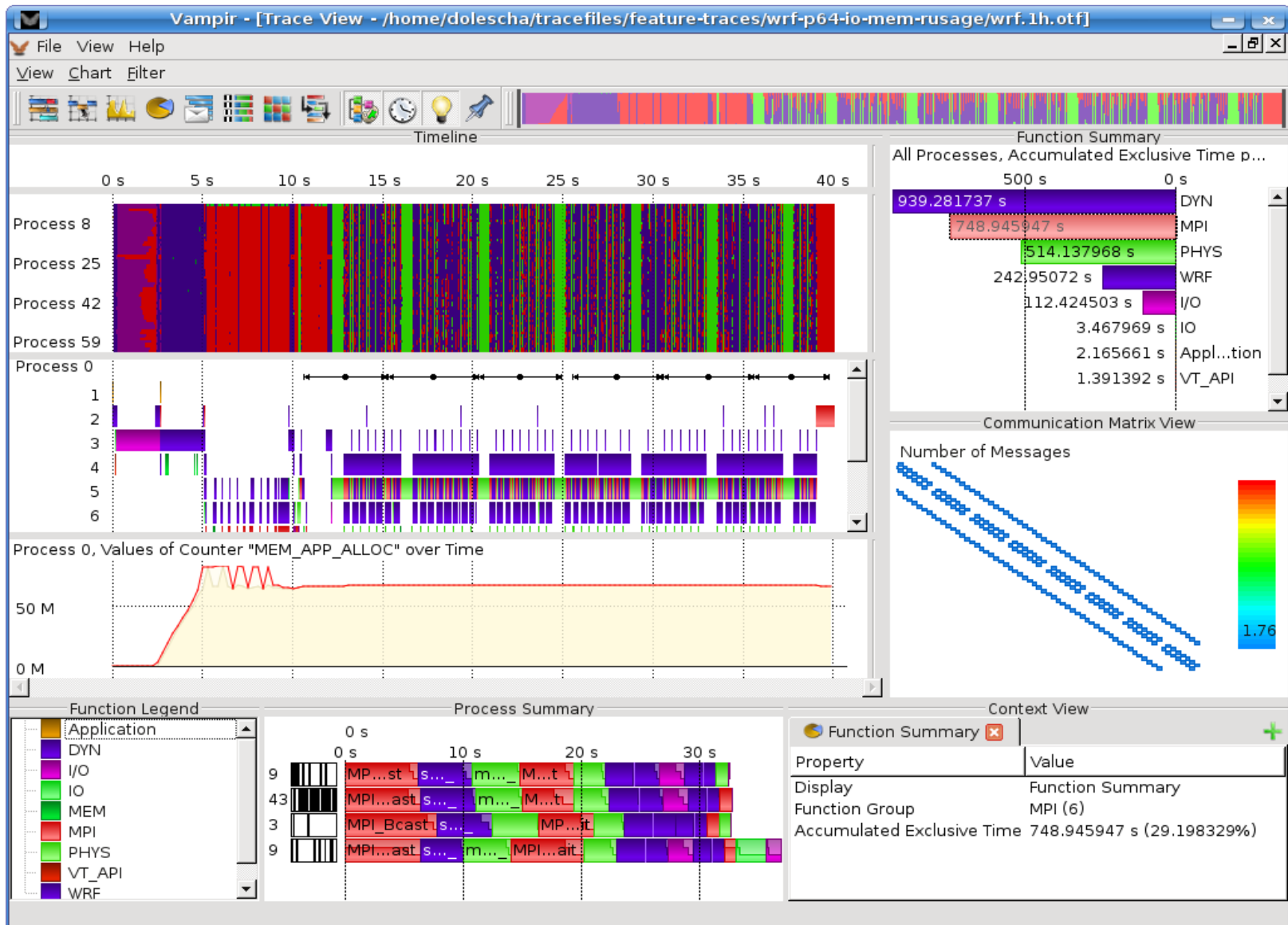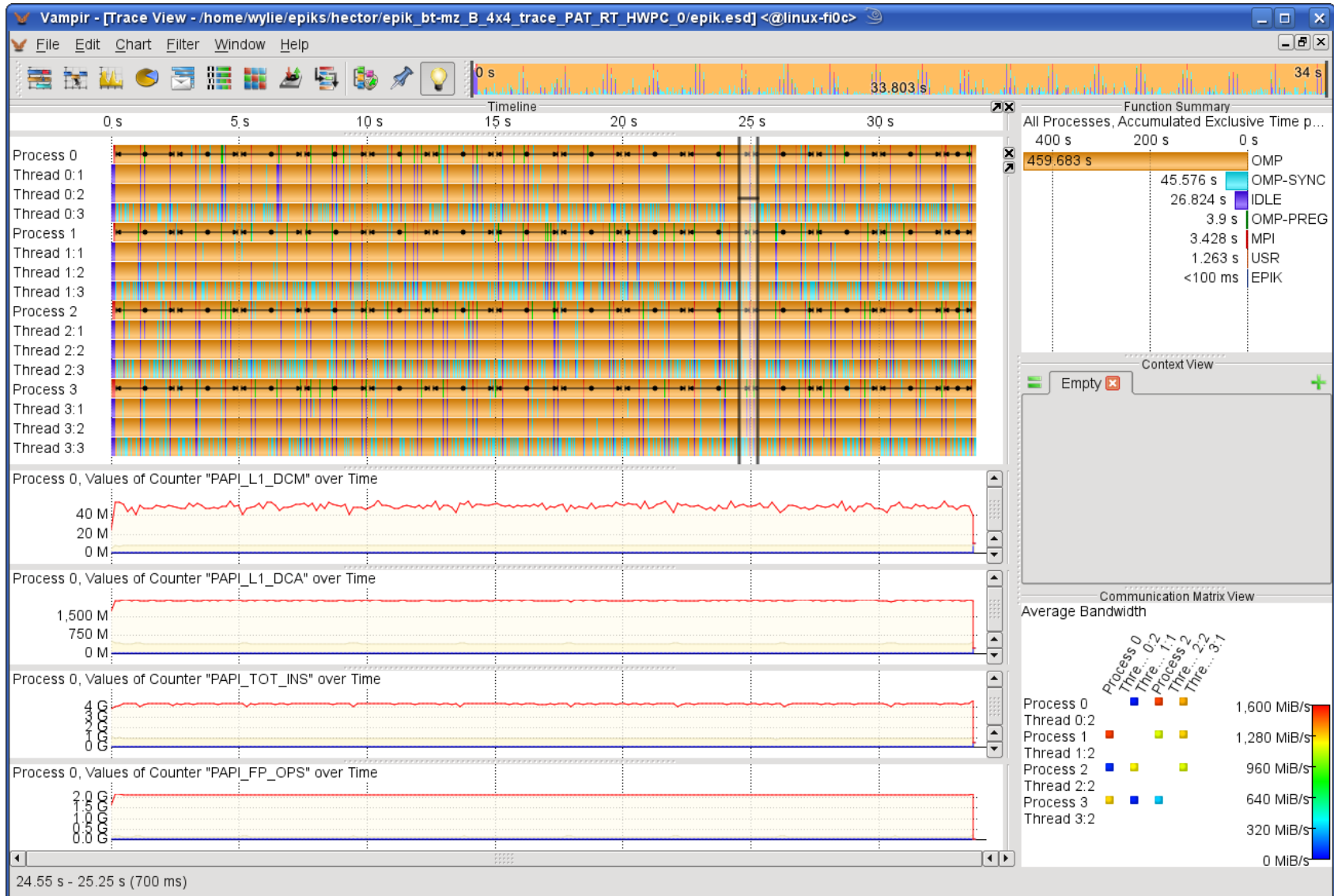    - event timeline chart for states & interactions of processes
    - Interactive browsing, zooming, selecting
  - Large variety of highly configurable analyses & displays
- Developed by Barcelona Supercomputing Center
  - Paraver trace analyser and Extrae measurement library
  - Open source available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

- Modular Assembler Quality Analyzer & Optimizer
  - Framework for binary manipulation
    - using plugins and scripting language
  - Tool exploiting framework to produce reports
    - fast prototyping and batch interface
  - STAN static performance model
  - MIL instrumentation language for dynamic analysis
    - building custom performance evaluation tools using HWCs
    - instrumentation of functions, loops, blocks & instructions

- Developed by UVSQ Exascale Computing Research lab
  - Supports Intel x86_64 microarchitecture
  - Available from www.maqao.org

- Key tool components also provided as open-source
  - Program development environment
    - Eclipse PTP ETFw, UNITE
  - Program/library instrumentation
    - COBI, OPARI, PDToolkit
  - Runtime measurement systems
    - PnMPI, Score-P, UniMCI
  - Scalable I/O
    - SIONlib
  - Libraries & tools for handling (and converting) traces
    - EPILOG, OTF, PEARL
  - Analysis algebra & hierarchical/topological presentation
    - CUBE

- Scalable performance measurement infrastructure
  - Supports instrumentation, profiling & trace collection, as well as online analysis of HPC parallel applications
  - Works with Periscope, Scalasca, TAU & Vampir prototypes
  - Based on updated tool components
    - CUBE4 profile data utilities & GUI
    - OA online access interface to performance measurements
    - OPARI2 OpenMP & pragma instrumenter
    - OTF2 open trace format

- Created by German BMBF SILC & US DOE PRIMA projects
  - JSC, RWTH, TUD, TUM, GNS, GRS, GWT & UO PRL
  - Available as BSD open-source from http://www.score-p.org/

- 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/

- Uniform integrated tool environment
  - Manages installation & access to program development tools
    - based on software environment management "modules"
    - commonly used on most cluster and HPC systems
    - configurable for multiple MPI libraries & compiler suites
  - Specifies how & where tools packages get installed
    - including integrating tools where possible
  - Defines standard module names and different versions
  - Supplies pre-defined module files
  - Configurable to co-exist with local installations & policies

- Developed by JSC, RWTH & TUD
  - Available as open-source from
    http://www.vi-hps.org/projects/unite/