# The ParMA Tools Package

**Bernd Mohr – Jülich Supercomputing Centre**

## ‖ ParMA

*Par*allel Programming for *M*ulti-core *A*rchitectures
**www.parma-itea2.org**

**ISC2010 ParMA BOF, Hamburg, May 31st 2010**

---

- **The ParMA Tools**

- Tool Integration

- The ParMA Tools Package

## Goals

- **Starting point of the ParMA project (2007)**
  - Established MPI debugging and performance analysis tools of German and UK partners
  - New innovative tools from French partners

- **Achieved results (2010)**
  - **Adaptation to + enhancements for multi-core systems**
    - Support for multi-threaded programming (e.g., OpenMP and POSIX threads)
  - General enhancements (e.g., MPI-2 and I/O support)
  - **Integration of tools to coherent environment**
  - Based on feedback from real-world, industrial applications

## Established Tools

**Universität Stuttgart**
- Open MPI / Peruse introspection interface
- **Marmot** MPI correctness checker

**Technische Universität Dresden / GWT**
- **Marmot** MPI correctness checker
- **VampirTrace** trace measurement system
- **Vampir** / **VampirServer** event trace visualizer

**Jülich Supercomputing Centre**
- **KOJAK** / **Scalasca** automatic trace-based performance analysis
- **CUBE** visualization browser

**Allinea (unfunded partner)**
- **DDT** parallel debugger

## New Innovative Tools

**UVSQ (France)**

MAQAO

- **MAQAO** code and memory analysis on x86_64 architecture
- **DECAN** decremental analysis of machine code instructions
- **MTL** memory tracing library
- **MADRAS** binary instrumentation

**It Sudparis (France)**

- **STEP** can automatically transform OpenMP into MPI (and hybrid MPI/OpenMP) programs
- Allows to move shared-memory applications (limited by the power of one node) on to a cluster

## DDT Parallel Debugger
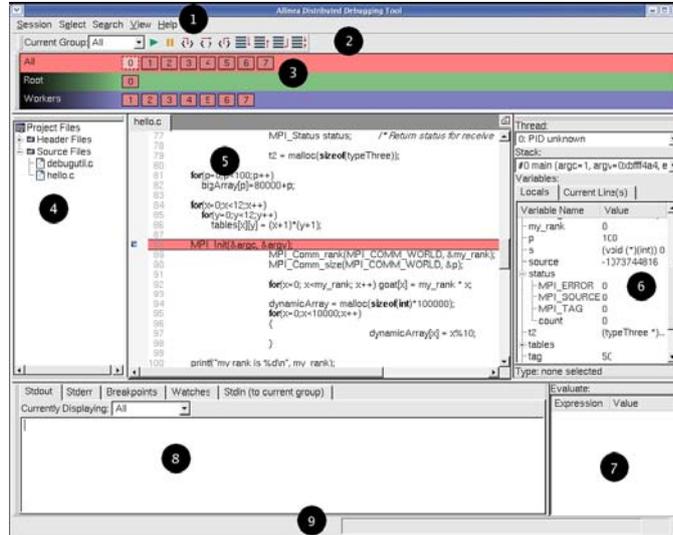
- **DDT**

  allinea
  SCALE TO NEW HEIGHTS

  – Graphical parallel debugger

- Commercial product of Allinea
- Portable across many platforms
- Support for parallel programs
  – Multi-tasking (MPI)
  – Multi-threading (OpenMP, pthreads)
  – GPUs
  – At scale

## Example: DDT Main Screen

❶ Menu bar
❷ Process controls
❸ Process group window
❹ File window
❺ Code window
❻ Variable window
❼ Evaluate window
❽ Output window
❾ Status bar

---

## MARMOT MPI Correctness Checker

- **MARMOT** checks
  - Correct usage of MPI calls
  - Portability
  - MPI resource usage

- Open source
- Developed in cooperation of ZIH and HLRS
- Supports
  - C, C++ and Fortran
  - MPI, or hybrid OpenMP/MPI

# Example: MARMOT Analysis

On Call: MPI_Recv for MPI-Standard information see:/usr/local/marmot/marmot_r827_openmpi-1.2.5-dbg/share/doc/marmot-2.0.0

144: Error from rank 2(Thread: 0) with Text: ERROR: MPI_Recv: At least a part of the specified buffer is still in use!

=== Detailed information ===
This buffer is specified by:
Starting address: 140734999716208
Count: 1
Extent of used datatype: 12
Resulting end address (first non used byte): 140734999716220
Other buffer is specified by:
Starting address: 140734999716212
Count: 1
Extent of used datatype: 12
Resulting end address (first non used byte): 140734999716224

On Call: MPI_Recv for MPI-Standard information see:/usr/local/marmot/marmot_r827_openmpi-1.2.5-dbg/share/doc/marmot-2.0.0

153: Error from rank 0(Thread: 0) with Text: ERROR: MPI_Recv: At least a part of the specified buffer is still in use!

=== Detailed information ===
This buffer is specified by:
Starting address: 140737469793568
Count: 1
Extent of used datatype: 12
Resulting end address (first non used byte): 140737469793580
Other buffer is specified by:
Starting address: 140737469793572
Count: 1
Extent of used datatype: 12
Resulting end address (first non used byte): 140737469793584

- Plain text

- HTML

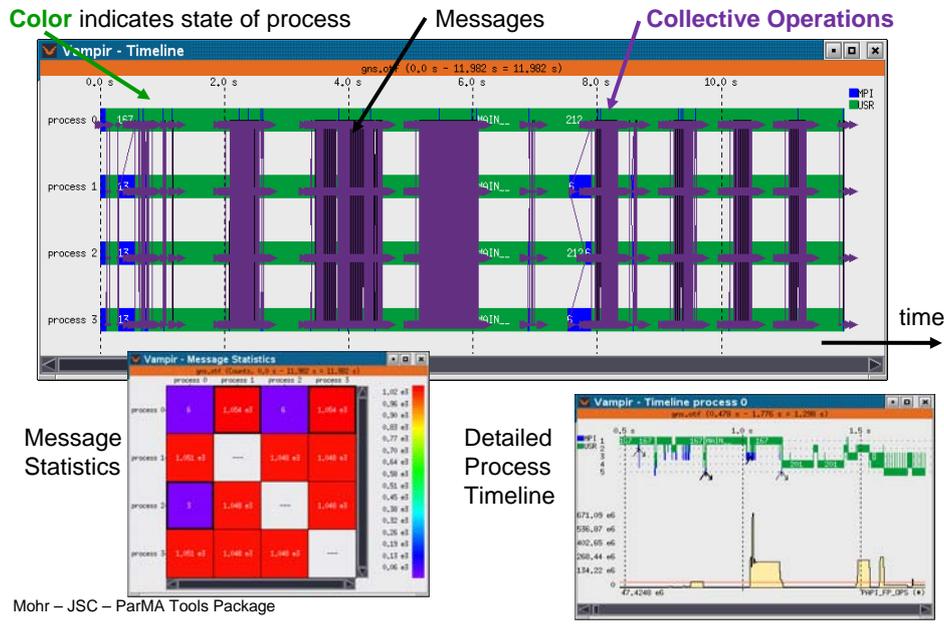| 140 | 0 | 0 | Warning | Text: WARNING: MPI_Issend: Count = 0! Call: MPI_Issend |
| 142 | 0 | 0 | Warning | Text: WARNING: MPI_Issend: datatype is for reduction functions (C)! Call: MPI_Issend |
| 144 | 3 | 0 | Error | Text: ERROR: MPI_Recv: At least a part of the specified buffer is still in use! The buffer is still used by a call to MPI_Issend. === Detailed information === This buffer is specified by: Starting address: 140737197187344 Count: 1 Extent of used datatype: 12 Resulting end address (first non used byte): 140737197187356 Other buffer is specified by: Starting address: 140737197187348 Count: 1 Extent of used datatype: 12 Resulting end address (first non used byte): 140737197187360 Call: MPI_Recv |
| | | | | Text: ERROR: MPI_Recv: At least a part of the specified buffer is still in use! The buffer is still used by a call to MPI_Issend. |

---

# Vampir Tool Environment

- **VampirTrace**
  - Open source
  - Event trace measurement system
  - Instruments C, C++, and Fortran
    - MPI, (simple) OpenMP, or hybrid, I/O
  - Collects event traces in OTF format
- Post-mortem visual analysis
  - Developed originally by Jülich and since 1997 by ZHR/ZIH of TU Dresden
  - Commercial distribution by GWT-TUD
  - **Vampir**: sequential visualizer
  - **VampirServer**: distributed client / parallel server architecture

VAMPIR

ZIH
Center for Information Services & High Performance Computing

GWT forschung+innovation

# Example: Vampir Display of GNS code

**Color** indicates state of process     Messages     **Collective Operations**



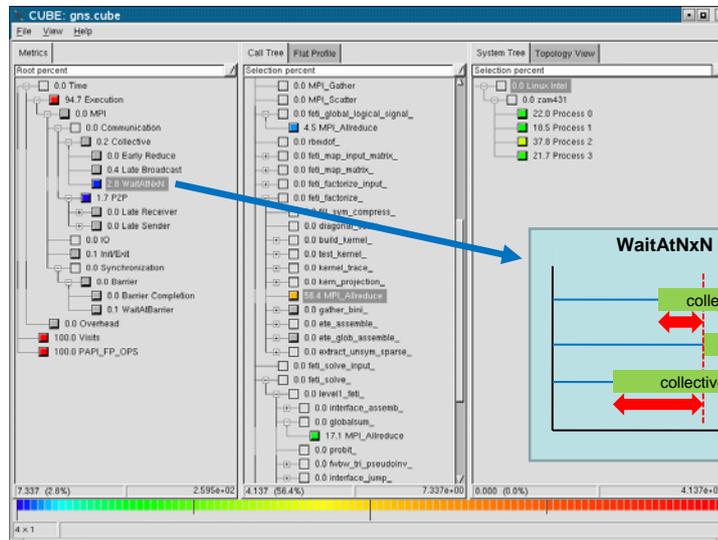time

Message Statistics

Detailed Process Timeline

---

# KOJAK / Scalasca Toolset

- Open source
- Jülich Supercomputing Centre

- Automatic performance analysis
  - Instrument C, C++, and Fortran parallel applications
    - Based on MPI, (simple) OpenMP, or hybrid
  - Collect call-path-profiles and EPILOG event traces
  - Scan trace for event patterns representing inefficiencies
    - **KOJAK**:     sequential analysis
    - **Scalasca**:     parallel analysis
  - Categorize and rank inefficiencies found
  - Visualize via **CUBE** browser

Example: KOJAK Analysis of GNS code

Identified Problems    Where in source?    Which processes are affected?

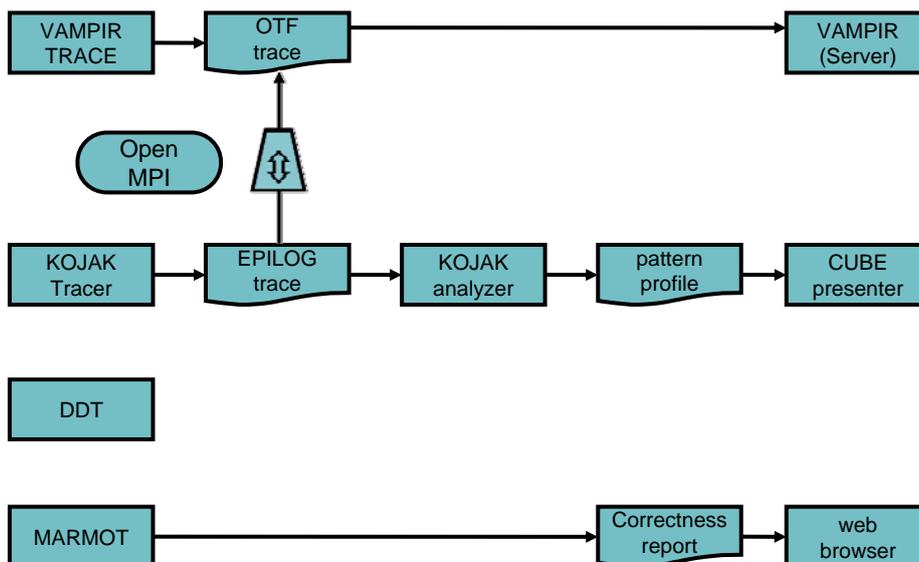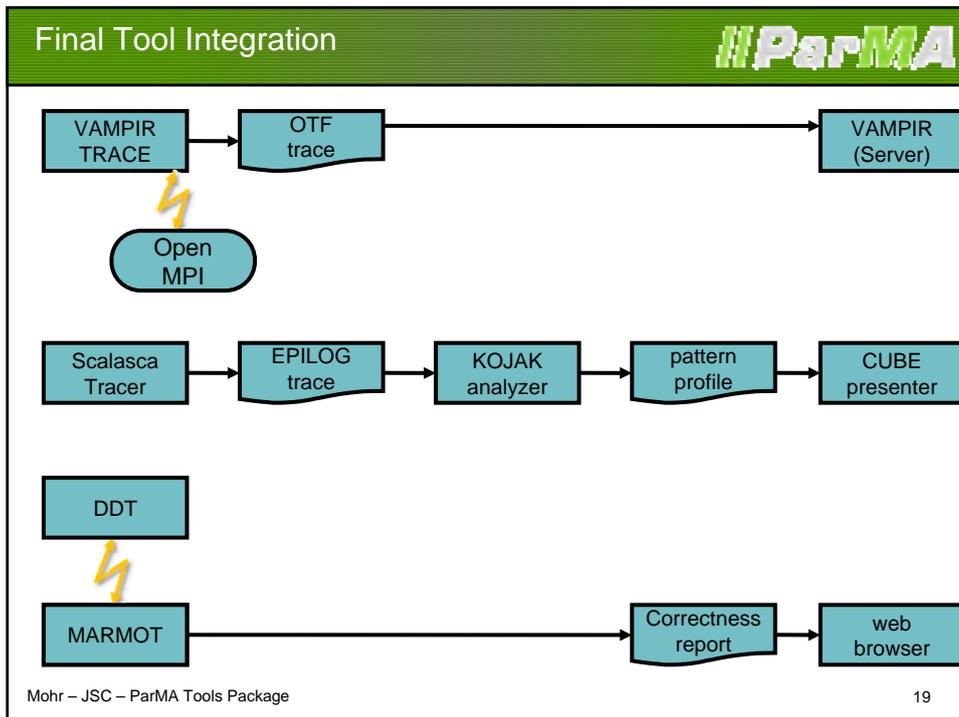WaitAtNxN Pattern

---

# MAQAO Toolset

- **MAQAO (Modular Assembly Quality Analyzer and Optimizer)**
  - Combines static analysis of assembly code with dynamic analysis of execution traces
  - Focus on loop-level
- **MADRAS (Multi Architecture Disassembler ReAssembler)**
  - Instrumentation of binary files for MAQAO
- **MTL (MAQAO Trace Library)**
  - Memory trace library integrated in MAQAO, for threaded codes
  - Detect important inefficiencies (e.g., false sharing, strided access) or opportunities for optimizations (setting thread affinity according to reuse among threads)
- **DECAN (Decremental Analysis)**
  - Automatic detection of performance anomalies (e.g., inefficient memory access) via iterative modification of machine code instructions of hot functions

- Université de Versailles St-Quentin-en-Yvelines

---

- DDT (Allinea) and Vampir (GWT-TUD)
  - Well established European HPC software vendors
- Marmot (HLRS / ZIH)
  - Competing tools either
    - vendor-specific only (NEC)
    - or of limited functionality (MPICH checker)
- Vampir (ZIH / GWT-TUD)
  - Competing tools either
    - vendor-specific only (Intel Trace Analyzer)
    - or less portable / un-supported (BSC Paraver)
- KOJAK/Scalasca (JSC)
  - World-leading product (no competitors)
- **Tool integrations** (DDT/Marmot, Vampir/Marmot, Scalasca/Vampir, …)
  - **UNIQUE!**

## Content

- The ParMA Tools

- **Tool Integration**

- The ParMA Tools Package

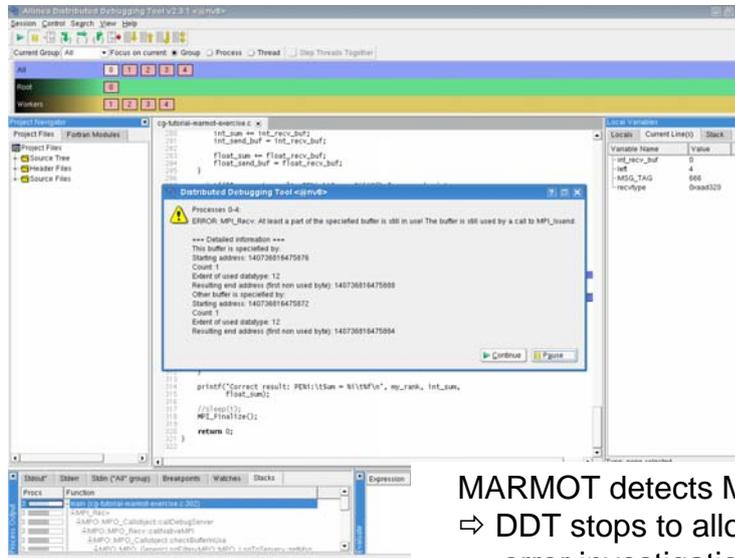## (NO) Integration at Start of ParMA

## Final Tool Integration

---

## Marmot  ⇨  DDT Integration I



- **Integration of MARMOT in DDT as plugin**

- MARMOT can be activated at start of a new DDT session

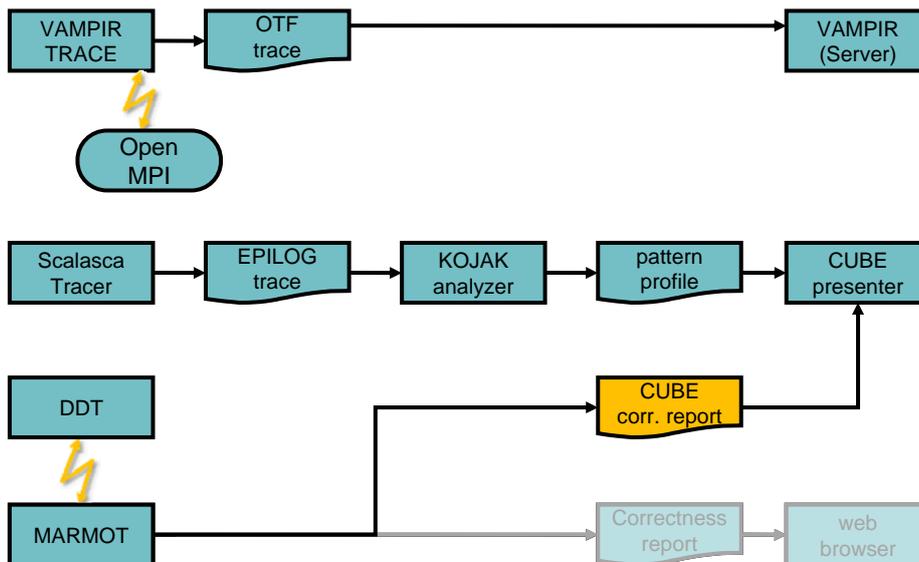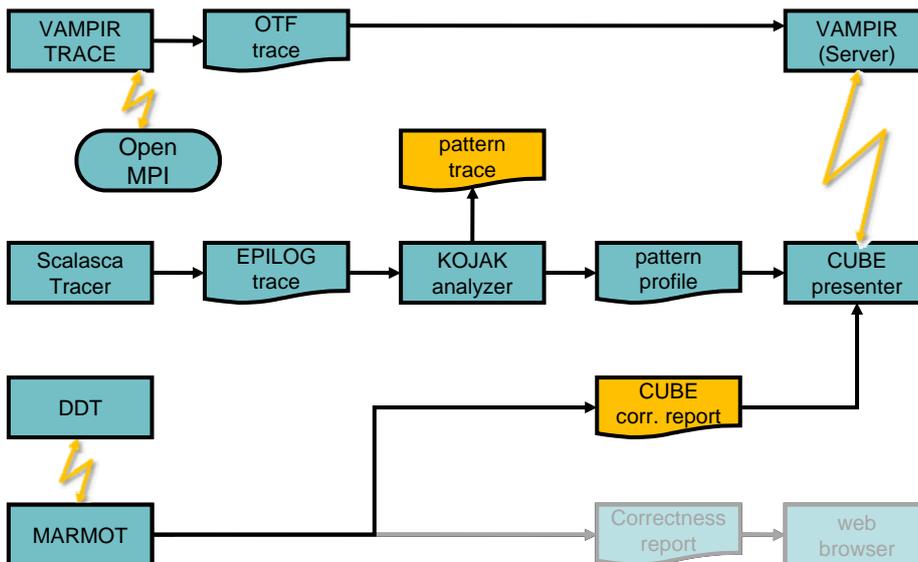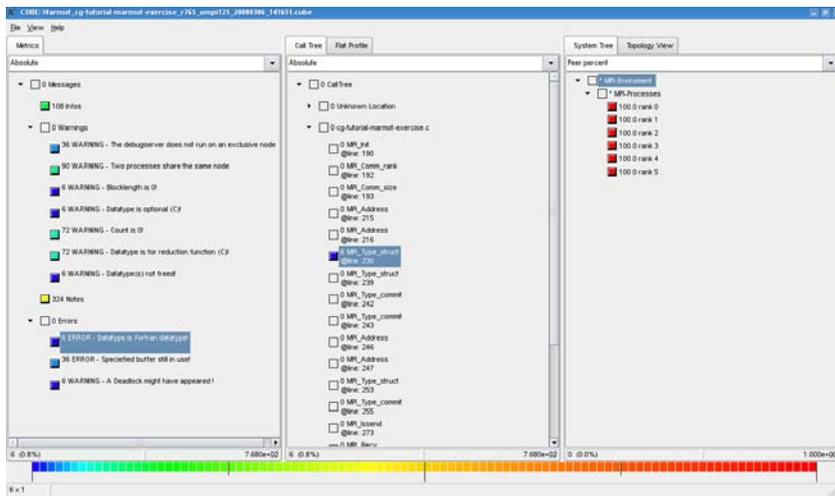## Marmot ⇨ DDT Integration II



MARMOT detects MPI issue
⇨ DDT stops to allow further
error investigation

## Final Tool Integration

Marmot ⇨ CUBE

- **MARMOT result analysis using CUBE browser**

Final Tool Integration

VAMPIR ⇨ KOJAK via Pattern Traces

**||ParMA**

Original Vampir event trace

Pattern trace generated by KOJAK analysis highlighting problematic areas
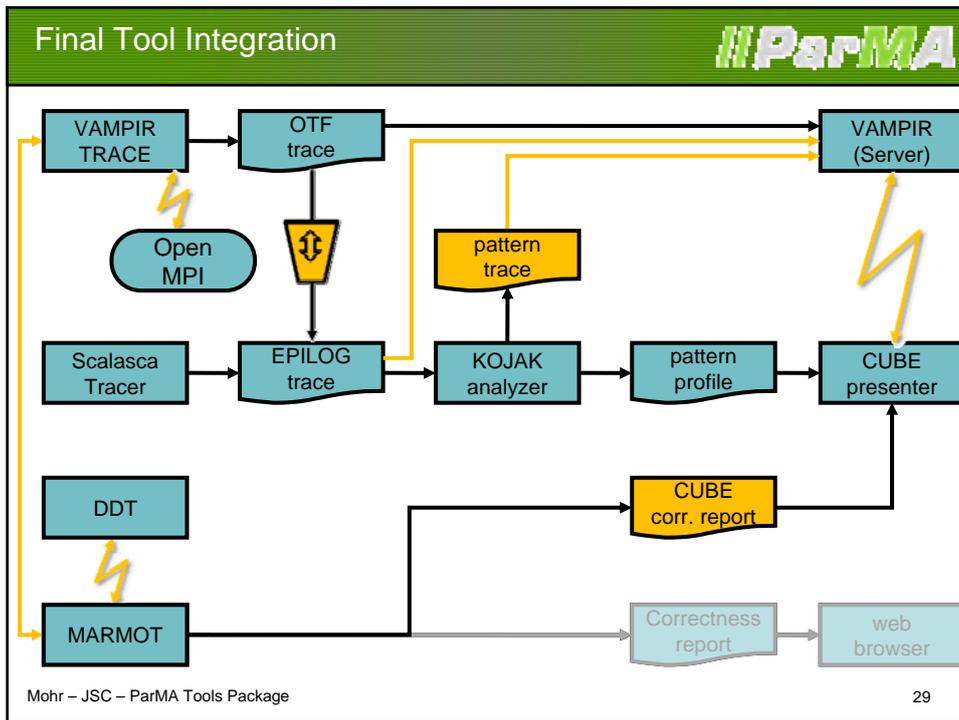
Mohr – JSC – ParMA Tools Package

25



KOJAK ⇨ Vampir Integration

**||ParMA**

❶ **Connect to trace browser**

❷ **Max severity in trace browser**

Mohr – JSC – ParMA Tools Package

## Final Tool Integration (D3.all.5)

## Vampir ⇔ KOJAK/Scalasca Integration

- VampirServer can read
  - KOJAK sequential EPILOG traces
  - Scalasca parallel EPILOG traces

  in addition to native OTF traces

- New OTF-to-EPILOG converter
  - Also: beta version of KOJAK's event pattern trace analyzer based on new Œ interface (common OTF/EPILOG reading interface)

- **New** BMBF SILC **project**:
  - **full integration** of **Scalasca** and **VampirTrace** measurement systems and trace formats

Final Tool Integration — ParMA

VAMPIR TRACE → OTF trace → VAMPIR (Server)

Open MPI

pattern trace

Scalasca Tracer → EPILOG trace → KOJAK analyzer → pattern profile → CUBE presenter

DDT

CUBE corr. report

MARMOT → Correctness report → web browser

Marmot ⇨ Vampir Integration — ParMA

- **Uni**versal **M**PI **C**orrectness **I**nterface for integration of MPI correctness tools into other tools (UniMCI)
- Implemented by Marmot (Producer) and VampirTrace (Consumer)

**∥ParMA**

- The ParMA Tools

- Tool Integration

- **The ParMA Tools Package**

---

**∥ParMA**

- **UN**ified **I**ntegrated **T**ool **E**nvironment
- UNITE website: **http://apps.fz-juelich.de/unite/**
- Lower bar for inexperienced users and admins
  - **Common** usage and installation **documentation**
  - **Download, build and install**
    all ParMA tools **in one package:**

| | |
|---|---|
| – UNITE | – Scalasca-1.3.1 |
|    package installer and | – Vampirtrace-5.8.2 |
|    module package | – UniMCI-1.0.1 |
| – OTF-1.6.5 | – Marmot-2.4 |
| – pdtoolkit-3.15 | – Vampir-5.x or 7.x **(*)** |
| – cube-3.3 | – VampirServer-1.x, 2.x **(*)** |

**(*)** Valid icense required   

# ParMA Tools Package II

- Extensively tested on
  - Itanium/IA32/x86_64 platforms with various MPI libraries (MPICH1, MPICH2, OpenMPI, Intel MPI, LAM, BullMPI, Parastation MPI, SGI MPT, ...)
  - AIX and Solaris clusters

- Already in use on Bull Nova and production machines of JSC, ZIH, RWTH, HLRN, (and soon LRZ, HLRS, …)

- Future work (beyond ParMA):
  - Integration of rest of ParMA tools (DDT, MAQAO, STEP)
  - Integration of other tools (Paraver, TAU, …)
  - More platforms (Cray XT, IBM BlueGene, NEC)