# pallas

## Competence in High Performance Computing

## Portable MPI Tools at Work
## - Cracking Performance Problems

Werner Krotz-Vogel, Pallas GmbH

<krotz@pallas.com>

42nd CUG Conference, May 02000

Pallas GmbH
Hermülheimer Straße 10
D-50321 Brühl, Germany

info@pallas.com
www.pallas.com

# Pallas Tools Software

State-of-the-art program development tools ...

... in detail:

- Vampir-2.5 (online-Demo), Vampirtrace-2.0, Dimemas

... briefly:

- Etnus TotalView 4.0 Multi-process Debugger
- PGI 3.1 x86 Compilers, Cluster Development Kit (CDK)
- KAP/Pro Toolset 3.8, OpenMP
- KAI C++ 3.4, ISO standard
- FORESYS - Fortran Restructuring Tool

... free open source:

- PMB - Pallas MPI Benchmark Suite (incl. "effective Bandwith")

# PGI Cluster Development Kit (CDK)

Compilers & Tools ...

- PGI 3.1 x86 compilers , C, C++, F77, F90, HPF, pgrof, pgdbg
- SMP/OpenMP support for C, C++, F77, F90

... plus convenient add-on´s:

- parallel ScaLAPACK
- optimized BLAS, LAPACK
- MPI/mpich
- PVM
- PBS - Portable Batch System
- Tutorial, examples
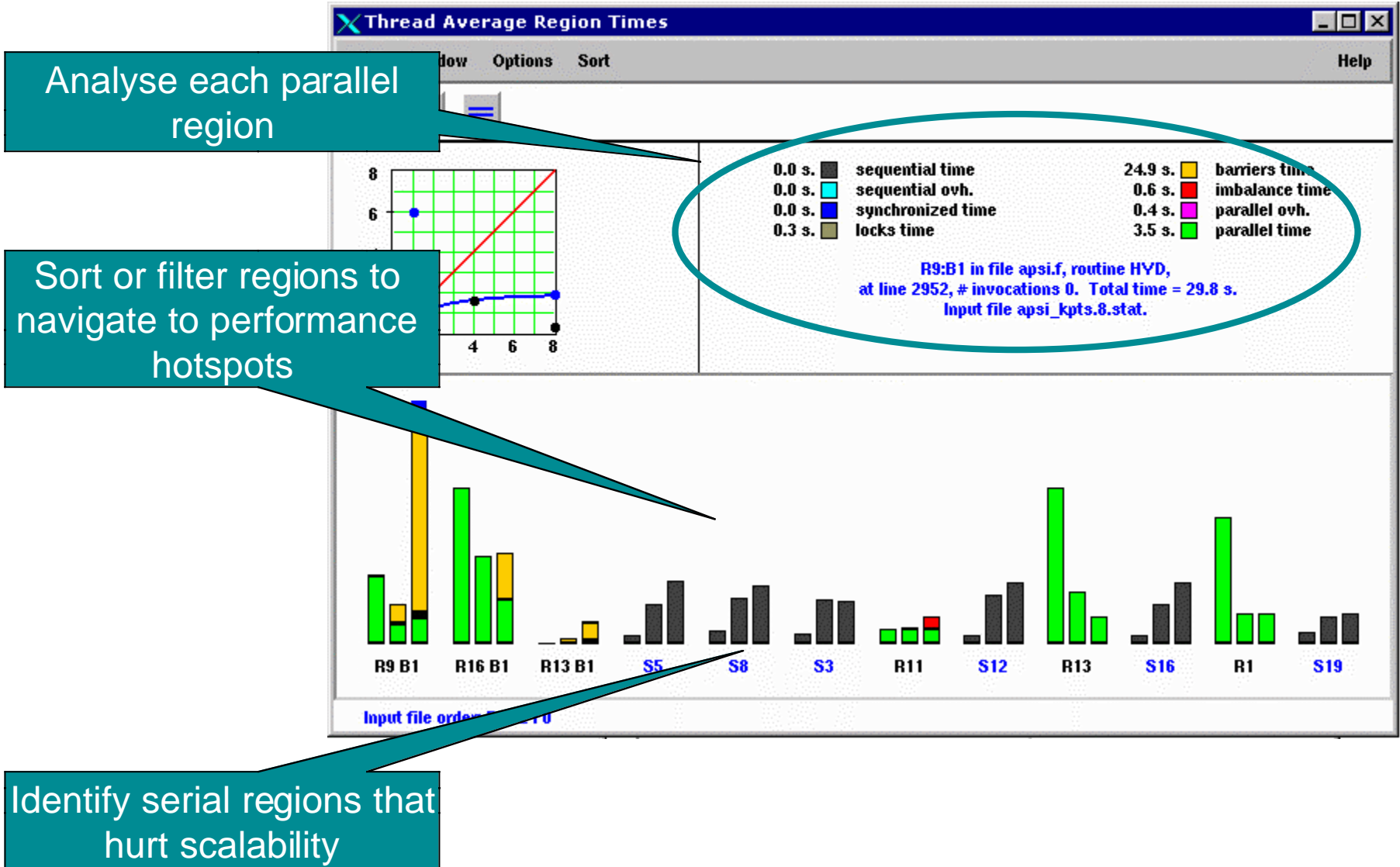- Cluster management utilities

# KAP/Pro Toolset - Assure Example



Source location

# KAP/Pro Toolset - GuideView Example



**Analyse each parallel region**

**Sort or filter regions to navigate to performance hotspots**

**Identify serial regions that hurt scalability**

Thread Average Region Times

Options  Sort  Help

| 0.0 s. | sequential time | 24.9 s. | barriers time |
| 0.0 s. | sequential ovh. | 0.6 s. | imbalance time |
| 0.0 s. | synchronized time | 0.4 s. | parallel ovh. |
| 0.3 s. | locks time | 3.5 s. | parallel time |

R9:B1 in file apsi.f, routine HYD,
at line 2952, # invocations 0.  Total time = 29.8 s.
Input file apsi_kpts.8.stat.

R9 B1   R16 B1   R13 B1   S5   S8   S3   R11   S12   R13   S16   R1   S19

Input file order

# KAI C++

The most modern, best performing, platform independant C++

- ISO C++ standard syntax, including exeptions and member templates

- ISO C++ standard class library

- multi-platform support

- meet C performance requirements

- thread safety (on most platforms)

# FORESYS

- Translates FORTRAN code ( F77 - F95 ) into abstract syntax tree (ForLib)

- FORTRAN code consistency checks (definitions of functions, common blocks etc.)

- Interactive visualization & analysis of inconsistencies

- Upgrading from FORTRAN 77 to FORTRAN 90

- Interactive/batch analysis of parallelization possibilities

- Automatic code quality analysis/improvements

# pallas

# Vampir 2.5

**V**isualization and
**A**nalysis of
**MPI**
**P**rograms

# Vampir

- New version: Vampir 2.5

- Significant new features
  - support for collective MPI operations
  - trace comparison
  - tracefile re–write
  - message–length histogram
  - local and global calling trees
  - source–code reference
  - support for MPI–2 I/O operations

# Vampir Features

- Offline trace analysis for MPI (and others ...)
- Traces generated by Vampirtrace tool (`ld ... -IVT -lpmpi -lmpi`)
- Convenient user–interface
- Scalability in time and processor–space
- Excellent zooming and filtering
- High–performance graphics
- Display and analysis of MPI and application events:
    - execution of MPI routines
    - point–to–point and collective communication
    - MPI–2 I/O operations
    - execution of application subroutines (optional)

- Easy customization

Vampir 2.5 main window



- **Tracefile loading** can be interrupted at any time
- **Tracefile loading** can be resumed
- **Tracefile** can be loaded starting at a specified time offset

- **Tracefile** can be re–written (re–grouped symbols)

# Vampir Displays

- **Global displays** show all selected processes
  - **Summary Chart:** aggregated profiling information
  - **Activity Chart:** presents per–process profiling information
  - **Timeline:** detailed application execution over time axis
  - **Communication statistics:** message statistics for each process pair
  - **Global Comm. Statistics:** collective operations statistics
  - **I/O Statistics:** MPI I/O operation statistics
  - **Calling Tree:** draws global or local dynamic calling trees

- **Process displays** show a single process per window
  - Activity Chart
  - Timeline
  - Calling Tree

# Summary Chart



VAMPIR – Summary Chart
lu-T3E-A.64.bpv: Summary Chart (Times, 47.775 ms–4.464 s)

| | |
|---|---|
| MPI | |
| Calculation | |
| Setup | 98.853 ms |
| Communication | 75.454 ms |
| Application | 10.662 ms |
| Verify | 9.344 ms |
| Tracing | 0.162 ms |
| VT_API | 3.572 us |

0.5 s     1.0 s

VAMPIR – Summary Chart
lu-T3E-A.64.bpv: Summary Chart (Times, 47.775 ms–4.464 s)

| | |
|---|---|
| MPI_Recv | 0.813 s |
| MPI_Bcast | 0.155 s |
| MPI_Wait | 0.115 s |
| MPI_Send | 74.611 ms |
| MPI_Allreduce | 29.973 ms |
| MPI_Irecv | 3.503 ms |
| MPI_Barrier | 48.077 us |
| MPI_Comm_rank | 7.08 us |
| MPI_Comm_size | 4.161 us |

0.2 s     0.4 s     0.6 s     0.8 s

MPI

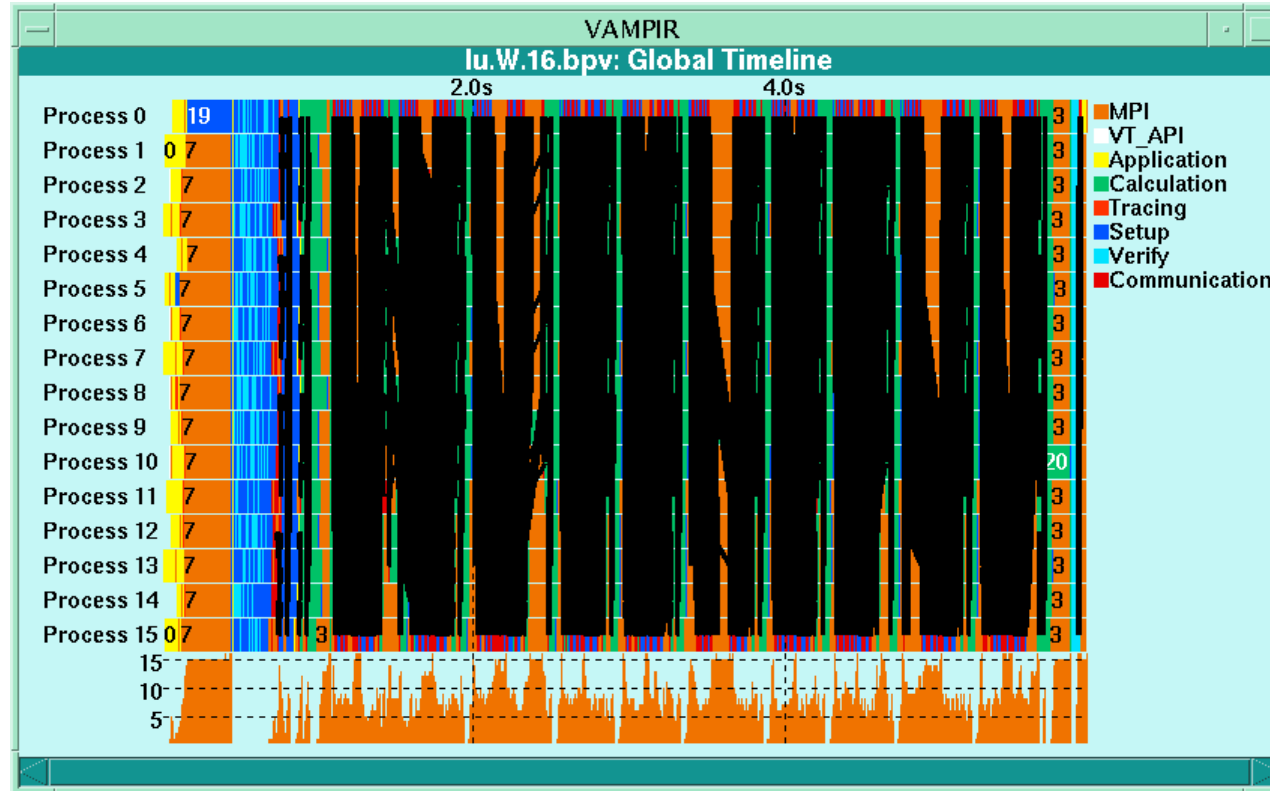- Aggregated profiling information
    - execution time
    - number of calls
- Inclusive or exclusive of called routines
- Look at all/any category or all states
- Values can be exported/imported
- Tracefiles can be compared1

© Pallas GmbH

# Timeline Display



- Now displays MPI collective and I/O operations
- To zoom, draw rectangle with the mouse
- Also used to select sub–intervals for statistics

# Timeline Display (Message Info)



See message details

Click on message line

Message send op

Message receive op

**Identified Message**

Message sent from Process 7 to Process 6
communicator: 0, type: 1
length: 21120
sent at 898.439 ms, received at 907.921 ms
Data rate: 2.227 MBytes/sec

Close

VAMPIR
lu.W.16.bpv: Global Timeline
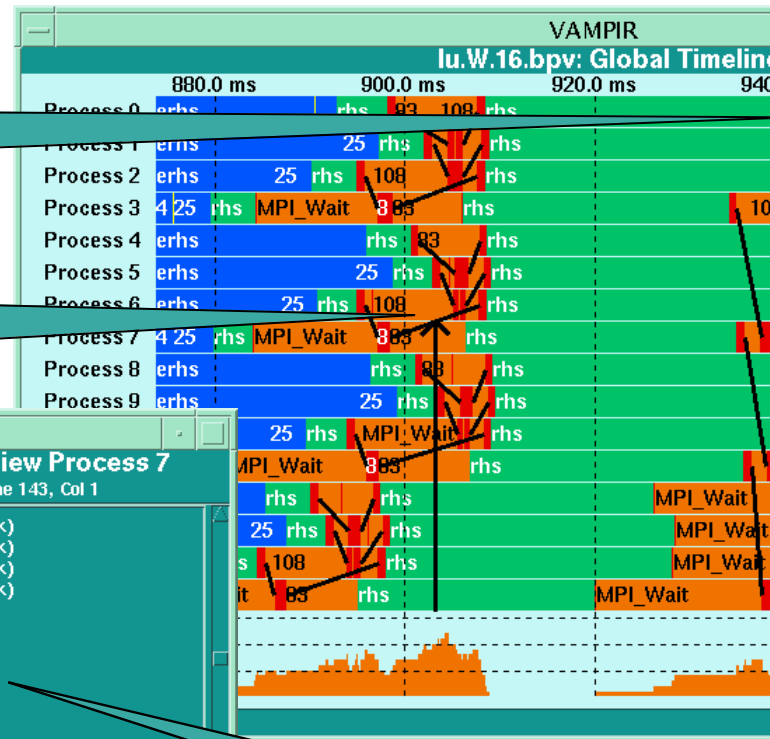880.0 ms    900.0 ms    920.0 ms    940.

VAMPIR
lu.W.16.bpv: Source View Process 7
/tmp/traces/NPB–LU/exchange_3.f: Line 143, Col 1

```
136        buf(2,ipos2) = g(2,1,j,k)
137        buf(3,ipos2) = g(3,1,j,k)
138        buf(4,ipos2) = g(4,1,j,k)
139        buf(5,ipos2) = g(5,1,j,k)
140      end do
141      end do
142
143      call MPI_SEND( buf,
144    >        10*ny*nz,
145    >        dp_type,
146    >        north,
147    >        from_s,
148    >        MPI_COMM_WORLD,
149    >        IERROR )
150      end if
```

VAMPIR
lu.W.16.bpv: Source View Process 6
/tmp/traces/NPB–LU/exchange_3.f: Line 156, Col 1

```
149    >        IERROR )
150      end if
151
152  c----------------------------------
153  c  receive from south
154  c----------------------------------
155      if (south.ne.–1) then
156        call MPI_WAIT( mid, STATUS, IERROR )
157
158      do k = 1,nz
159      do j = 1,ny
160        ipos1 = (k–1)*ny + j
161        ipos2 = ipos1 + ny*nz
162        x+2,j,k) = buf1(1,ipos1)
163        2,j,k) = buf1(2,ipos1)
```

# Communication Statistics



- Message statistics for each process pair:
  - Byte and message count
  - min/max/avg message length
  - min/max/avg bandwidth
- Filter for message tags or communicators

# Message Histograms



- Message statistics by length, tag or communicator
  - Byte and message count
  - min/max/avg bandwidth
- Filter for message tags or communicators

# Collective Operations

- For each process: mark operation locally

**MPI_Gather**

Stop of op

Start of op

Data being sent

Data being received

- Connect start/stop points by lines

**MPI_Gather** | User_Code

User_Code

User_Code

Connection lines

# Collective Operations

See global timing info

**VAMPIR – Identified Global Operation**

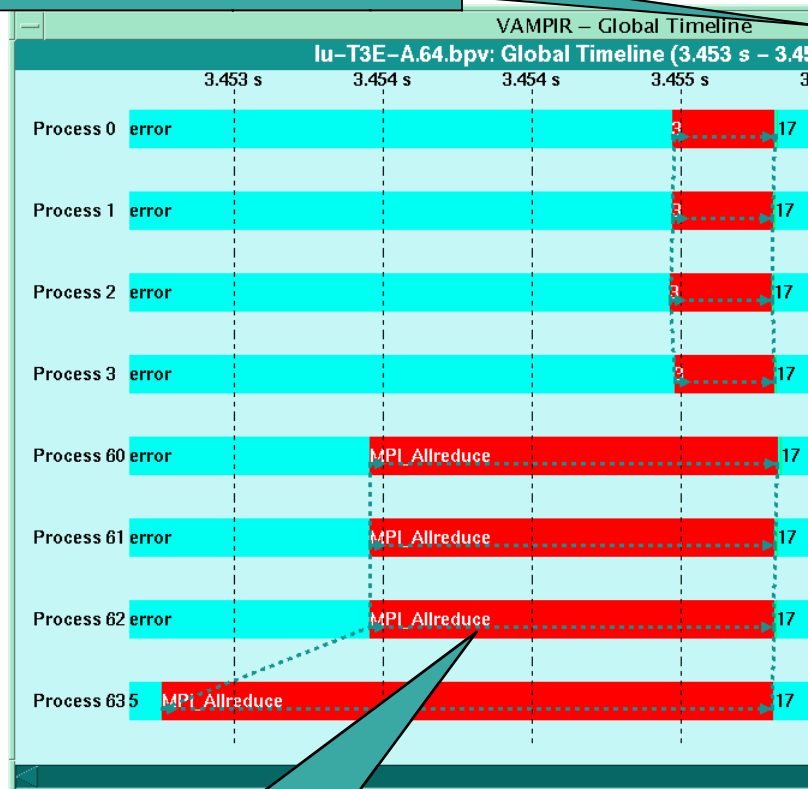| | |
|---|---|
| **Root:** | Process 0 |
| **Participants:** | Process(s) 0–63 |
| **Operation:** | MPI_Allreduce |
| **Communicator:** | 0 |
| **Interval:** | 3.453258 s – 3.455321 s |
| **Duration:** | 2.06312 ms |
| **Length:** | 2560 bytes / 2560 bytes |
| **Send rate:** | 1.183 Mbytes/s |

**Local Values**   **Close**

**VAMPIR – Global Timeline**

**lu–T3E–A.64.bpv: Global Timeline (3.453 s – 3.4...**

| | 3.453 s | 3.454 s | 3.454 s | 3.455 s | 3 |
|---|---|---|---|---|---|
| Process 0 error | | | | 3 | 17 |
| Process 1 error | | | | 3 | 17 |
| Process 2 error | | | | 3 | 17 |
| Process 3 error | | | | 3 | 17 |
| Process 60 error | | MPI_Allreduce | | | 17 |
| Process 61 error | | MPI_Allreduce | | | 17 |
| Process 62 error | | MPI_Allreduce | | | 17 |
| Process 63 5 | MPI_Allreduce | | | | 17 |

MPI_Wait 108

**VAMPIR – Identified Global Operation**

| | |
|---|---|
| **Root:** | Process 0 |
| **Location:** | Process 62 |
| **Operation:** | MPI_Allreduce |
| **Communicator:** | 0 |
| **Interval:** | 3.453956 s – 3.455311 s |
| **Duration:** | 1.35444 ms |
| **Length:** | 40 bytes / 40 bytes |
| **Send rate:** | 28.84 Kbytes/s |

**Global Values**   **Close**

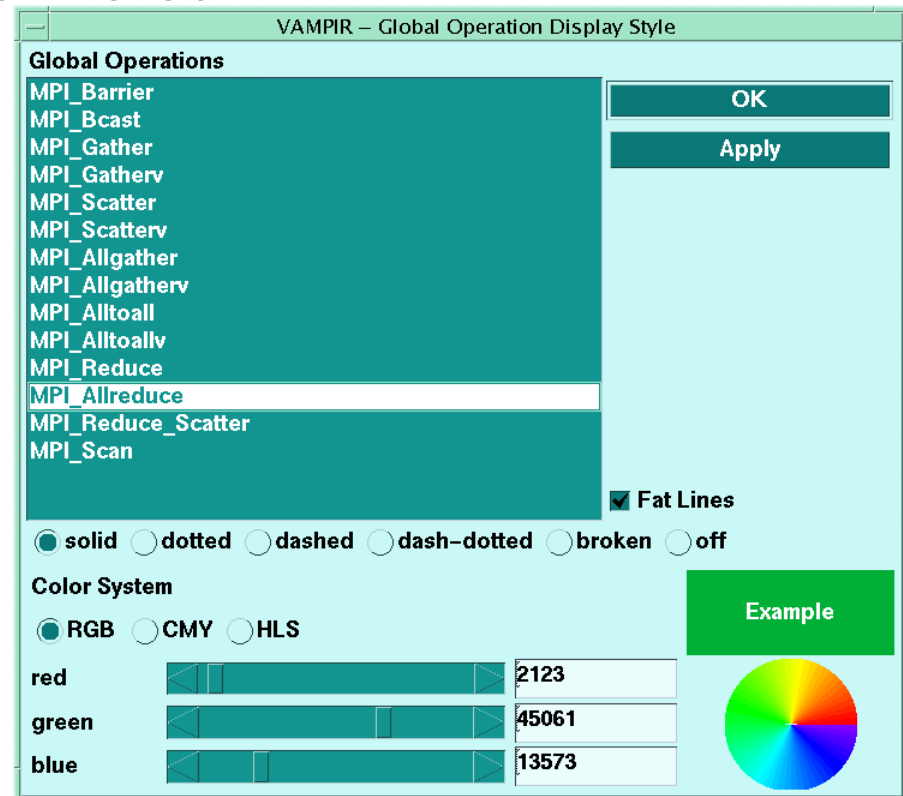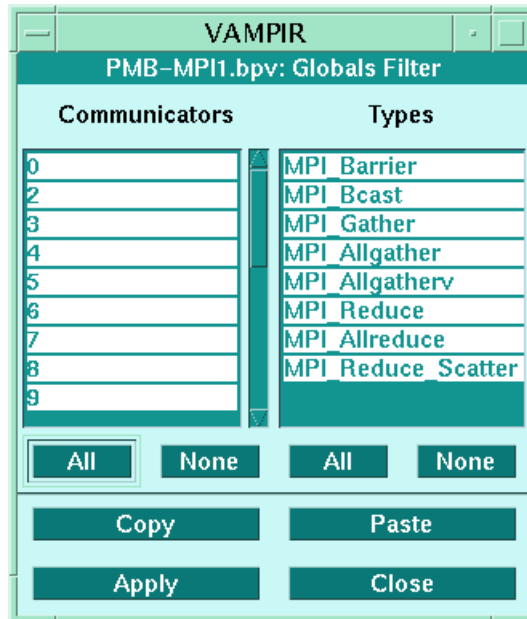MPI_Allreduce

Click on collective
operation display

See local timing info

# Collective Operations

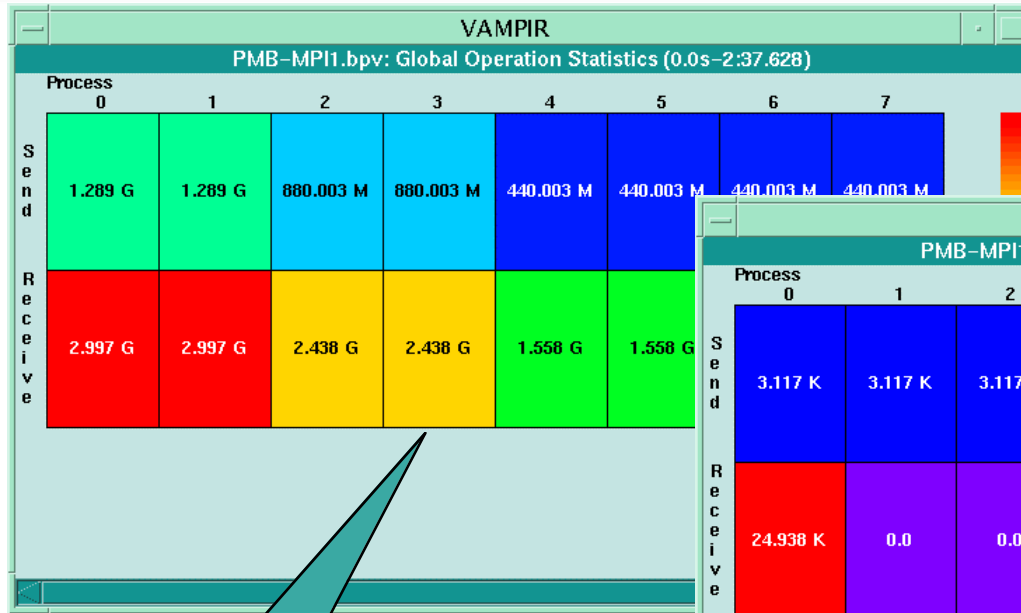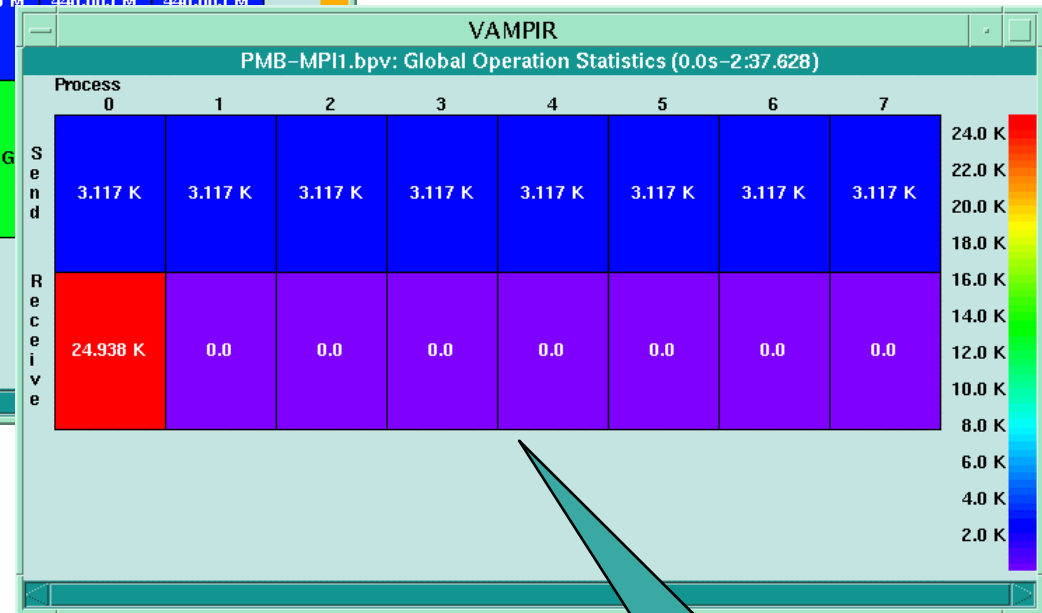- **Collective operations can be filtered**



- **The display style can be**
  **adapted for each collective operation**

# Global Communication Statistics



All collective operations

MPI_Gather only

- ■ Statistics for collective operations:
  - – operation counts, Bytes sent/received
  - – transmission rates
- ■ Filter for collective operation

# Vampirtrace

Tracing of
MPI and
Application
Events

# Vampirtrace

- New version: Vampirtrace 2.0

- Significant new features:
  - records collective communication
  - enhanced filter functions
  - extended API
  - records source–code information (selected platforms)
  - support for shmem (Cray T3E)
  - records MPI–2 I/O operations

- Available for all major MPI platforms

# MPI–I/O Operations

See detailed I/O information

**VAMPIR – Identified File I/O Event**

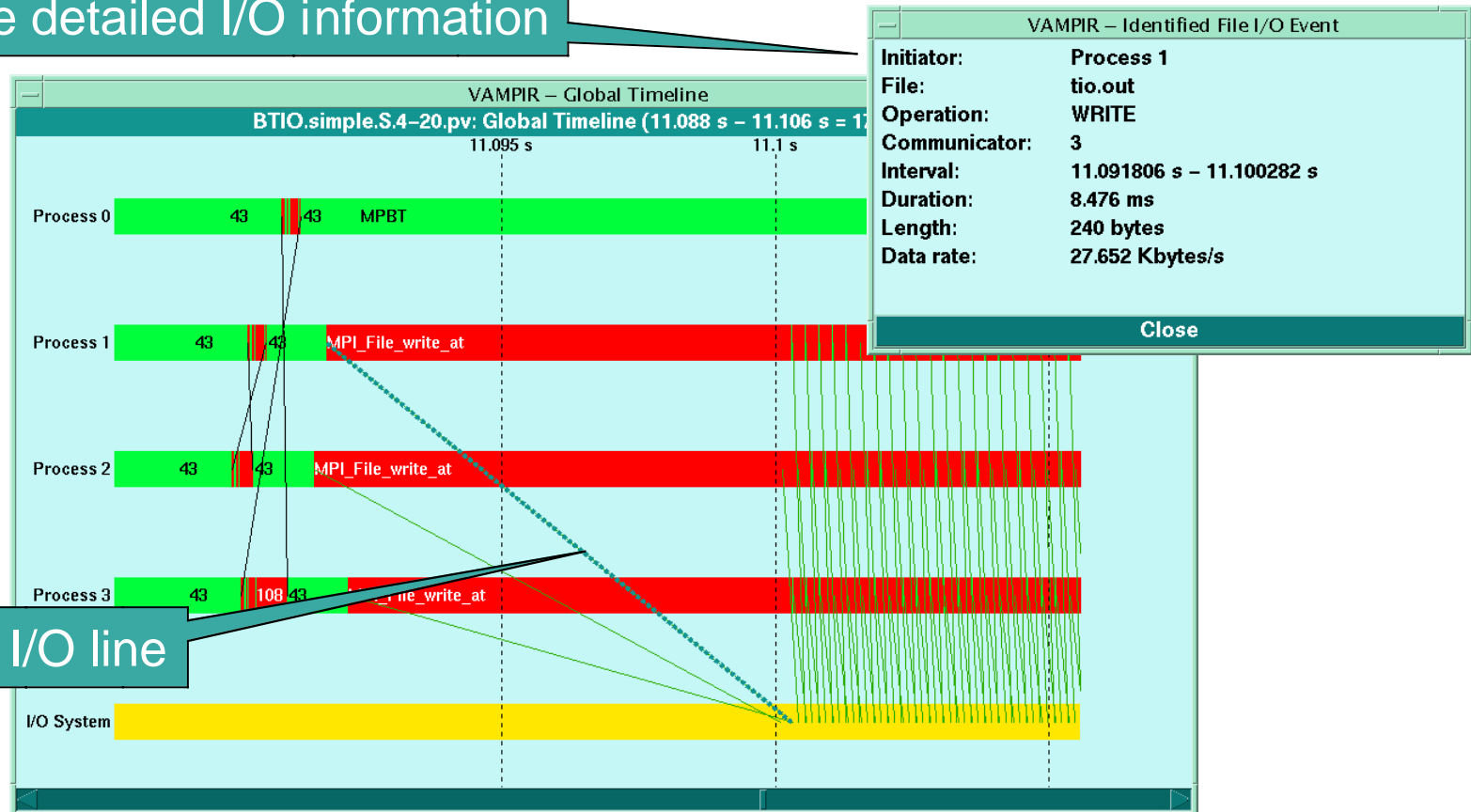| | |
|---|---|
| Initiator: | Process 1 |
| File: | tio.out |
| Operation: | WRITE |
| Communicator: | 3 |
| Interval: | 11.091806 s – 11.100282 s |
| Duration: | 8.476 ms |
| Length: | 240 bytes |
| Data rate: | 27.652 Kbytes/s |

Close

**VAMPIR – Global Timeline**

BTIO.simple.S.4–20.pv: Global Timeline (11.088 s – 11.106 s = 1?

11.095 s          11.1 s

Process 0    43    43    MPBT

Process 1    43    43    MPI_File_write_at

Process 2    43    43    MPI_File_write_at

Process 3    43    108  43    _ile_write_at

Click on I/O line

I/O System

- I/O transfers are shown as lines

# shmem Operations

See details

Click on collective op

Click on "message" line

See details

**VAMPIR**
**shmem–8.bpv: State Dialog**
On: Process 1
Activity: SHMEM 31 = shmem_barrier
From: 35.959ms To: 35.98ms (21.053us)
Transfer: 0 bytes sent, 0 bytes received
**Close**

**VAMPIR – Identified Message**
Message sent from Process 5 to Process 4
Communicator: 0, Type: 2
Length: 65600
sent at 36.127ms, received at 36.407ms (diff. 0.279ms)
Data rate: 234.744MBytes/sec
**Close**

shmem–8.bpv: Global Timeli

36.0ms    36.1ms

Process 0  check_get  128  reset_target  shmem
Process 1  cb  128  reset_target  shmem_get
128  reset_target  shmem_get
Process 3  check_get  reset_target  shmem_get
Process 4  check_get  128  reset_target  shmem_get
Process 5  check_get  128  reset_target  shmem_get
Process 6  check_get  res  shmem_ge
128  reset_target  shmem_get

- Display one–sided transfers as messages
- Display shmem global operations
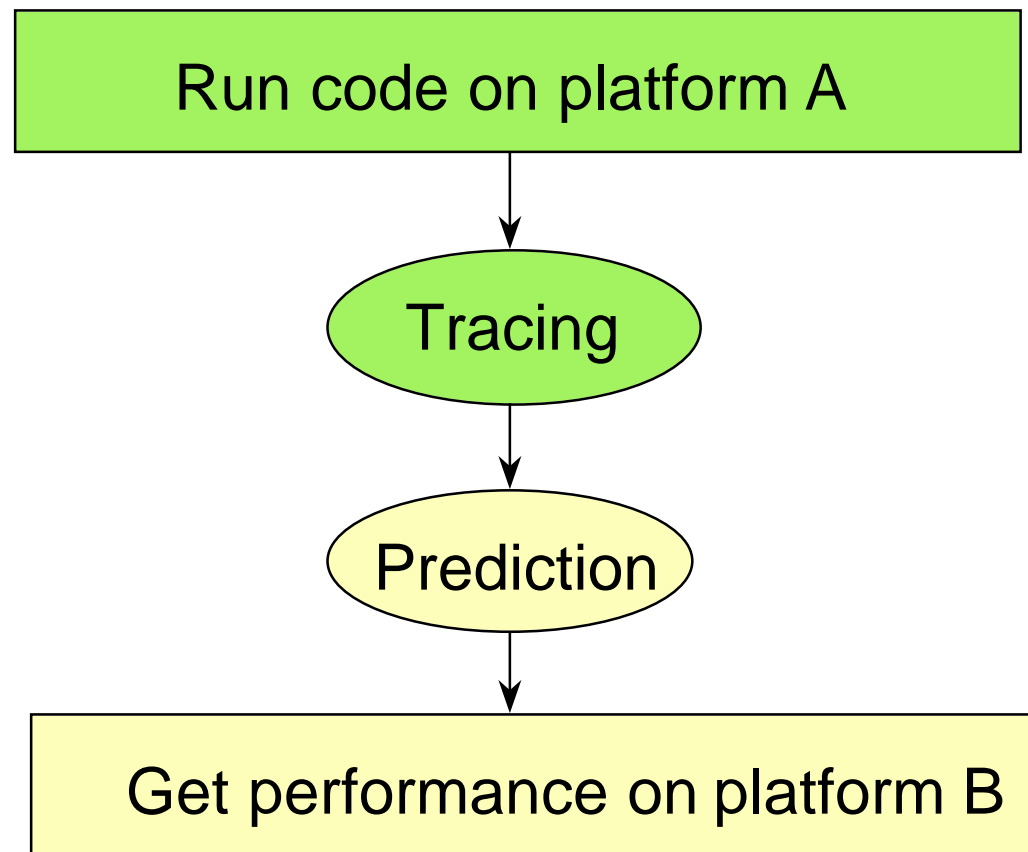
# Vampir Track Record

- Reference customers: ARL, ARSC, CEWES, LANL, LLNL, MHPCC, NASA, NERSC, NSA, Cornell TC, Oregon Univ., CEA, DWD, ECMWF, GMD, HLRS, LRZ, PC$^2$, RUKA, ...

- URLs:
  - www.tc.cornell.edu/Edu/Tutor/Vampir
  - www.llnl.gov/sccd/lc/DEG/vampir/vampir.html
  - www.uni-karlsruhe.de/~Vampir
  - www.lrz-muenchen.de/services/software/parallel/vampir
  - www.hlrs.de/structure/support/parallel_computing/tools/performance/vampir.html

**pallas**

**Dimemas**

**P**erformance
**P**rediction
**M**ade
**E**asy

Run code on platform A
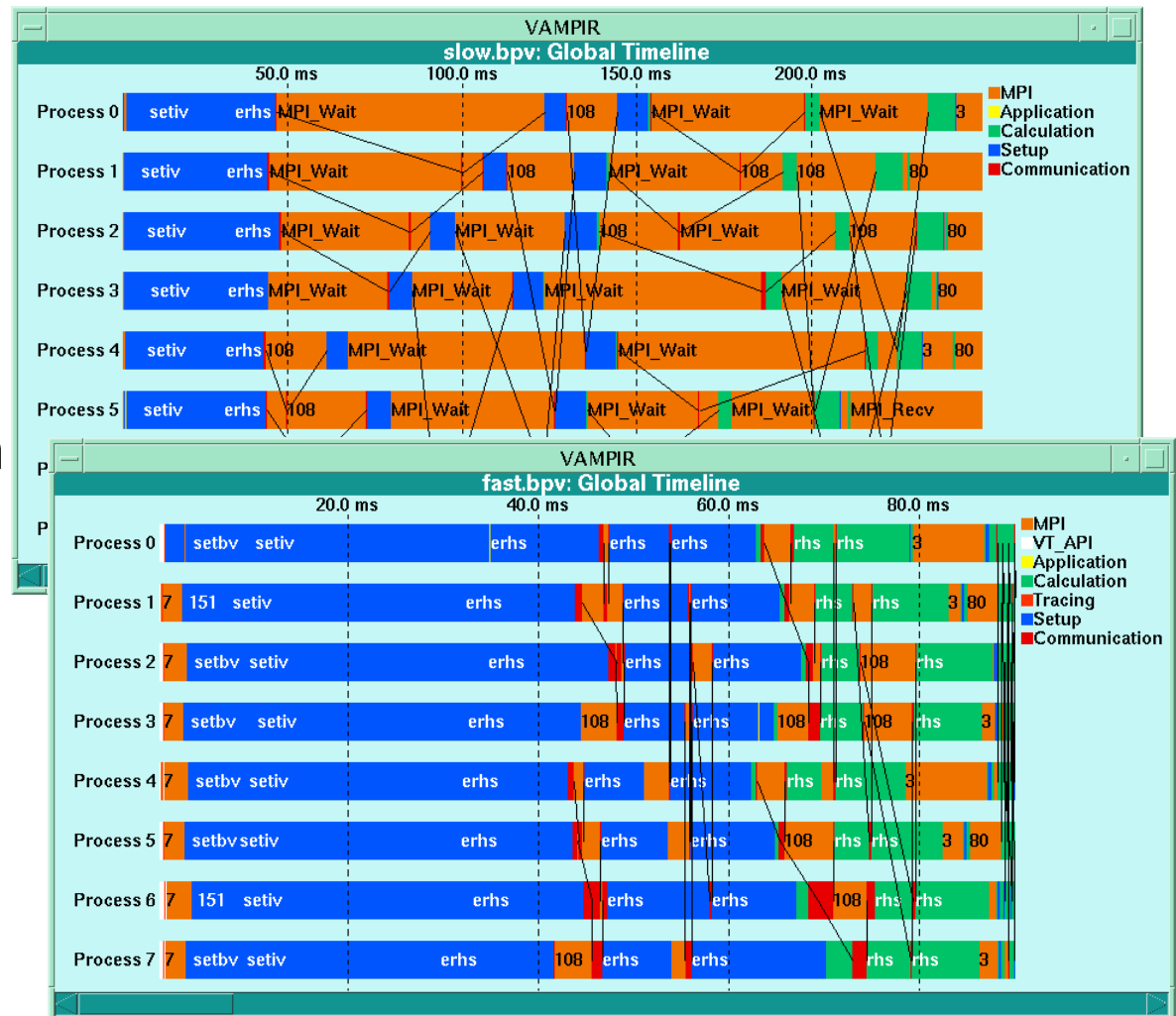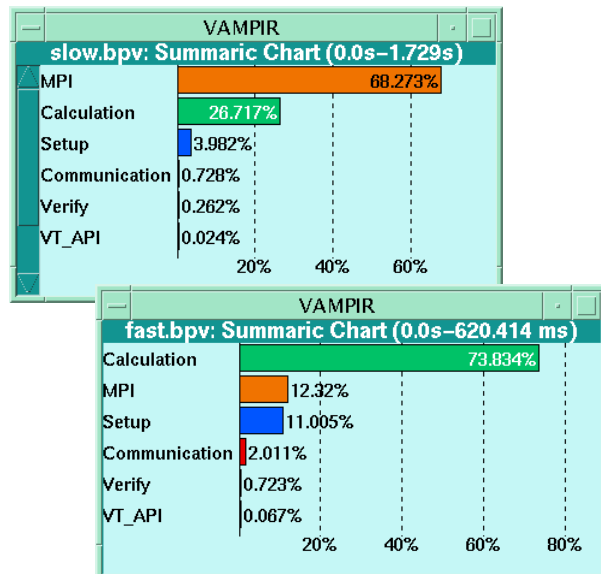
Tracing

Prediction

Get performance on platform B

# Dimemas Outputs – Vampir Tracefiles

- Actual program run

  vs.

- Ideal communication



© Pallas GmbH
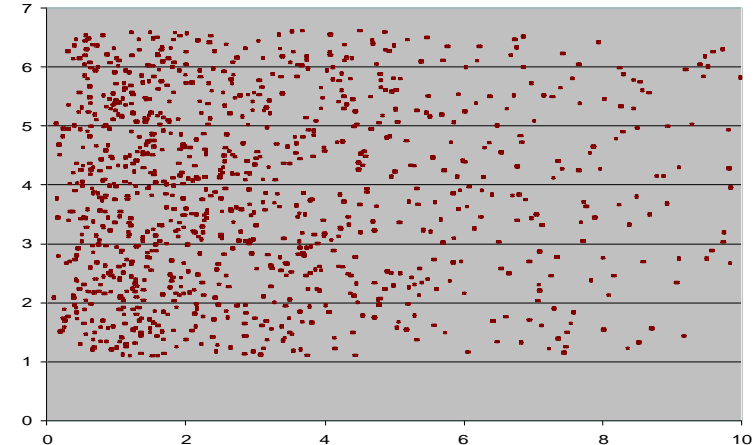
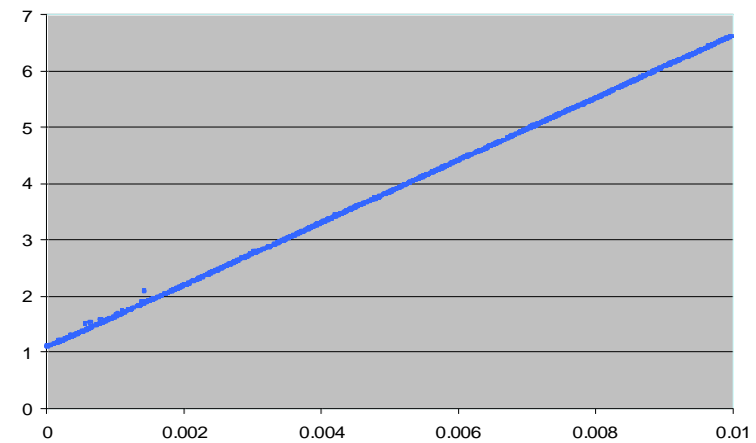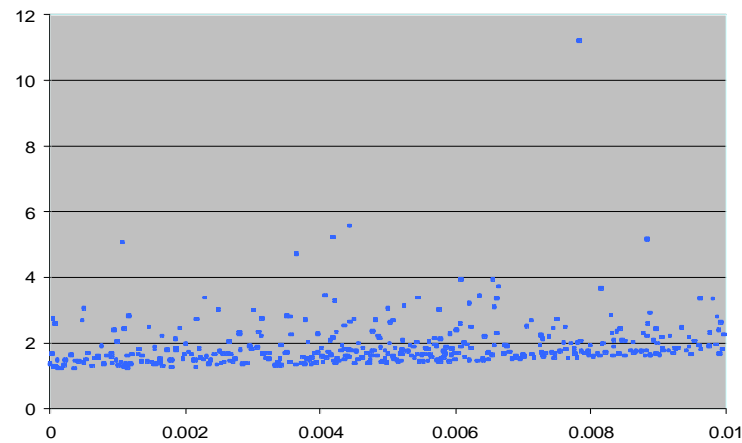# Dimemas - Prediction studies

## FFT

## PDE



Communication bandwidth (MB/s)

Communication latency  (s)

# Example Results – Performance Prediction

- NAS LU benchmarks

| Benchmark | Class | Time | Simulated | |
|-----------|-------|------|-----------|------|
| BT | A | 316.580000 | 319.858708 | 1.04% |
| BT | W | 6.410000 | 5.548973 | 13.43% |
| CG | A | 4.480000 | 4.210193 | 6.02% |
| CG | B | 342.860000 | 335.641787 | 2.11% |
| CG | W | 1.130000 | 1.210509 | 7.12% |
| EP | A | 18.970000 | 18.843001 | 0.67% |
| EP | B | 75.990000 | 77.122498 | 1.49% |
| EP | W | 2.380000 | 2.391074 | 0.47% |
| LU | A | 131.660000 | 129.139735 | 1.91% |
| LU | B | 793.070000 | 785.966349 | 0.90% |
| LU | W | 16.550000 | 17.415359 | 5.23% |
| MG | A | 16.850000 | 16.290057 | 3.32% |
| MG | B | 59.480000 | 64.974789 | 9.24% |
| MG | W | 1.052000 | 1.112022 | 5.71% |
| SP | A | 152.020000 | 151.821937 | 0.13% |
| SP | W | 16.950000 | 15.237959 | 10.10% |

## Future plans - Vampir, Dimemas

- Towards automatic performance analysis
    - improve user guidance in Vampir and Dimemas
    - add "assistant" module for inexperienced users
- Support for clustered shared–memory systems
    - support shared–memory programming models (threads, OpenMP)
    - expose cluster structure
    - aggregate information on SMP nodes
- Support for (very) large systems
    - new structured tracefile format
    - fine–grain interactive control over tracing
    - scalable displays
    - new Vampir structure (can exploit parallelism)

# pallas

## Access to Pallas Tools

**Download free evaluation copies**

**http://www.pallas.com**

# Thanks for your attention!

pallas

Pallas GmbH
Hermülheimer Straße 10
D–50321 Brühl,
Germany

info@pallas.com
www.pallas.com