



TECHNISCHE
UNIVERSITÄT
DRESDEN

Zentrum für Informationsdienste und Hochleistungsrechnen

HRSK Overview Phobos

Zellescher Weg 12

Tel. +49 351 - 463 - 35450

Dresden, 4 Juli 2006

Wolfgang E. Nagel, Matthias S. Müller



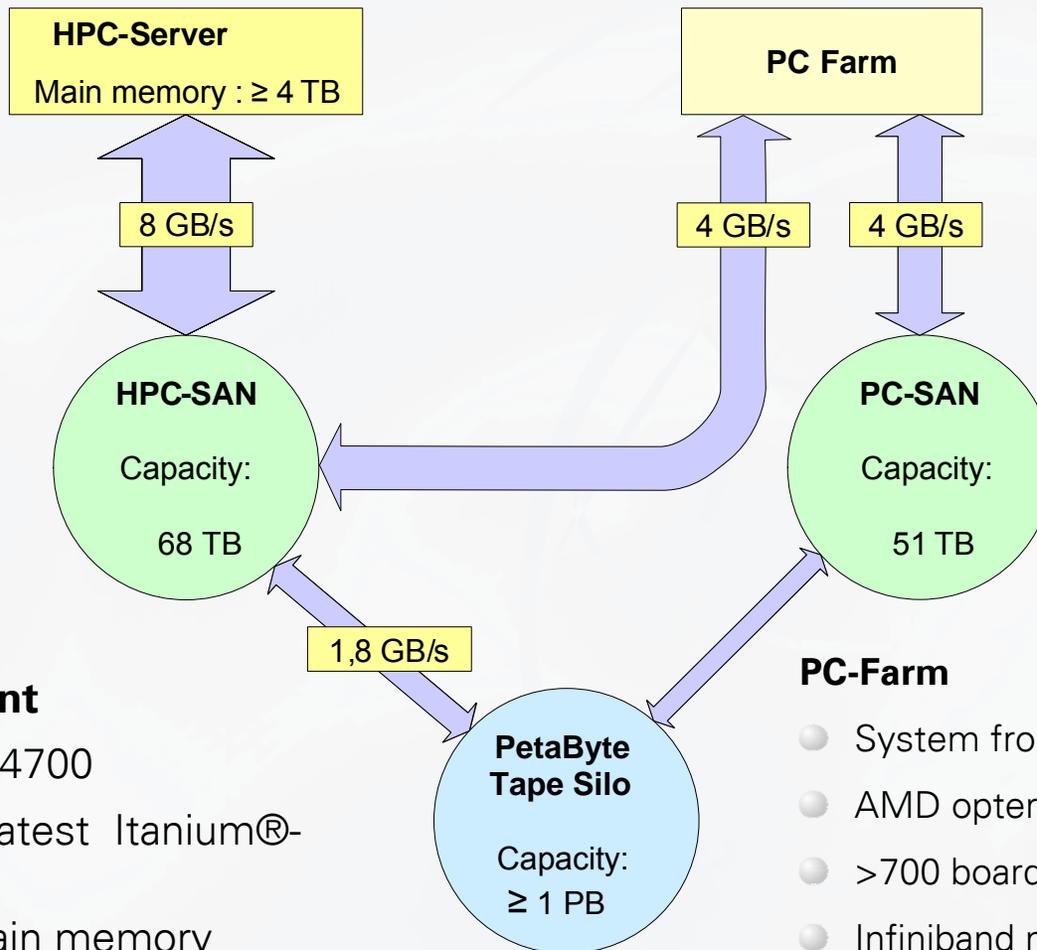
Center for Information Services and HPC (ZIH)

- Central Scientific Unit at TU Dresden
- Merged institution: TUD Computing Center (URZ) and Center for High Performance Computing (ZHR)
- Competence Center for „Parallel Computing and Software Tools“
- Strong commitment to support *real users*
- Development of algorithms and methods: Cooperation with users from all departments

Responsibilities of ZIH

- Providing infrastructure and qualified service for TU Dresden and Saxony
- Research topics
 - Architecture and performance analysis of High Performance Computers
 - Programming methods and techniques for HPC systems
 - Software tools to support programming and optimization
 - Modeling algorithms of biological processes
 - Mathematical models, algorithms, and efficient implementations
- Role of mediator between vendors, developers, and users
- Pick up and preparation of new concepts, methods, and techniques
- Teaching and Education

Procurement: Overall Infrastructure / Future Directions



HPC-Component

- SGI® Altix® 4700
- 2048 of the latest Itanium®-Cores
- 6.5 TByte main memory

PC-Farm

- System from Linux Network
- AMD opteron CPUs
- >700 boards with >2500 cores
- Infiniband networks between the nodes

Timeline

	2005						2006									
	Jul	Aug	Sep	Oct	Nov	Dec	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	
Machine Room Upgrade																
Installation Stage 1a (Test operation)																
Installation Stage 1b																
Installation Stage 2																



HPC-Component

- Next generation Shared-Memory-Plattform SGI® Altix® (Codename 'Tollhouse')
- more than 1500 of the latest Intel®-Itanium®-Cores
- 6 TByte main memory

- "Capability Computing" system
- Focus: „In-Memory“-Computing
- New research instrument to extract knowledge from data
- "Accelerator of the Theory"

PC-Farm

- System from Linux Networx
 - CPUs: Always latest most powerful chips from AMD
 - Overall more than 700 boards with more than 2000 cores
 - Infiniband networks between the nodes
-
- “Capacity Computing” system
 - Focus: Throughput with excellent performance
 - Heterogeneous diversity is part of the design

Stage1a

Zellescher Weg 12

Tel. +49 351 - 463 - 35450

Wolfgang E. Nagel, Matthias S. Müller

HRSK Stage 1a – HPC-Server



- SGI Altix 3700 Bx2
- 192x 1.5GHz/4MB L3 Cache Itanium2 CPUs
- 1152 GFlops/s Peak Performance
- 768 GB Shared Memory (4 GB/CPU), NUMA
- 1 TB lokal discs + 34 TB SAN
- SuSE SLES 9 inkl. SGI ProPack 4
- Intel Compiler and Tools:
 - C/C++ , Fortran Compiler
 - MKL
 - Vtune, Trace Collector, Trace Analyzer
- Alinea DDT Debugger
- Batchsystem LSF

merkur.hrsk.tu-dresden.de

HPC-SAN Stage 1a



- 1x DDN RAID System S2A9500:
 - 2x 38/42U Racks (air cooled system)
 - 2x S2A9500 Couplet (5GB Cache, 8x FC4 Host Ports)
 - 16x Optical SFP Interface Module
 - 20x Blue 16-Slot Dual-Loop FC 2Gb JBOD,
 - 20x Rack Mount Kits
 - 292x 146GB 10k RPM FC Disks (4 hot spare)
 - 34 TB net capacity
 - DIRECT-GUI

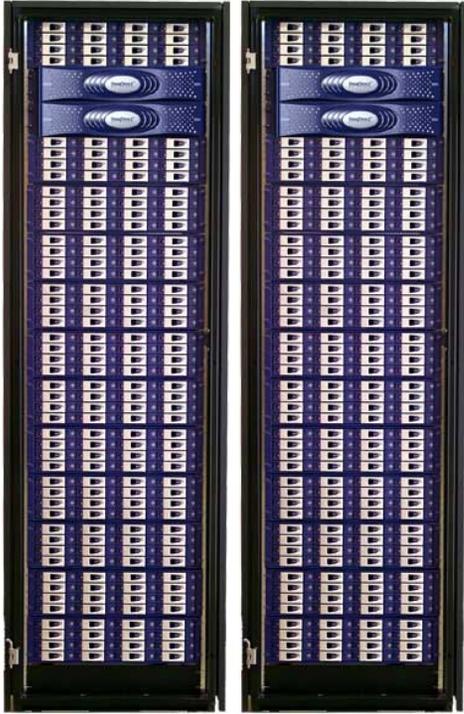
PC-Farm Stage 1a



phobos.hrsk.tu-dresden.de

- 64 dual CPU nodes
- 128 AMD Opteron DP248 2.2 GHz (Single-Core) CPUs
- 563,2 GFlops/s peak performance
- 256 GB main memory (4GB per node)
- SUSE operating system
- Infiniband 4x Interconnect
- 80 GB local disc per node
- 21,2 TB shared disc space:
 - 2x DDN RAID System S2A9500

PC-SAN Stage 1a



- DDN RAID System S2A9500:
 - 2x 38/42U Racks (air cooled),
 - 1x S2A9500 Couplet (5GB Cache, 6x FC4 Host Ports),
 - 8x Optical SFP Interface Module,
 - 19x Blue 16-Slot Dual-Loop FC 2Gb JBOD,
 - 19x Rack Mount Kits,
 - 184x 146GB 10k RPM FC discs (4 Hot-Spare),
- 21,2 TB netto capacity

Stage 1b,2

Zellescher Weg 12

Tel. +49 351 - 463 - 35450

Wolfgang E. Nagel, Matthias S. Müller

Comparison of systems in 2005

	HPC Server		PC Farm	
Peak Performance	1152	+	563	+
#CPUs	192		128	
Nodes	1	+	64	
Total Memory	768 GB		256 GB	
Total Disc	34 TB		21 TB	
MPI		+		+
OpenMP		++		
Usable Shared Memory	768 GB	++	4 GB	-
MPI Bandwidth	~ 3.2 GB/s	++	~ 1GB/s	+
MPI Latency	~ 1.5-3 microsecond	++	~ 4-5 microseconds	+

PC-Farm Stage 1b



- Compute Partition:
 - 128 single CPU, Opteron dual core
 - 104 dual CPU, Opteron dual core
 - 24 quad CPU, Opteron dual core
- AMD Opteron dual core 2.6 GHz
- 3801,6 GFlops/s peak performance
- 864GB main memory (2GB per CPU chip)
- SUSE operating system
- Infiniband 4x Interconnect
- 21,2 TB shared disc space:
 - DDN RAID System S2A9500

PC-Farm Stage 2



- Compute Partition:
 - 256 single CPU, Opteron dual core
 - 128 dual CPU, Opteron dual core
 - 88 quad CPU, Opteron dual core
- AMD Opteron dual core 2.6 GHz
- 8294,4 GFlops/s peak performance
- 3456 GB main memory (2GB per core)
- SUSE operating system
- Infiniband 4x Interconnect
- 51 TB total shared disc space:
 - DDN RAID System S2A9500

PC-Farm final config



- Compute Partition:
 - 128+256 single CPU, Opteron dual core
 - 104+128 dual CPU, Opteron dual core
 - 24+88 quad CPU, Opteron dual core
- AMD Opteron dual core 2,6 GHz
- 13478 GFlops/s peak performance
- 4320GB main memory
 - 1-2GB per core
 - 2-32GB per node
- SUSE operating system
- Infiniband 4x Interconnect
- 51 TB total shared disc space:
 - DDN RAID System S2A9500

HPC-Server Stage 2



- SGI Tollhouse
- 1024 MonteCito CPUs (2048 Cores)
- 12288 GFlops/s peak performance



- 6 TB Shared Memory, NUMA
- 68 TB Disk
- SUSE SLES 9
- Intel Compiler and Tools:



- C/C++ , Fortran Compiler
- MKL
- Vtune, Trace Collector, Trace Analyzer
- Alinea DDT debugger
- LSF batch system

Comparison of systems in 2006

	HPC Server		PC Farm	
Peak Performance	12288	+	13478	+
#Cores	2048		2592	
Nodes	4	+	728	
Total Memory	6 TB		4 TB	
Total Disc	68 TB		50 TB	
MPI		+		+
OpenMP		++		
Usable Shared Memory	2048 GB	++	<=32 GB	-
MPI Bandwidth	~ 1-2 GB/s	++	~ 1GB/s	+
MPI Latency	~ 1.5-3 microsecond	++	~ 4-5 microseconds	+

Main Contractor: SGI

HPC-component

- 2048 cores (Intel Itanium 2, Montecito)
- 6 TB Main memory with high bandwidth (major decision criteria)



SAN for HPC / PC

- 68 / 51 TB Capacity



Tape Silo

- 1 PB Capacity

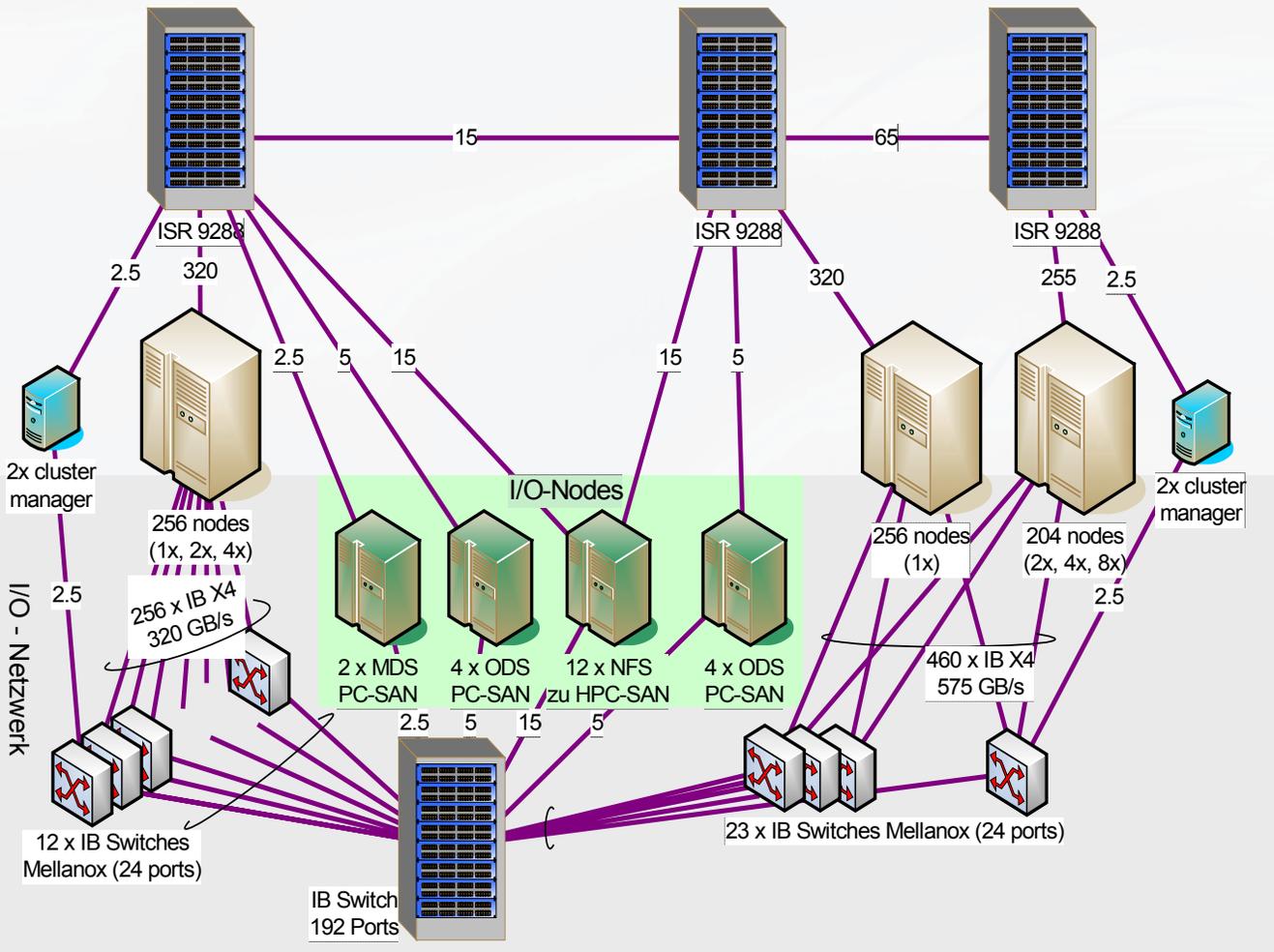


PC-Farm

- > 700 System boards (AMD Opteron)
- Infiniband network build from 288-port switches

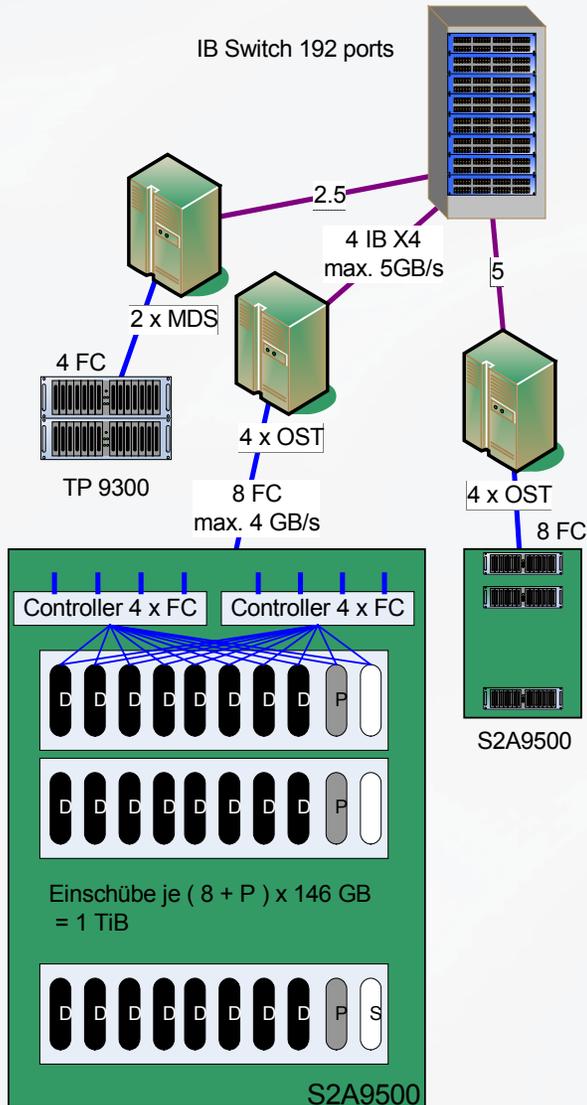


PC-Farm - Details



- Interconnect: IB X4
- Compute-Knoten connected via three switches (288 ports)
- connection between switches
- 2 x 2 Cluster Manager
- Connection to HPC-SAN with 12 NFS-Server (CXFS-Clients)
- cascaded I/O-net:
 - 35 IB X4 Switches (24 ports)
 - each with 4 x IB-Uplink to...
 - IB-Switch (192 ports)

PC-San - Details



- Lustre FS
- IB Switch (192 ports)
- 2 x DDN S2A 9500
- 4 Controller with 4 x FC 4Gb/s each
- capacity: 50,9 TB
- 4 Racks
- 440 discs with 146 GB each
- 2 x MDS je 2 dual-core Opteron 275, 4 GB RAM, dual-port FC
- 4 x ODS (OSS,OST) je 2 dual-core Opteron 275, 4 GB RAM, 2 x dual-port FC
- 4 x ODS (OSS,OST) je 2 dual-core Opteron 280, 4 GB RAM, 2 x dual-port FC (Stufe 2)
- TP 9300 (MDS Storage Subsystem) 16 x 146 GB für Metadaten

User Support at ZIH

A few examples

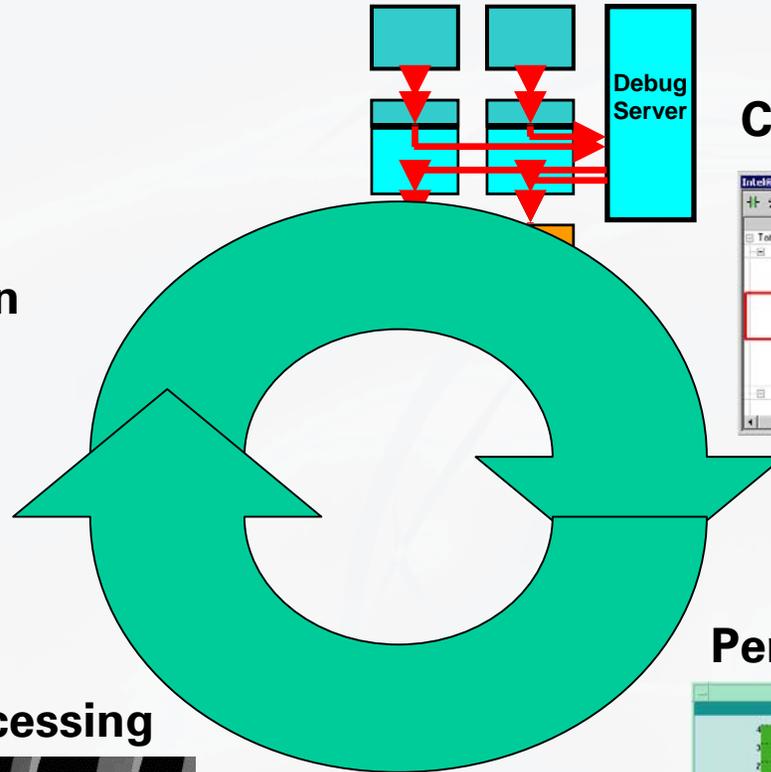
Zellescher Weg 12

Tel. +49 351 - 463 - 35450

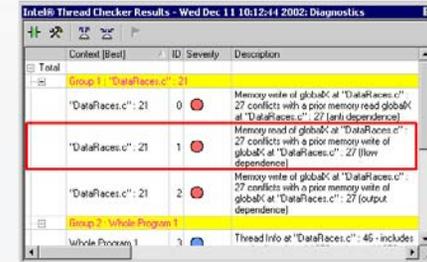
Wolfgang E. Nagel, Matthias S. Müller

Evolution of a parallel application

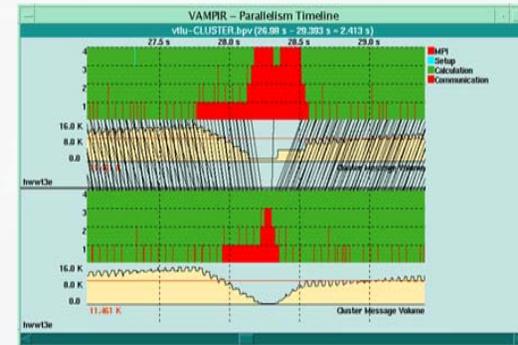
Code Modification



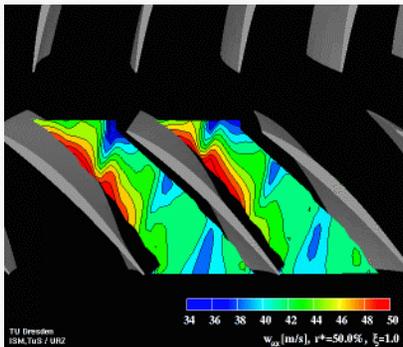
Correctness

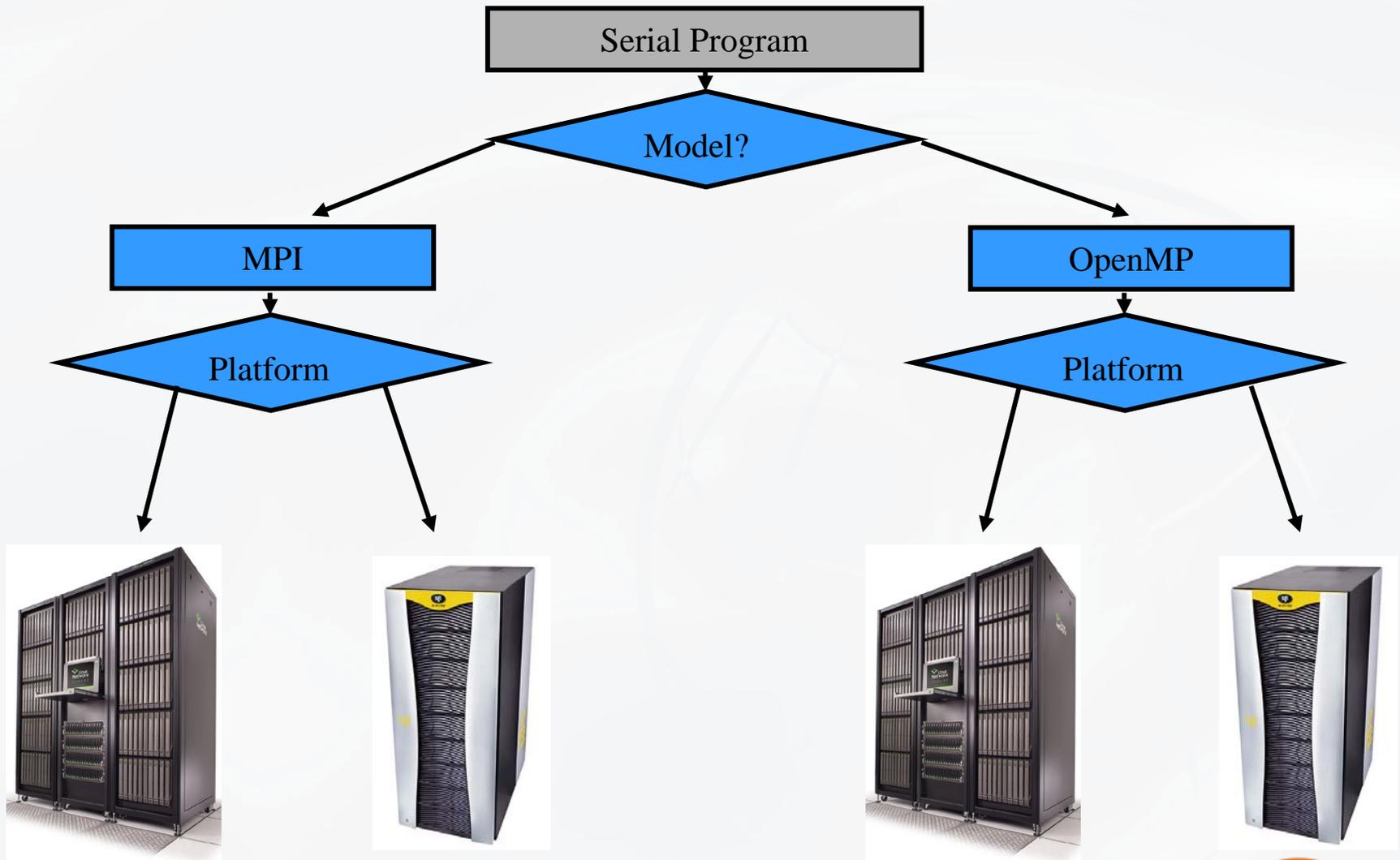


Performance



Post processing





MPI code checking with MARMOT

Tool for the development of MPI applications

Automatic runtime analysis of the application:

- Detect incorrect use of MPI
- Detect non-portable constructs
- Detect race conditions and deadlocks

MARMOT does not require source code modifications, just relinking

C and Fortran binding of MPI -1.2 is supported, also C++ and mixed C/Fortran code

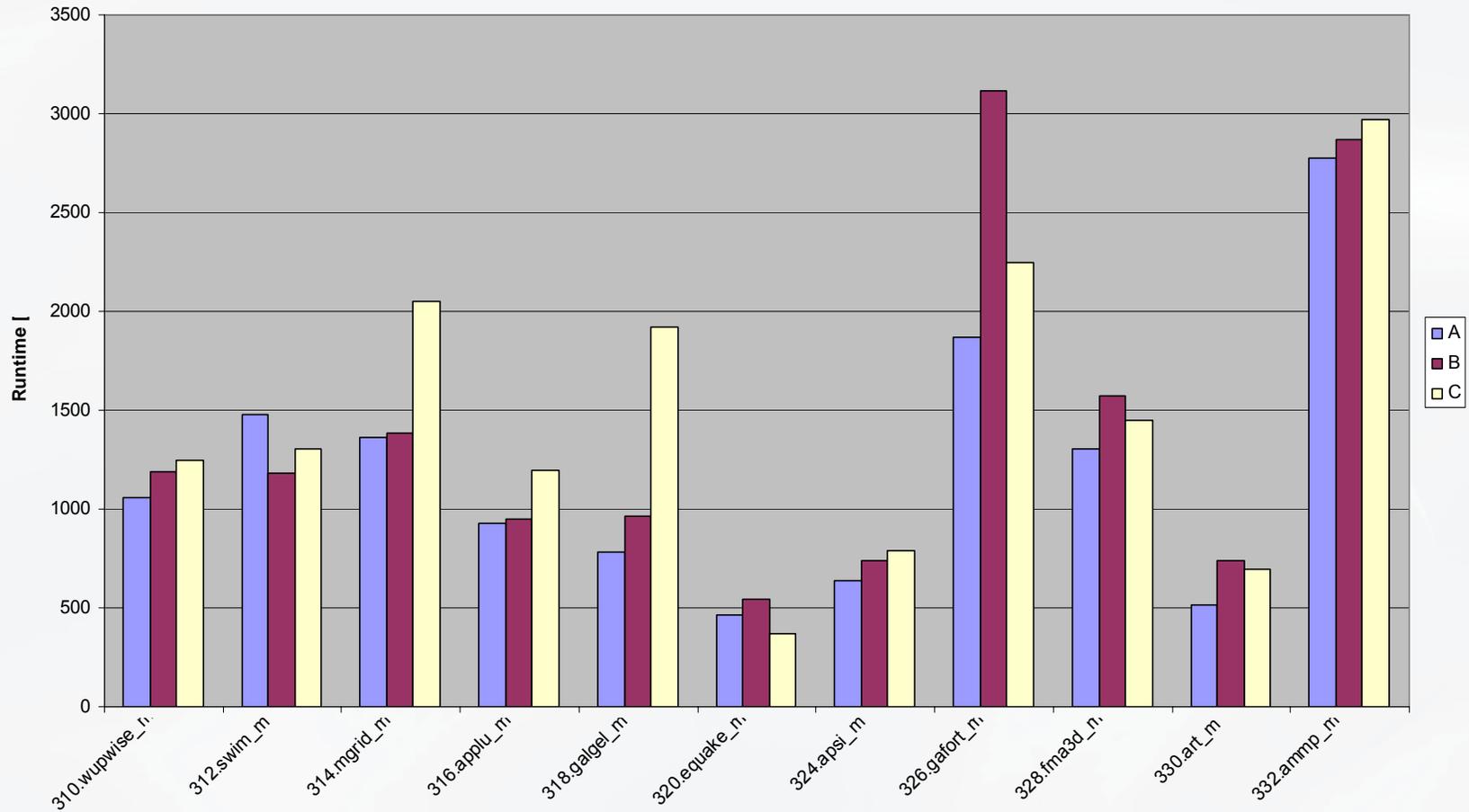
Development is still ongoing (not every possible functionality is implemented yet...)

Tool makes use of the so-called *profiling interface*

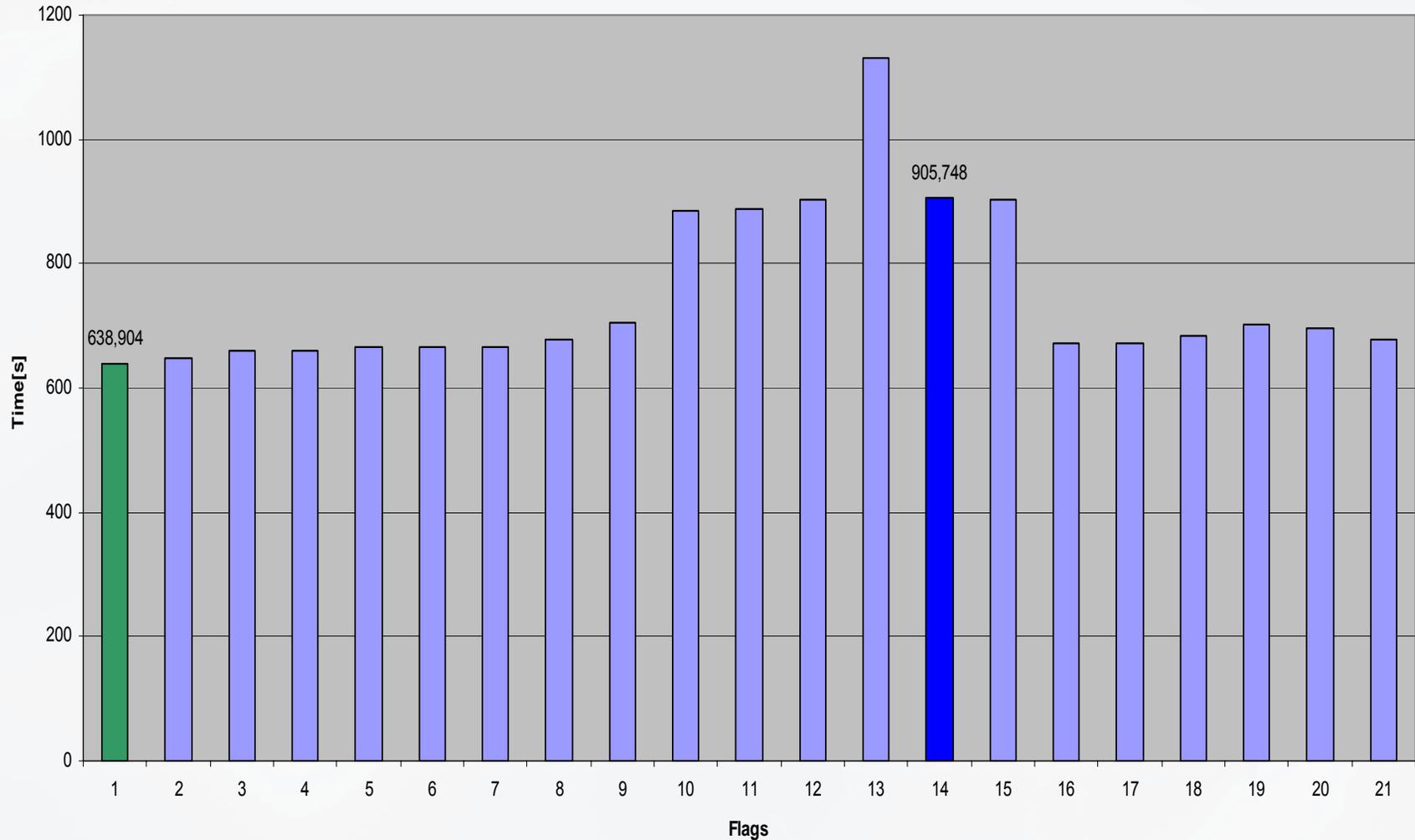
MARMOT will continue to be developed jointly by ZIH and HLRS



Compare different Compilers with SPEC OMPM2001



Code Tuning: different compiler flags



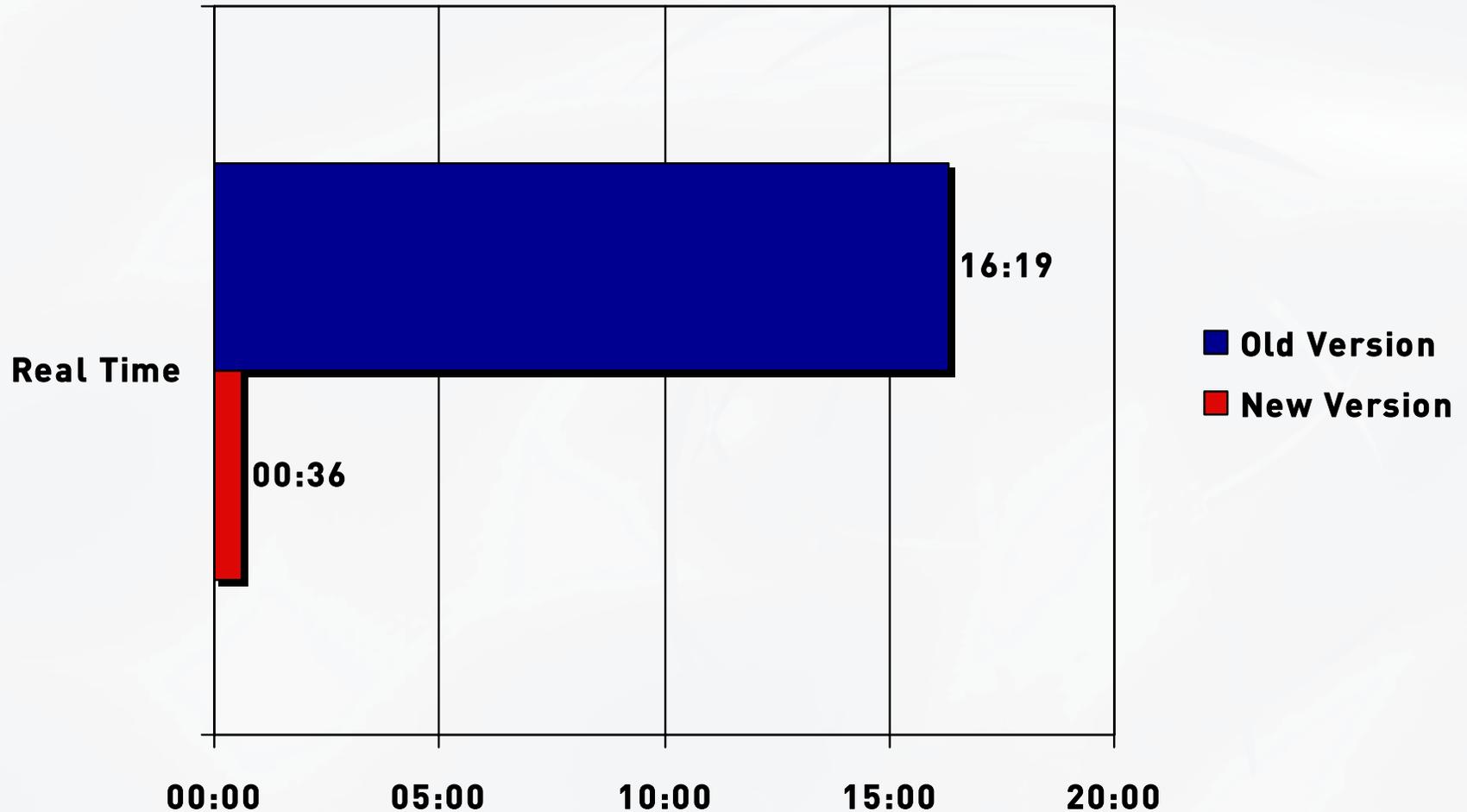
A few success stories

Zellescher Weg 12

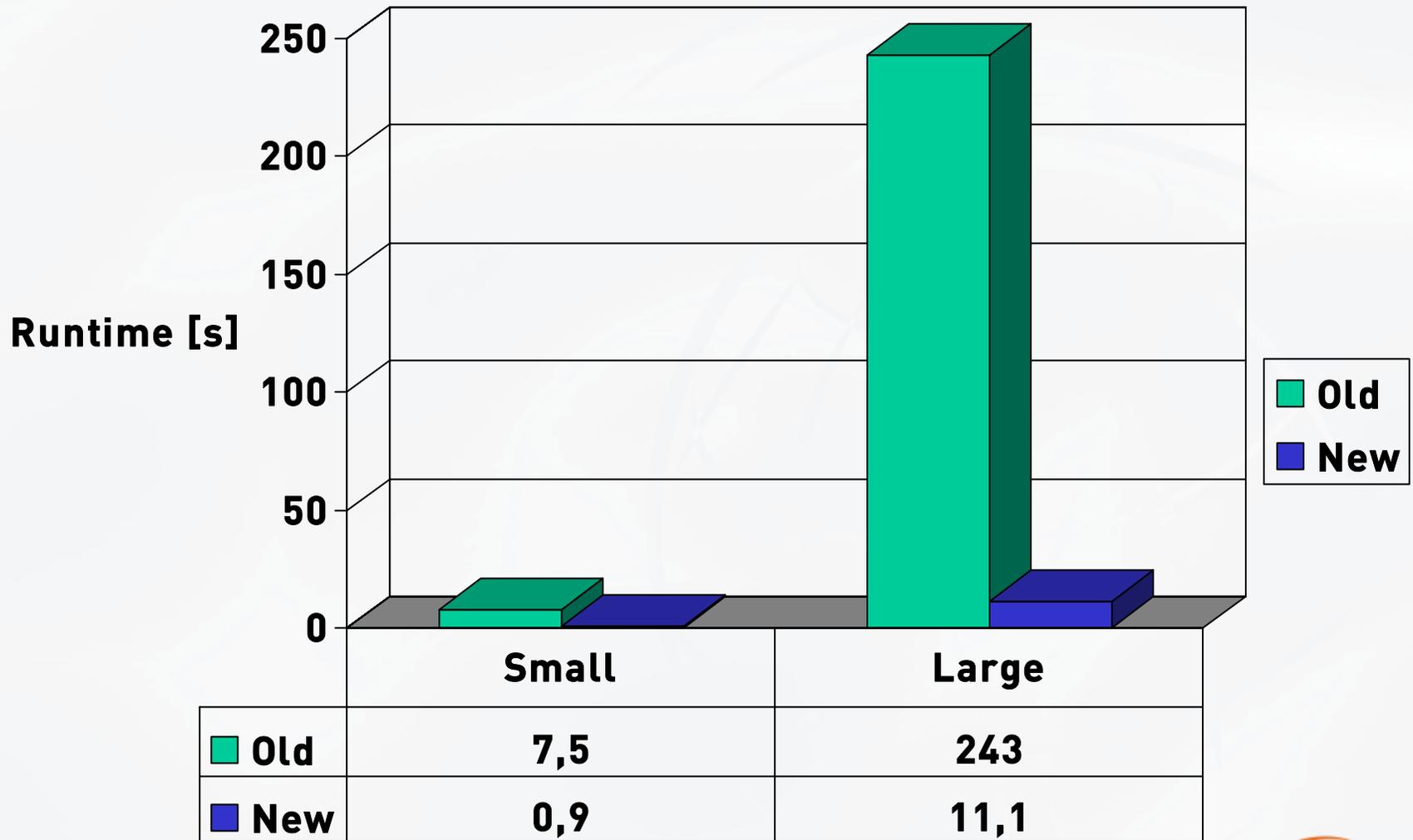
Tel. +49 351 - 463 - 35450

Wolfgang E. Nagel, Matthias S. Müller

Example: Serial Optimization

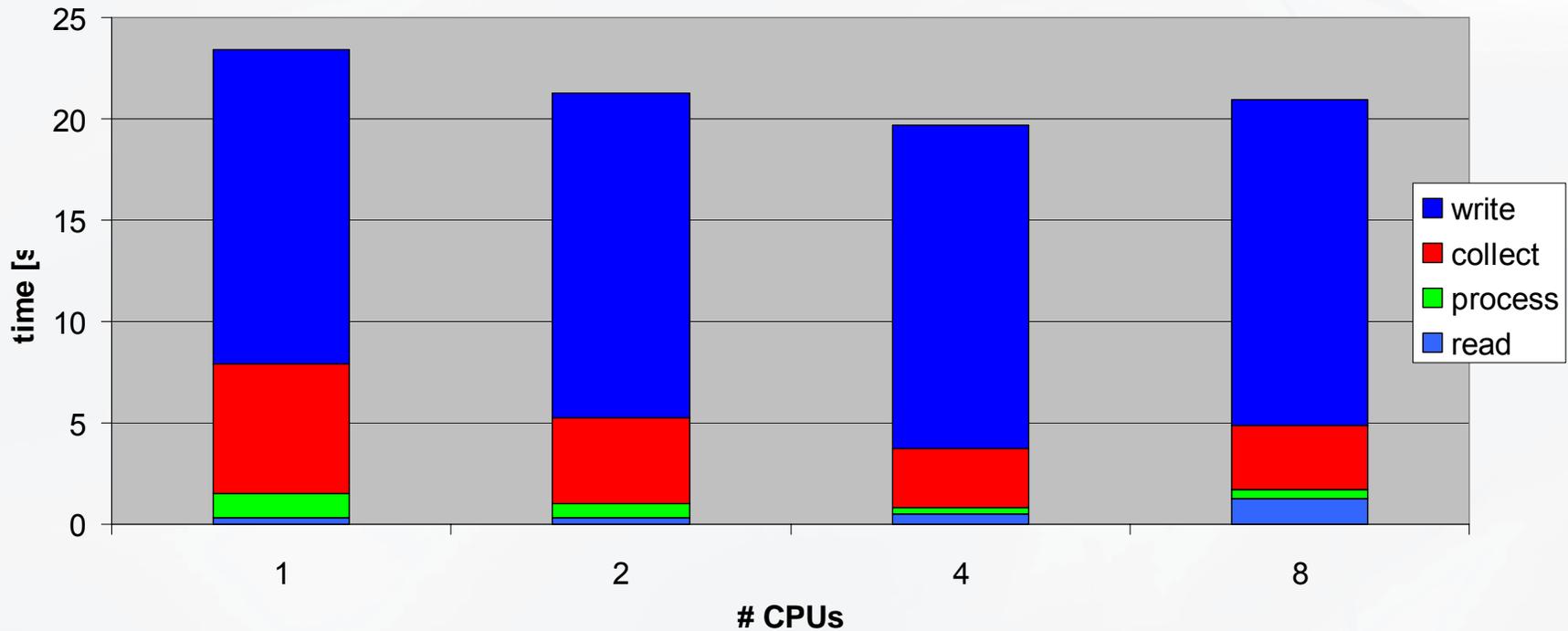


Example: Matlab Calculation



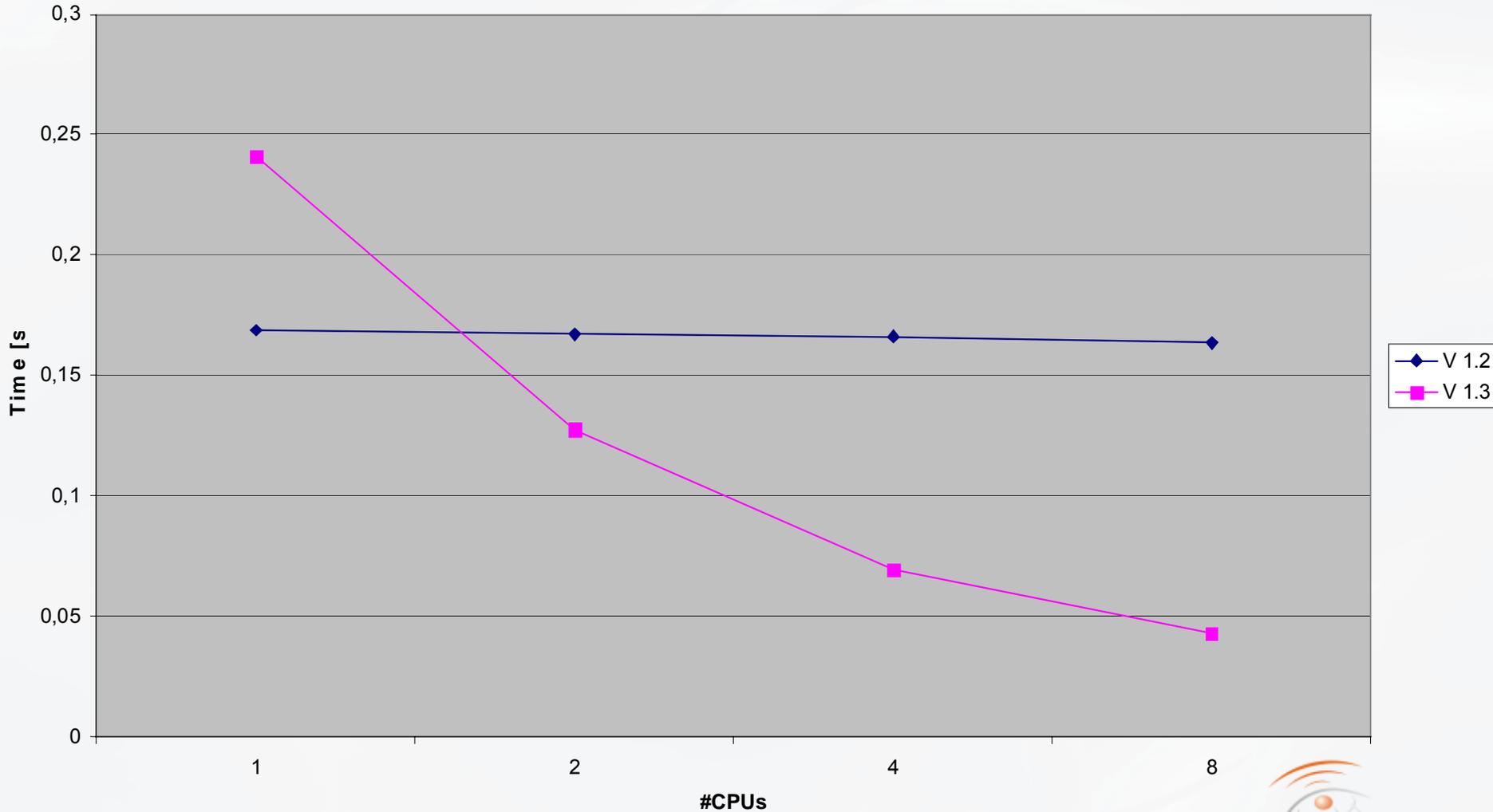
Overall performance of Geneindex 1.2

Performance on Altix (32 Mbp)



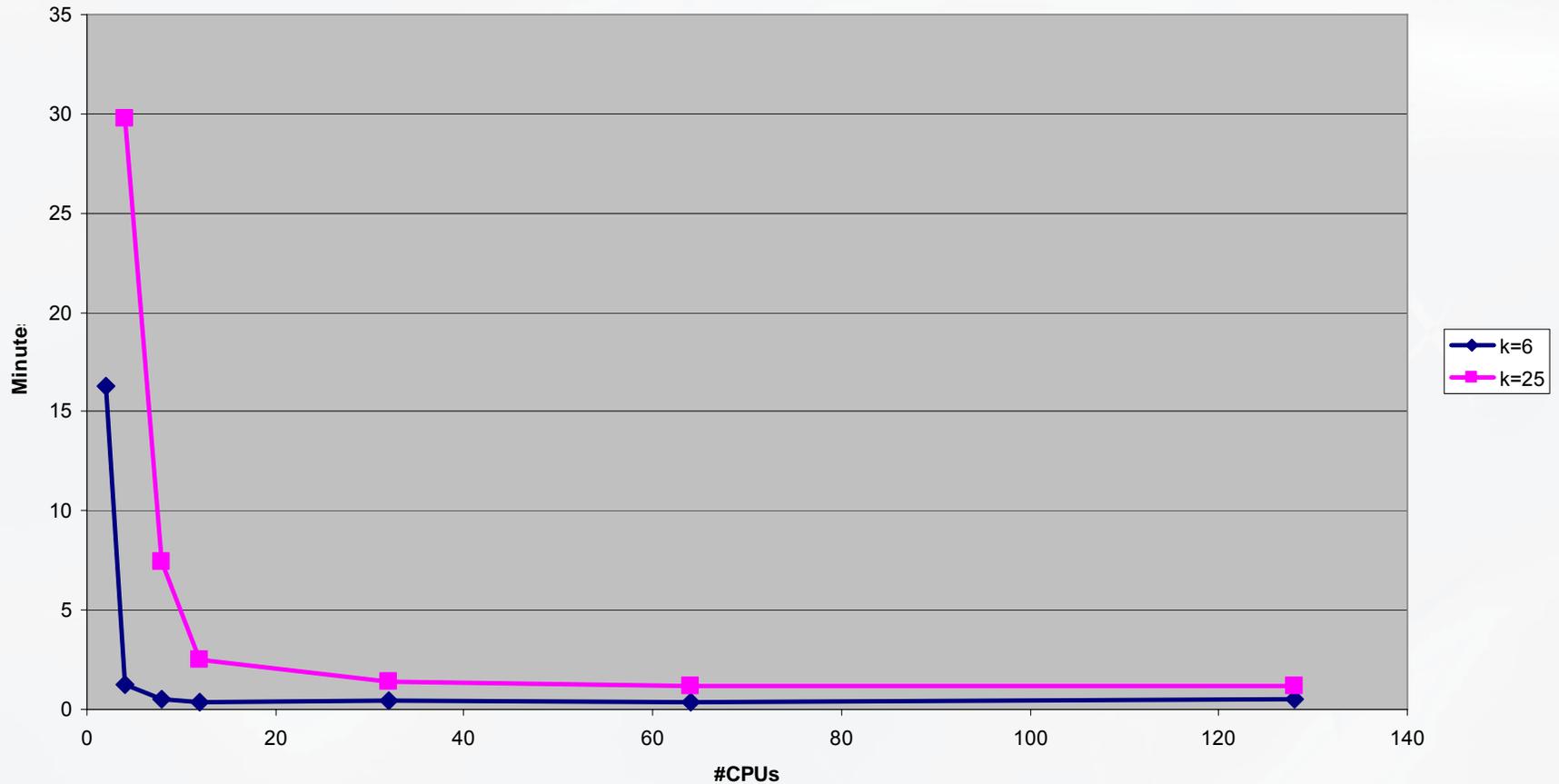
Read Performance on GPFS (davidd-o) with 32 Mbps

Read Time



Analysis of Fruit Fly Genome (165 Mbps)

Processing Time



Recent Improvements

- Printed Documentation:
Every new project receives one copy. Also available for download
- Mailing list for information of users:
zih-hpcnews
- Mailing list for discussions between users:
zih-hpcforum
- Trouble Ticket System for support (currently in beta):
HPCsupport@zih.tu-dresden.de

Summary and Conclusion

- The new installation offers a large opportunity
- It is still a long way to go
- We hope we can be of some help
- Enjoy the user training and make the best out of it



TECHNISCHE
UNIVERSITÄT
DRESDEN

Zentrum für Informationsdienste und Hochleistungsrechnen

Tools on Phobos

Zellescher Weg 12

Tel. +49 351 - 463 - 33640

Dresden, June 21 2006

Ulf Markwardt



Overview

Compilers

- Installed compilers
- Selected compiler options
- Comparison between compilers

Libraries

- Basic libraries
- Advanced mathematical libraries
- Comparison of FFT libraries (Daniel Hackenberg)

Debuggers

Profiling tools

Misc

- LSF
- Applications

Installed Compilers

- GCC, g77 version 3.3
on the way to get mature – but not yet competitor in performance

Commercial Compilers – Fortran 95/90/77, C 90, C++

- PathScale Compiler version 2.4
- PGI Compiler (Portland Group) version 6.0
- Intel Compiler version 9.0

Intel is trying to compete
currently: EM64T-code not optimized for Opteron CPU

Selected Compiler Options

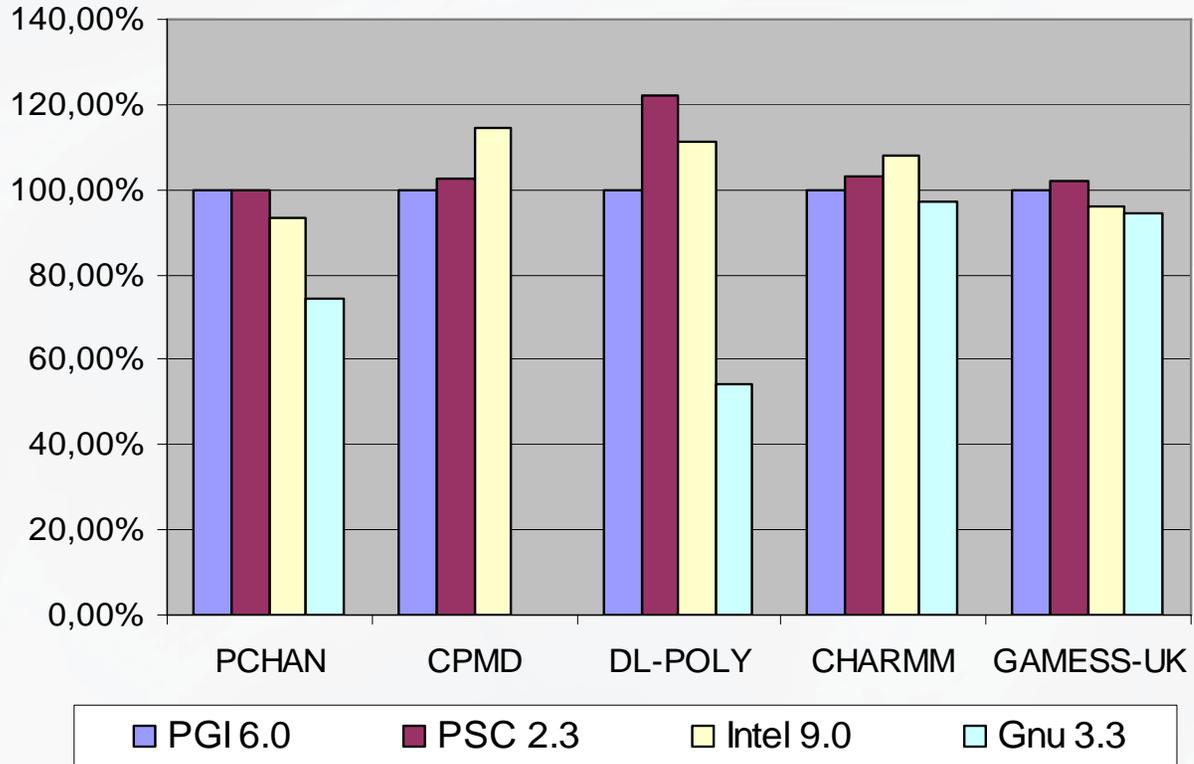
Intel	PGI	Pathscale	Description
-O0 ... -O3	-O0 ... -O3	-O0 ... -O3	Optimization level. (-O2 is a good start)
-ipo	-Mipa	-ipa	inter procedure optimization (across files)
-openmp	-mp	-mp	OpenMP support
-parallel	-Mconcur	-apo	Auto-parallelizer (SMP)
-mp	-Kieee	-fno-unsafe-math-optimizations	use this flag to limit floating-point optimizations and maintain declared precision
-mp1	-Knoieee	-ffast-math	some floating-point optimizations are allowed, less performance impact than IEEE conformity

Selected Compiler Options (continued)

Intel	PGI	Pathscale	Description
<code>-xW</code> (em64t-sup)	<code>-tp amd64</code>	<code>-mcpu=opteron</code>	optimize for Opteron processor
	<code>-fastsse</code>	<code>-msse2</code>	Use SSE instructions (increased performance for single precision)
<code>-prof-gen</code>	<code>-Mphi</code>		Generate profiling information during the run
<code>-prof-use</code>	<code>-Mpfo</code>		Use profiling information for compiling

Comparisons

Performance Comparison with PGI=100%



Data from „Performance comparison of compilers for AMD Opteron. Kozin, I.N., Deegan, M.J. et al. Computational Science an Engineering Dep. CCLRC Daresbury Laboratory, Warrington UK. 12/17/05

Basic Libraries

- MPICH over Infiniband: MVAPICH version 1.2.6
built for Intel, PGI, PathScale, Gnu – correct version for loaded module will be chosen:
 - module load pgi
 - mpicc, mpicxx, mpif77, mpif90
- OpenMPI
<http://www.open-mpi.org/>
built for Intel and Gnu
- Performance Application Programming Interface – PAPI 3.2.1
<http://icl.cs.utk.edu/papi/>
usage of performance counters
- Tracing Libraries
vampirtrace, ...

Mathematical Libraries

- Automatically Tuned Linear Algebra Software - ATLAS 3.6.0
<http://www.netlib.org/atlas/>
BLAS, LAPACK
- Intel Math Kernel Library - MKL 8.0
<http://www.intel.com/cd/software/products/asm-na/eng/perflib/mkl/index.htm>
BLAS, LAPACK, FFT, scalar and vector functions, sparse and dense matrices
- AMD Core Math Library – ACML 2.7.0
<http://developer.amd.com/acml.aspx>
BLAS, LAPACK, FFT, scalar and vector functions , sparse and dense matrices
- Gnu Scientific Library - GSL 1.8
<http://www.gnu.org/software/gsl>
functions above + differential eq., numeric differentiation, random distributions, Eigensystems, interpolation, etc.
- Numeric Algorithms Group – NAG Library Mark 8
<http://nag.co.uk/numeric>
functions optimization, integration, random distributions
- Fastest Fourier Tranfor in the West FFTW 3.1
<http://www.fftw.org>
FFT – might become obsolete with fftw2mkl – wrappers MKL 8.1

Comparison of Libraries

- Performance depends on problem size
- Users are strongly recommended to test the best

Debuggers

- gdb – GNU Debugger
GUI: /usr/bin/ddd
no threading support
- Intel: idb
- Pathscale: pathdb

```
DDD: /work/home/mark/ftw/ftw-3.1.1/libbench2/timer2.c <@phobos>
File Edit View Program Commands Status Source Data Help
0: Tl
2: Tj 0
3: Tk 1.77964905e-43
4: Tl 5.88545355e-44

float TG;
{
    float T1, T2, Tj, Tk;
    T1 = ri[0];
    T2 = ri[WS(is, 4)];
    T3 = T1 + T2;
    Tn = T1 - T2;
    {
        float Tg, Th, T4, T5;
        Tg = ii[0];
        Th = ii[WS(is, 4)];
        Ti = Tg + Th;
        TC = Tg - Th;
        T4 = ri[WS(is, 2)];
        T5 = ri[WS(is, 6)];
        T6 = T4 + T5;
        TB = T4 - T5;
    }
    Tj = ii[WS(is, 2)];
    Tk = ii[WS(is, 6)];
    Tl = Tj + Tk;
}

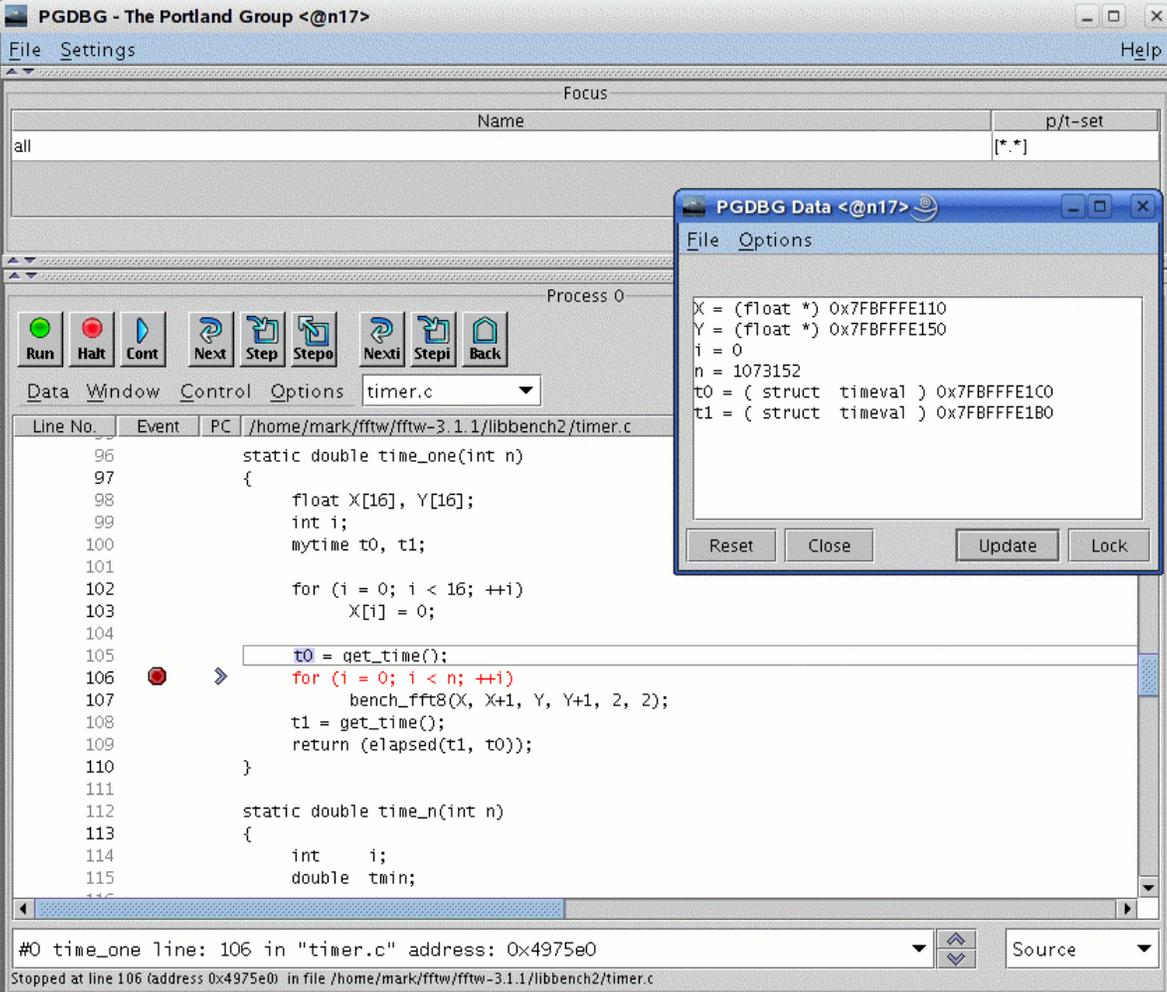
(gdb) step
(gdb) step
(gdb) graph display *T1
!!!Disabling display 1 to avoid infinite recursion
(gdb) list 'T1'
Line number not known for symbol "T1"
(gdb) graph display Tj
(gdb) graph display Tk
(gdb) step
(gdb) graph display Tl
(gdb) !

Display 4: Tl (enabled, scope bench_fft8, address 0x7fbfffe4ac)
```


Debuggers

pgdgb

- module load pgi
- Multi-threaded and OpenMP applications
- MPI applications **NOT YET** debuggable



PGDBG - The Portland Group <@n17>

File Settings Help

Focus

Name	p/t-set
all	[*.*]

Process 0

Run Halt Cont Next Step Stepo Nexti Stepi Back

Data Window Control Options timer.c

Line No.	Event	PC	/home/mark/fftw/fftw-3.1.1/libbench2/timer.c
96			static double time_one(int n)
97			{
98			float X[16], Y[16];
99			int i;
100			mytime t0, t1;
101			
102			for (i = 0; i < 16; ++i)
103			X[i] = 0;
104			
105			t0 = get_time();
106	●		for (i = 0; i < n; ++i)
107			bench_fft8(X, X+1, Y, Y+1, 2, 2);
108			t1 = get_time();
109			return (elapsed(t1, t0));
110			}
111			
112			static double time_n(int n)
113			{
114			int i;
115			double tmin;
116			

PGDBG Data <@n17>

File Options

```
X = (float *) 0x7FBFFFE110
Y = (float *) 0x7FBFFFE150
i = 0
n = 1073152
t0 = ( struct timeval ) 0x7FBFFFE1C0
t1 = ( struct timeval ) 0x7FBFFFE1B0
```

Reset Close Update Lock

#0 time_one line: 106 in "timer.c" address: 0x4975e0

Stopped at line 106 (address 0x4975e0) in file /home/mark/fftw/fftw-3.1.1/libbench2/timer.c

Profiling

- PathScale:
 - Compile code with `-pg`
 - Run the program
 - Start `pathprof`
- Gnu:
 - Compile code with `-pg`
 - Run the program
 - Start `gprof`

Profiling with pgprof

PGPROF - The Portland Group <@phobos> <2>

File Settings Help

Profiled: a.out on Mon Jul 03 10:21:47 CEST 2006 for 0.033764 seconds

Processes View Sort Search

View	Line	bench.c@mkplan_complex_interleaved	Count	Time
<input type="checkbox"/>	419	{	14%	0.000008 = 0%
<input type="checkbox"/>	420	FFTW(plan) pln;		
<input type="checkbox"/>	421	bench_tensor *sz = p->sz, *vecsz = p->vecsz;	14%	0.000009 = 0%
<input type="checkbox"/>	422			
<input type="checkbox"/>	423	if (vecsz->rnk == 0 && tensor_unitstridep(s	14%	0.000007 = 0%
<input type="checkbox"/>	424	goto api_simple;	14%	0.000007 = 0%
<input type="checkbox"/>	425			
<input type="checkbox"/>	426	if (vecsz->rnk == 1 && expressible_as_api_r		
<input type="checkbox"/>	427	goto api_many;		
<input type="checkbox"/>	428			
<input type="checkbox"/>	429	goto api_guru;		
<input type="checkbox"/>	430			
<input type="checkbox"/>	431	api_simple;	14%	0.000001 = 0%

Sort By Line No

Process	Count	Time
0 (419)	14%	0.000008 = 0%

Sort By Process

Profile: /work/home/mark/fftw/fftw-3.1.1/tests/pgprof.out

Batch System LSF

Commands

<code>bsub</code>	job submission
<code>bjobs</code>	display status of (default: running or pending) jobs (use <code>-l <id></code> for more detailed information)
<code>/usr/local/bin/lsf_info</code>	display all running and pending jobs
<code>bkill <id></code>	cancel job with given ID
<code>bqueues</code>	display queue information (use <code>-l <queue></code> for more detailed information)
<code>bhosts</code>	display status of compute nodes

Submission examples

- Single CPU job (<10 min) submit: `bsub -W 10 -o output.txt -e error.txt ./a.out`
- MPI-parallel job submit (using default queue): `bsub -n 16 -a mvapich mpirun.lsf ./a.out`
- OpenMP job submit: `bsub -q small -n 2 -x ./a.out` (exclusively usage of a node)
- MPI-parallel, use only 1 CPU per node: `bsub -q single_cpu -n 2 -x ./a.out`
- MPI-parallel, use only 1 CPU per node: `bsub -q single_cpu -n 2 -x ./a.out`
- Job dependencies
`bsub -J "Job1" run1`
`bsub -w 'done(Job1)' run2`

Applications

- Abaqus 6.5.4
- Ansys
- LS-DYNA

- CFX 10.0
- Fluent 6.2

- Gaussian 3.0
- Gromacs
- NAMD 2.6

- MatLab 7.2



TECHNISCHE
UNIVERSITÄT
DRESDEN

Zentrum für Informationsdienste und Hochleistungsrechnen (ZIH)

FFT-Berechnungen auf phobos

Daniel Hackenberg – daniel.hackenberg@zih.tu-dresden.de



FFT-Bibliotheken auf phobos

- Intel MKL:
 - include files in `/opt/cluster/intel/mkl/mkl80/include`
 - libraries in `/opt/cluster/intel/mkl/mkl80/lib/em64t`
 - alternativ: 'module load mkl/default' und dann Umgebungsvariablen `$MKL_INC` und `$MKL_LIB` nutzen
- AMD ACML:
 - include files in `/opt/cluster/acml/2.7.0/include`
 - libraries in `/opt/cluster/acml/2.7.0/lib`
 - alternativ: 'module load acml/default' und dann Umgebungsvariablen `$ACML_INC` und `$ACML_LIB` nutzen
- FFTW (Fastest Fourier Transform in the West):
 - freie Software (GPL), erhältlich unter fftw.org
 - auf phobos nicht verfügbar, muss manuell installiert werden

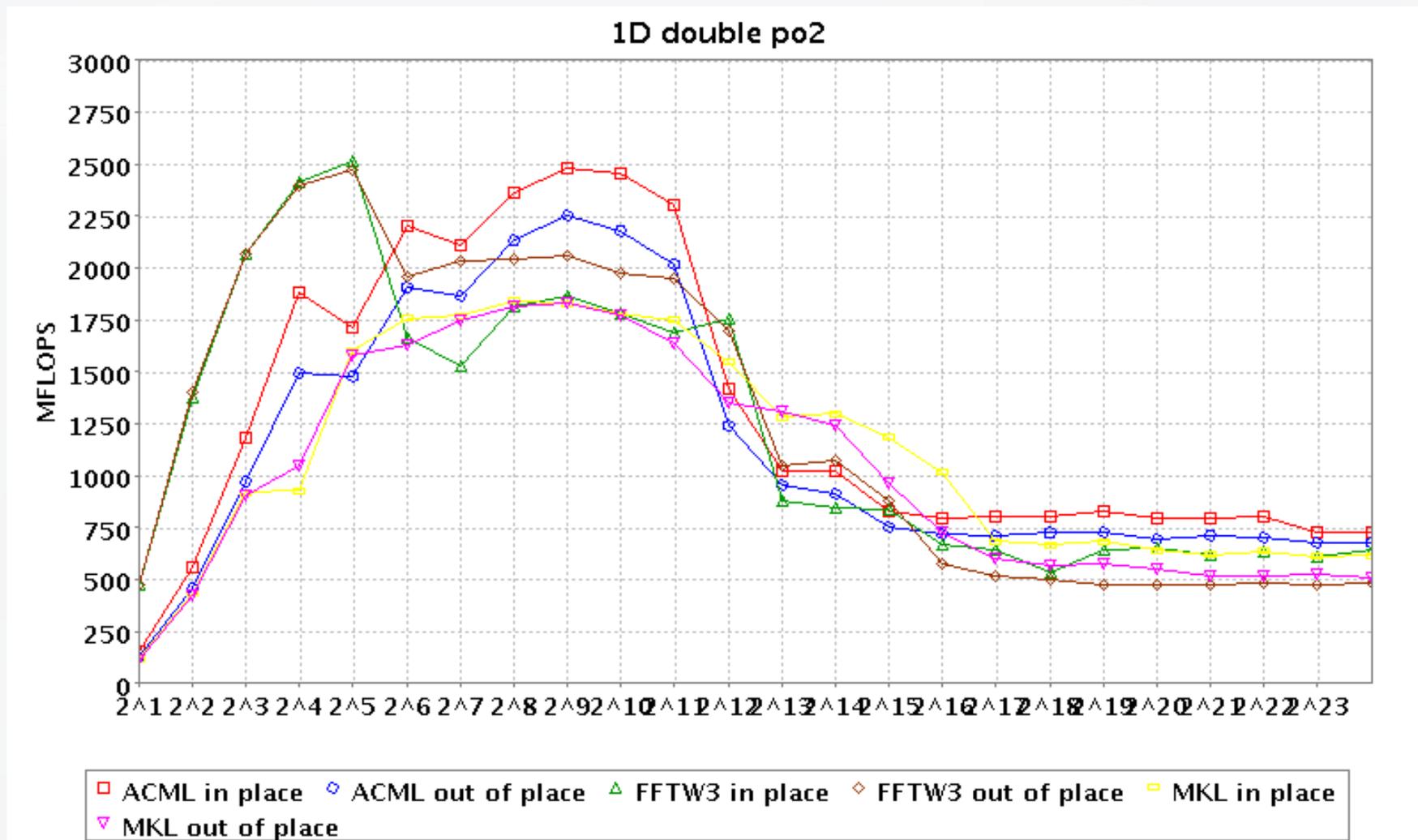
FFTW Installation

- Download unter <http://www.fftw.org/download.html> (derzeit Version 3.1.1)
- `./configure CC="icc" --prefix="/home/username/fftw3.1.1"`
 - für single precision: `--enable-float`
 - für long-double precision: `--enable-long-double`
- `make, make install`

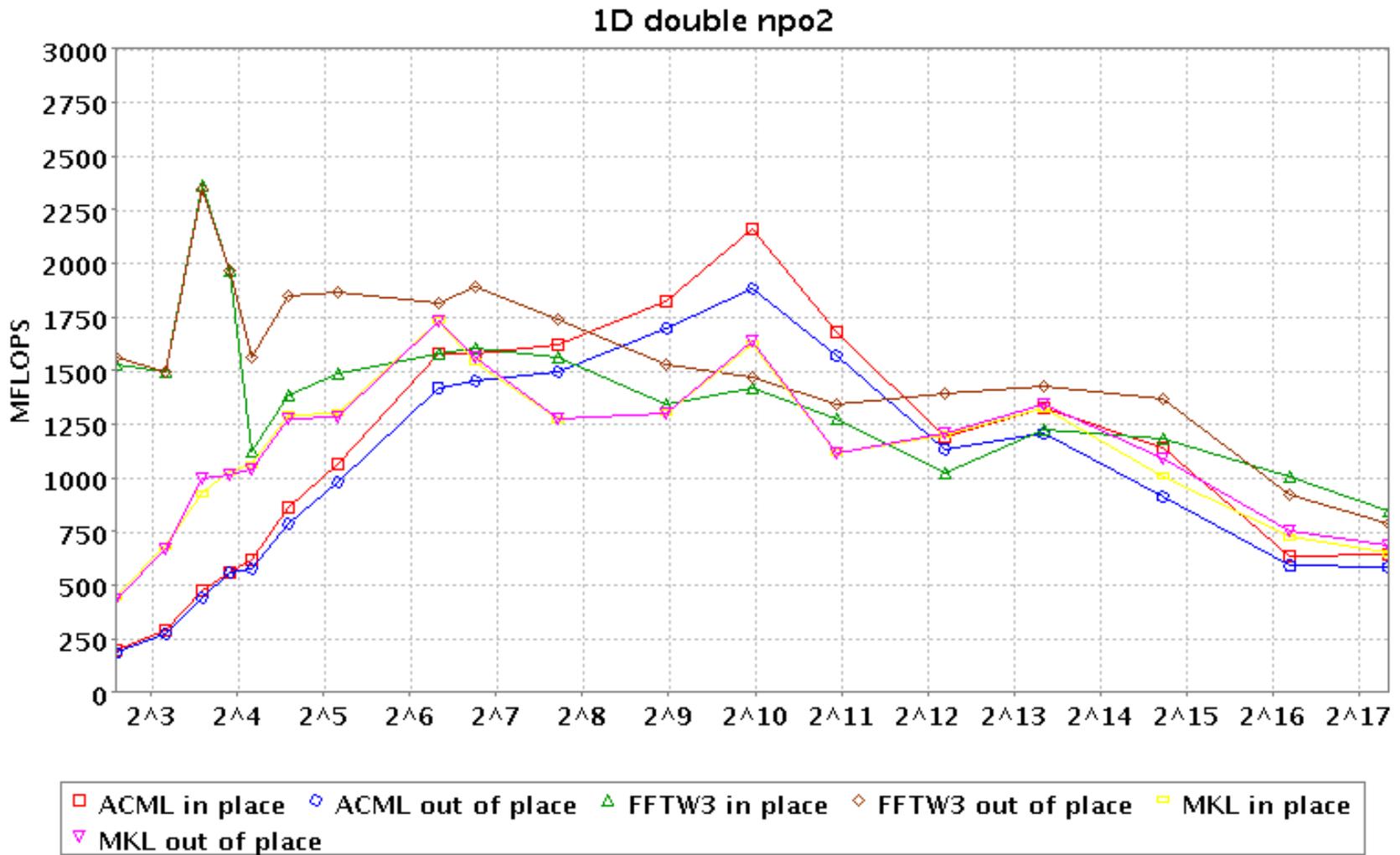
Dokumentation

- ACML: /opt/cluster/acml/2.7.0/Doc/acml.pdf (S. 19 ff.)
- MKL: /opt/cluster/intel/mkl/mkl80/doc/mklman.pdf (S. 11-1 ff.)
- FFTW: http://www.fftw.org/fftw3_doc/

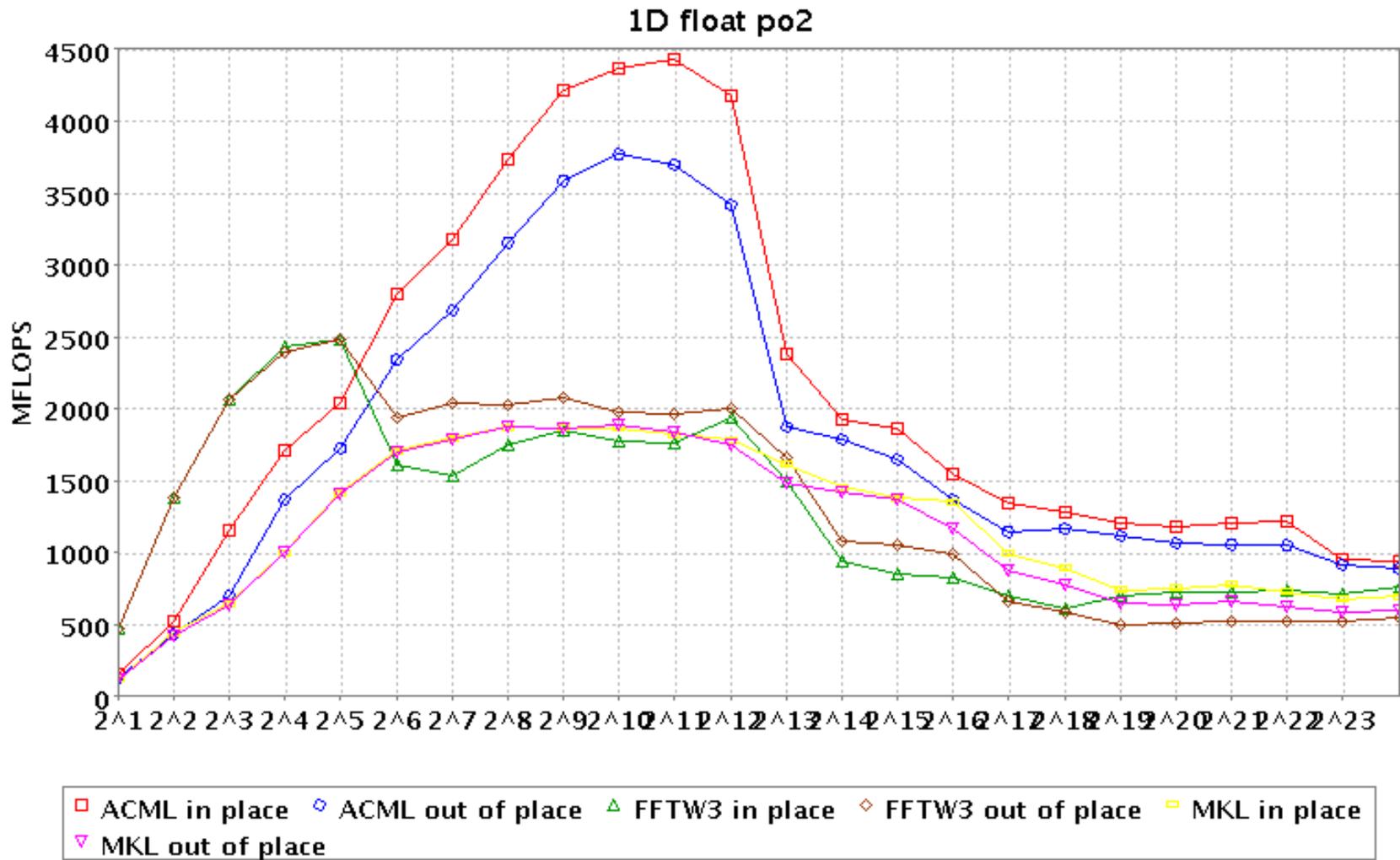
FFT: 1D, double precision, powers of two



FFT: 1D, double precision, non-powers of two



FFT: 1D, single precision, powers of two



Hinweise

- bei sehr vielen gleichartigen FFTs:
 - FFTW: Performance bei Plangenerierung mit “FFTW_MEASURE” statt “FFTW_ESTIMATE” testen (sucht besten Algorithmus durch Messung, Plangenerierung dauert entsprechend länger)
 - ACML: Initialisierung mit MODE=100, Effekt wie oben
- bei 2D/3D-FFTs: OpenMP-Parallelisierung mit MKL auf 2 CPUs möglich (OMP_NUM_THREADS=2)
- FFT ist für Zweierpotenzen am schnellsten, ggf. möglichst kleine Primfaktoren nutzen
- ACML bei Primzahlen zum Teil um mehr als Faktor 100 langsamer als MKL; bei nicht kontrollierbarer Problemgröße somit unbrauchbar!



TECHNISCHE
UNIVERSITÄT
DRESDEN

Center for Information Services and High Performance Computing (ZIH)

VAMPIR Performance Visualization

Software-Tuning on the Phobos PC-Farm
July 4th, 2006

Zellescher Weg 12

Willers-Bau A25 rechts

Tel. +49 351 - 463 - 35048

Holger Brunst (holger.brunst@tu-dresden.de)



Center for Information Services &
High Performance Computing

Vampir History

PARvis at Research Center Jülich

- Vampir at Research Center Jülich

<http://www.top500.org/reports/1995/vampir/vampir.html>

- Vampir at ZIH, TU Dresden

- was commercially available via Pallas GmbH, later via Intel
- Vampir (GUI) was developed by ZIH, Dresden
- Vampirtrace was developed by Pallas/Intel

- Successor: Vampir NG (next Generation)

- client/server architecture
- distributed storage, more scalable

Vampir Versions

Vampir+Vampirtrace 3.0

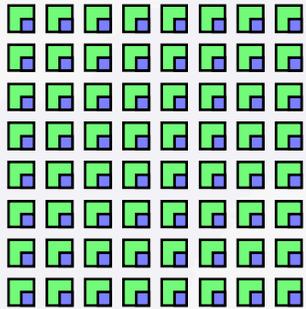
**Vampir+Vampirtrace 4.0 /
Intel Trace Analyzer 4.0+Intel Trace Collector 4.0**

**Vampir+Vampirtrace 5.0
Vampir Server+Client 1.5**

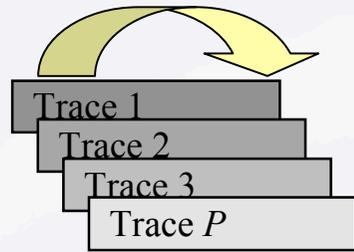
**Intel Trace Analyzer 6.0
Intel Trace Collector 6.0**

Vampir: technical components

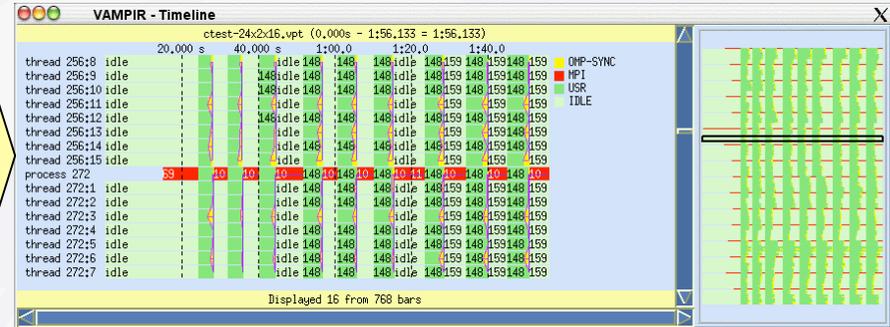
VampirTrace



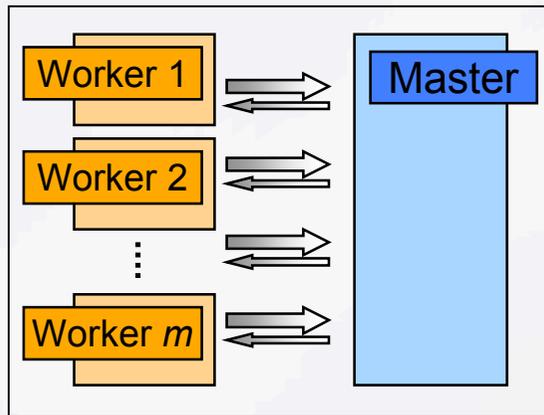
Vampir File (.otf/.vtf)



Vampir

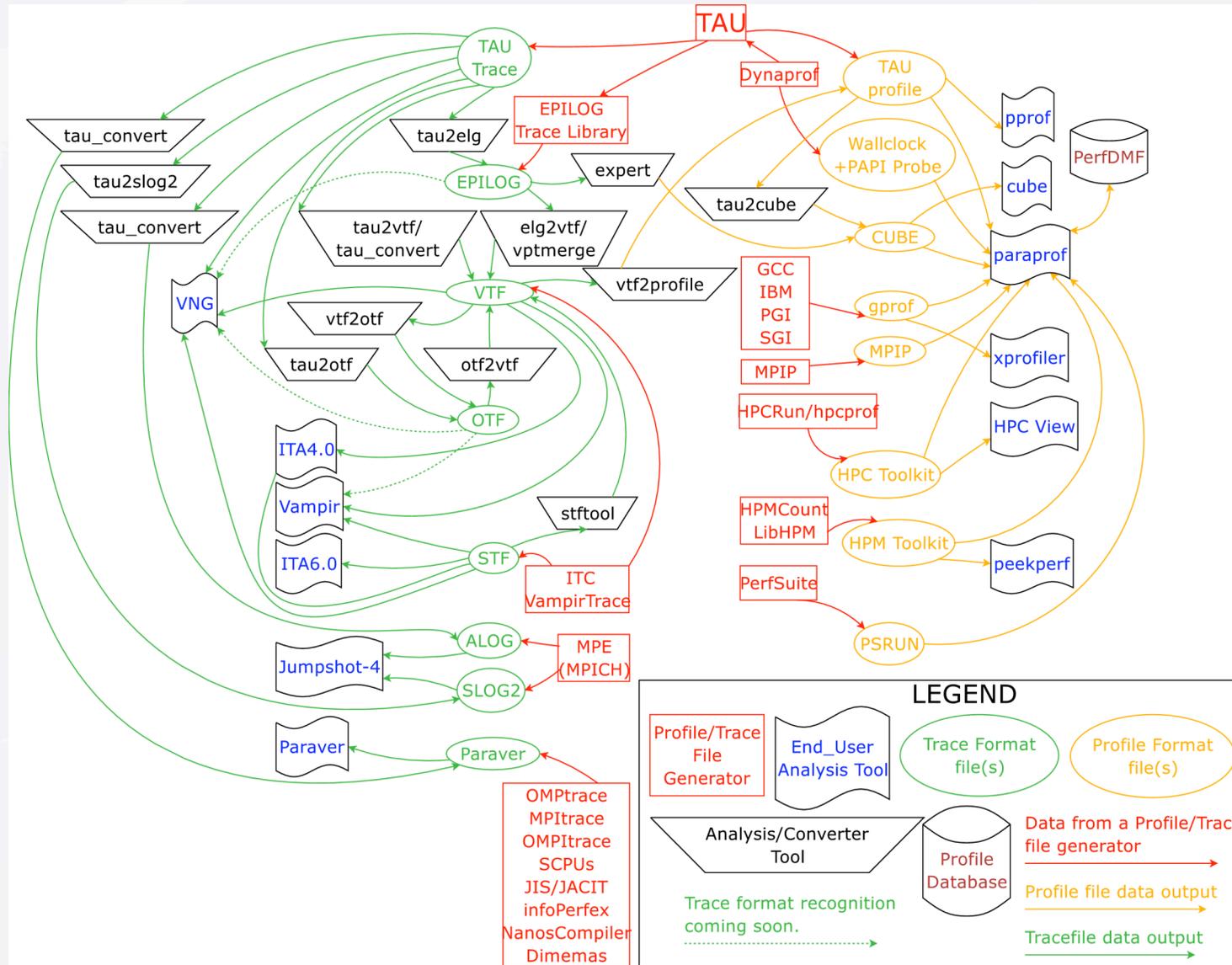


Vampir Server



1. Trace generator
2. Trace file(s)
3. Alignment and conversion tools
4. Viewer and analyzer
5. Parallel server engine

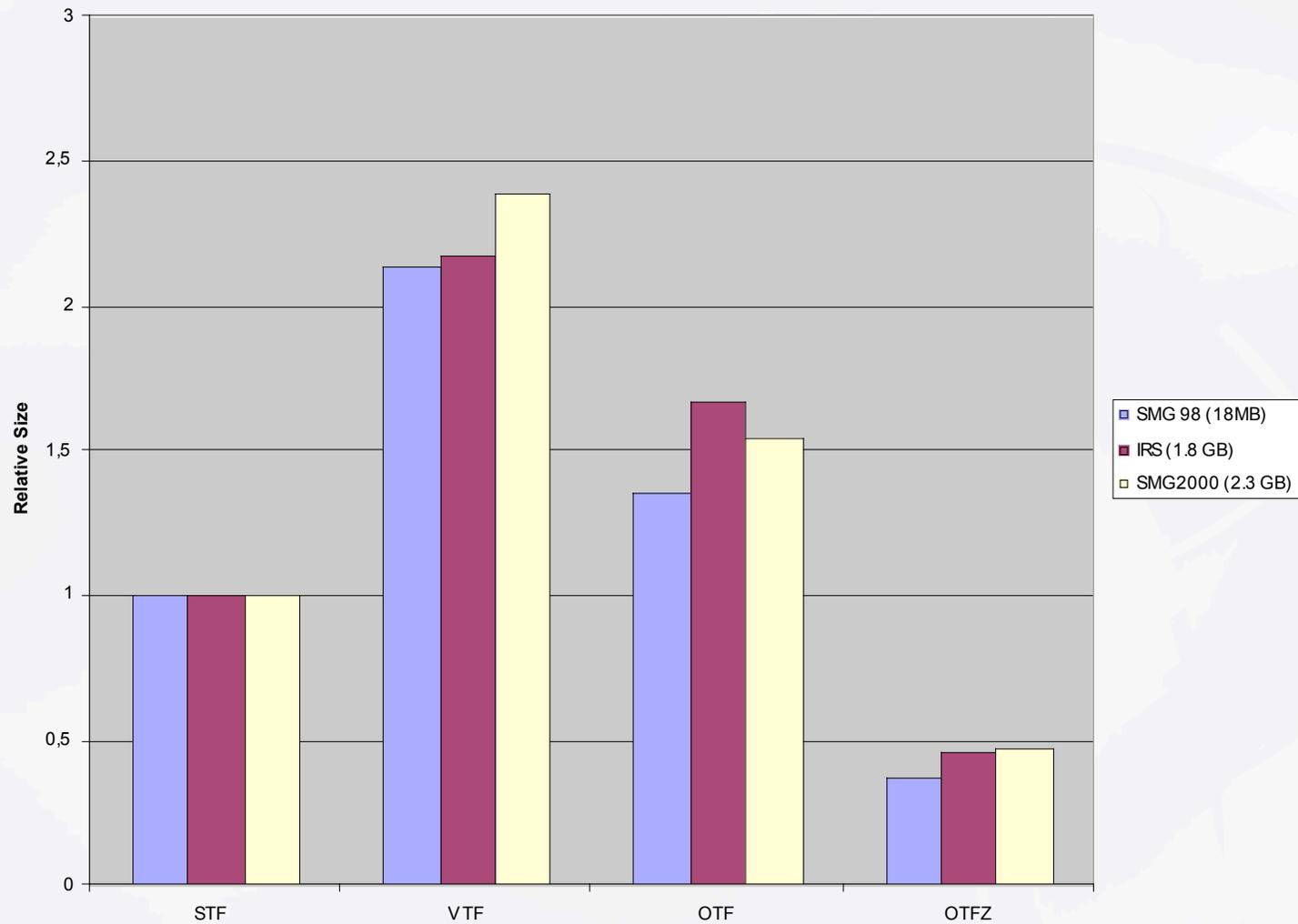
Many trace formats to choose from ...



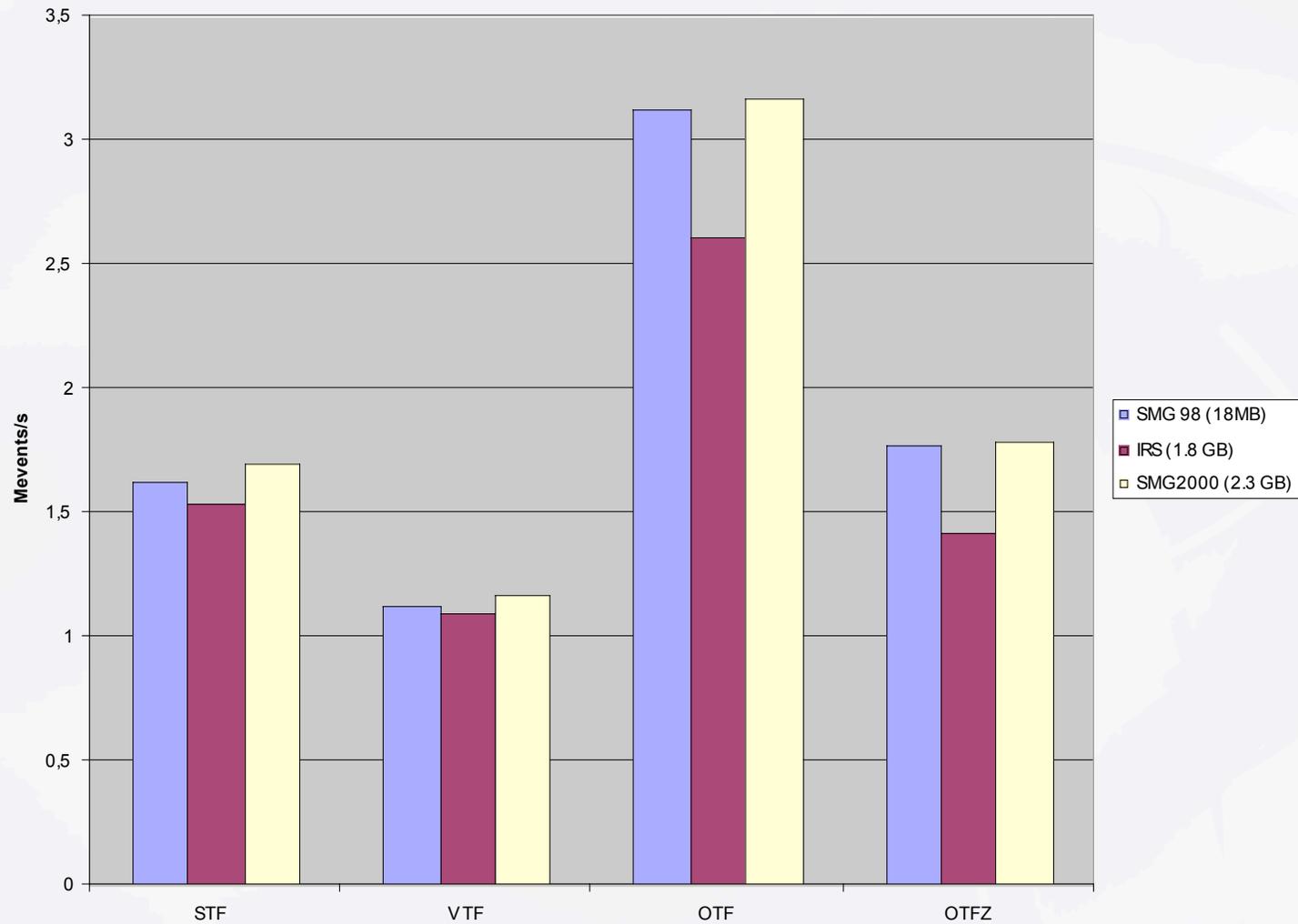
OTF Features

- Fast and efficient sequential and parallel access
- Platform independent
- Selective access to
 - Processes
 - Time intervals
- API / Interfaces
 - High level interface for analysis tools
 - Read/write complete traces with multiple files
 - Supports filtering and parallel I/O
 - Low level interface for trace libraries

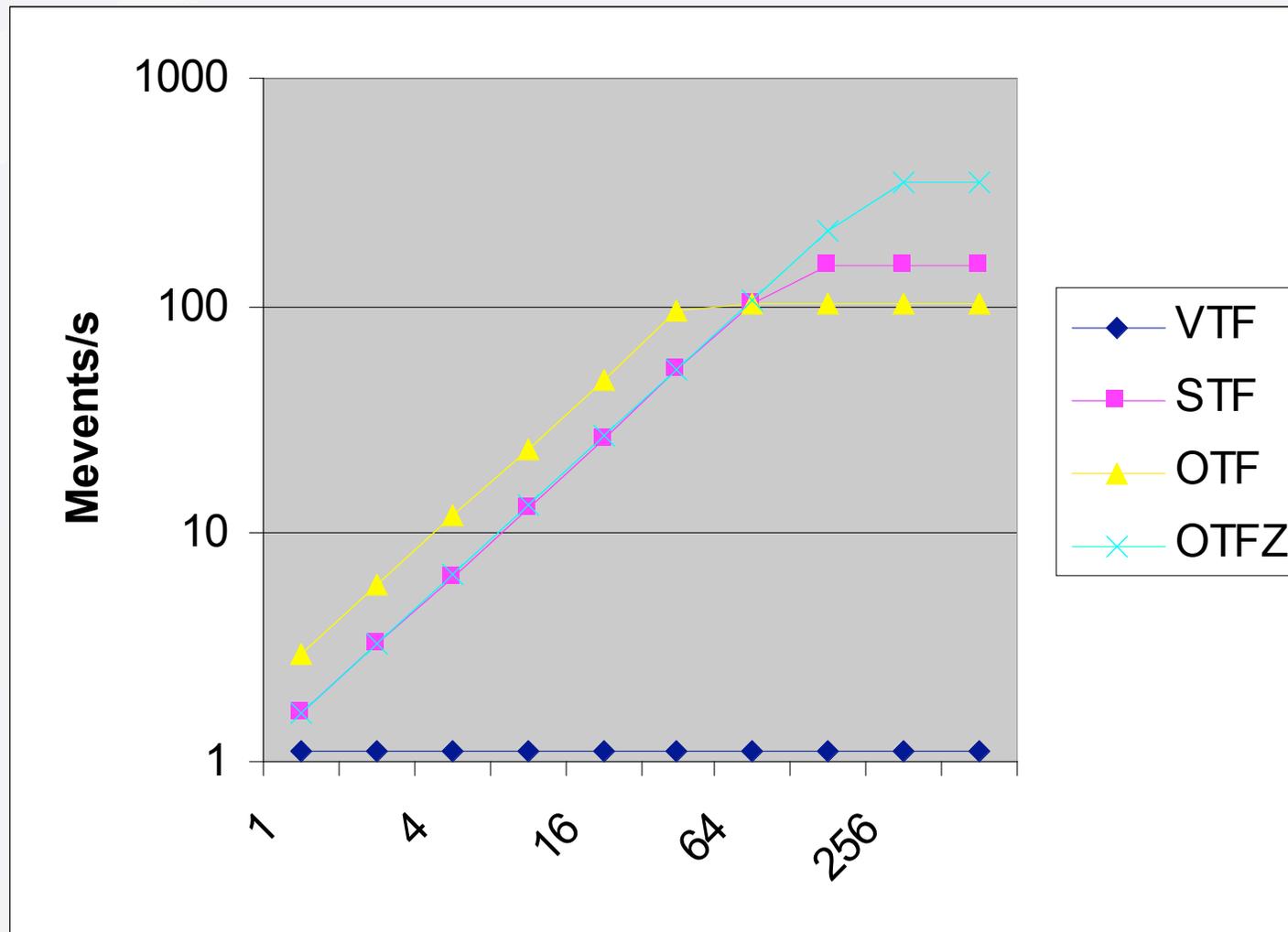
Relative File Size



Read Performance



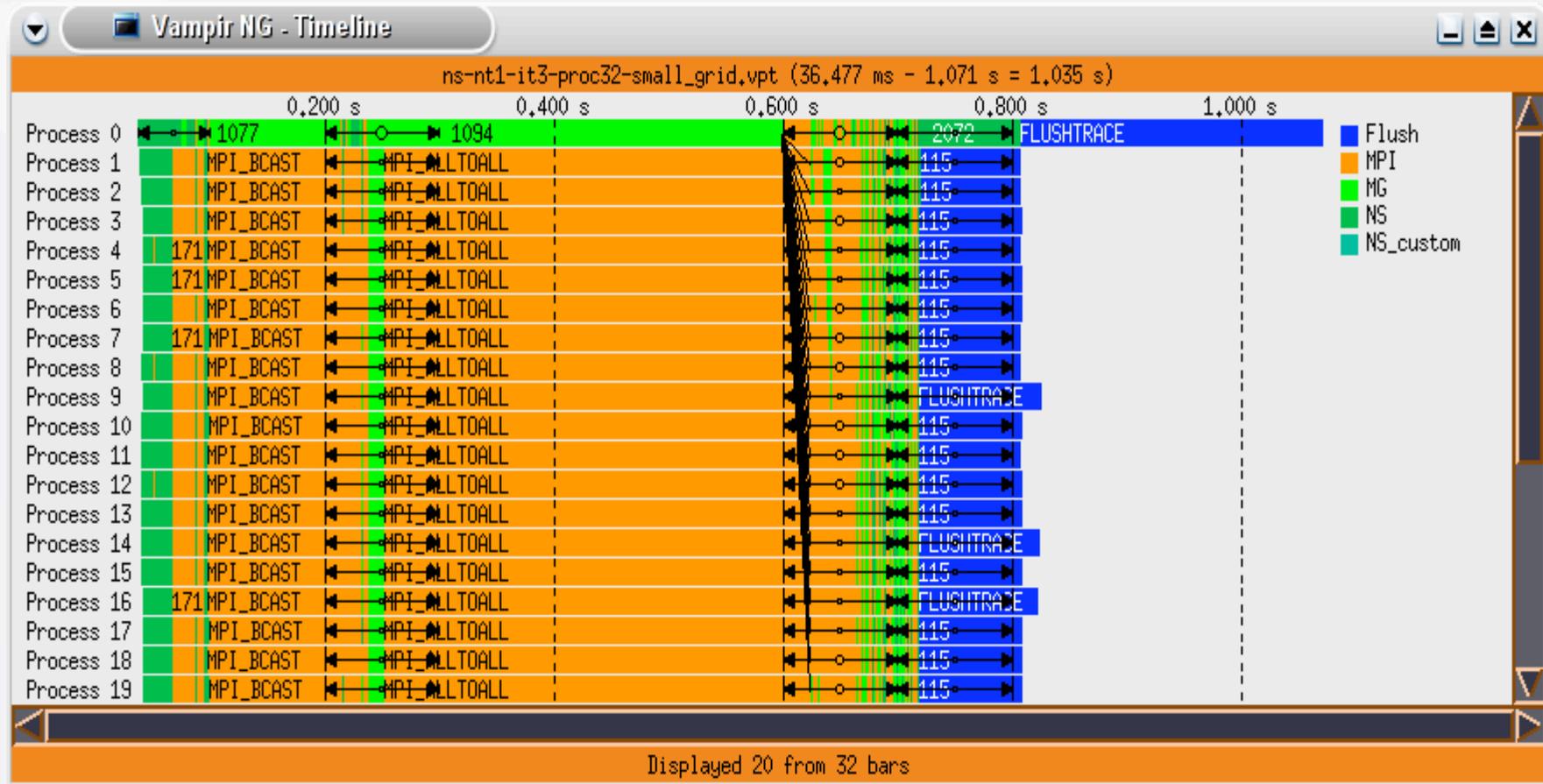
Performance Scalability



Vampir Displays

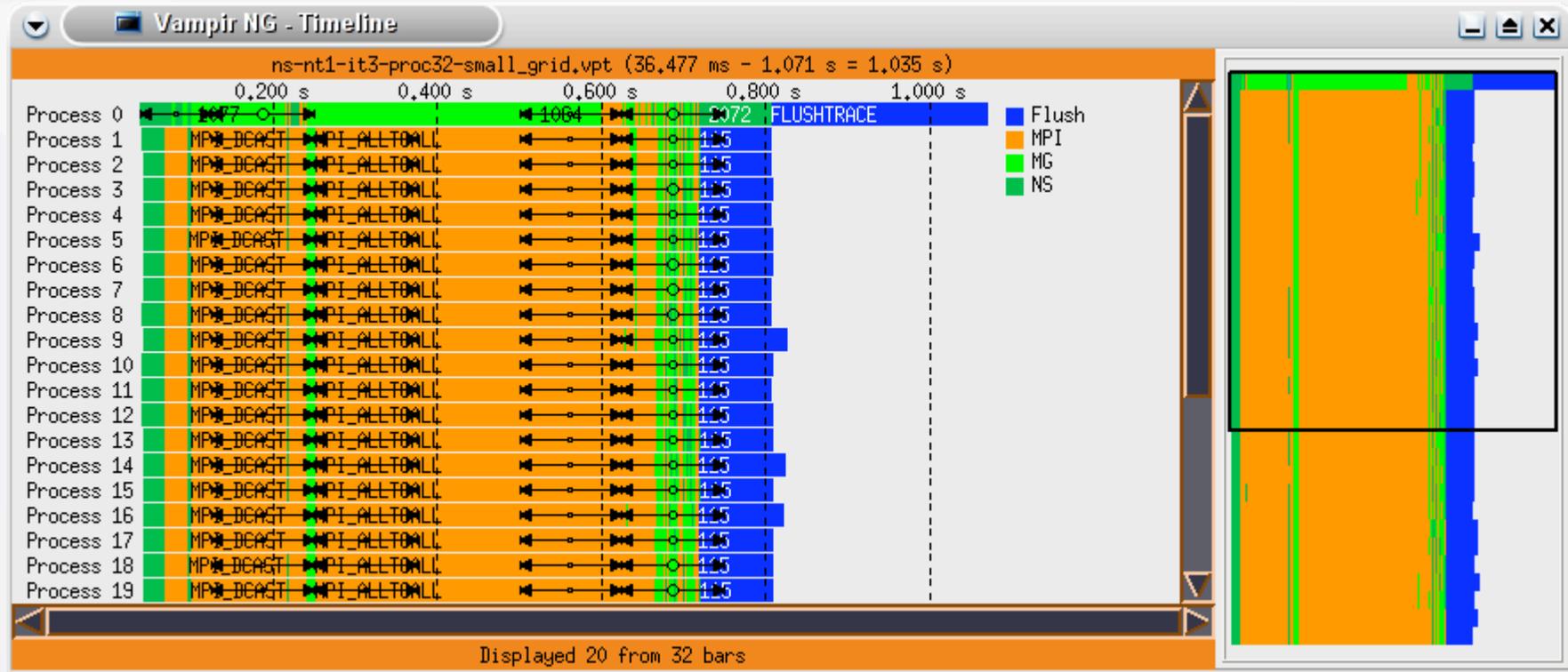
- Global Timeline
- Process Timeline
- Summary Chart
- Summary Timeline
- Counter Timeline
- Process Profile
- Call Tree
- Message Statistics
- Navigator
- Global Filters
- OpenMP

Vampir Displays - Global Timeline



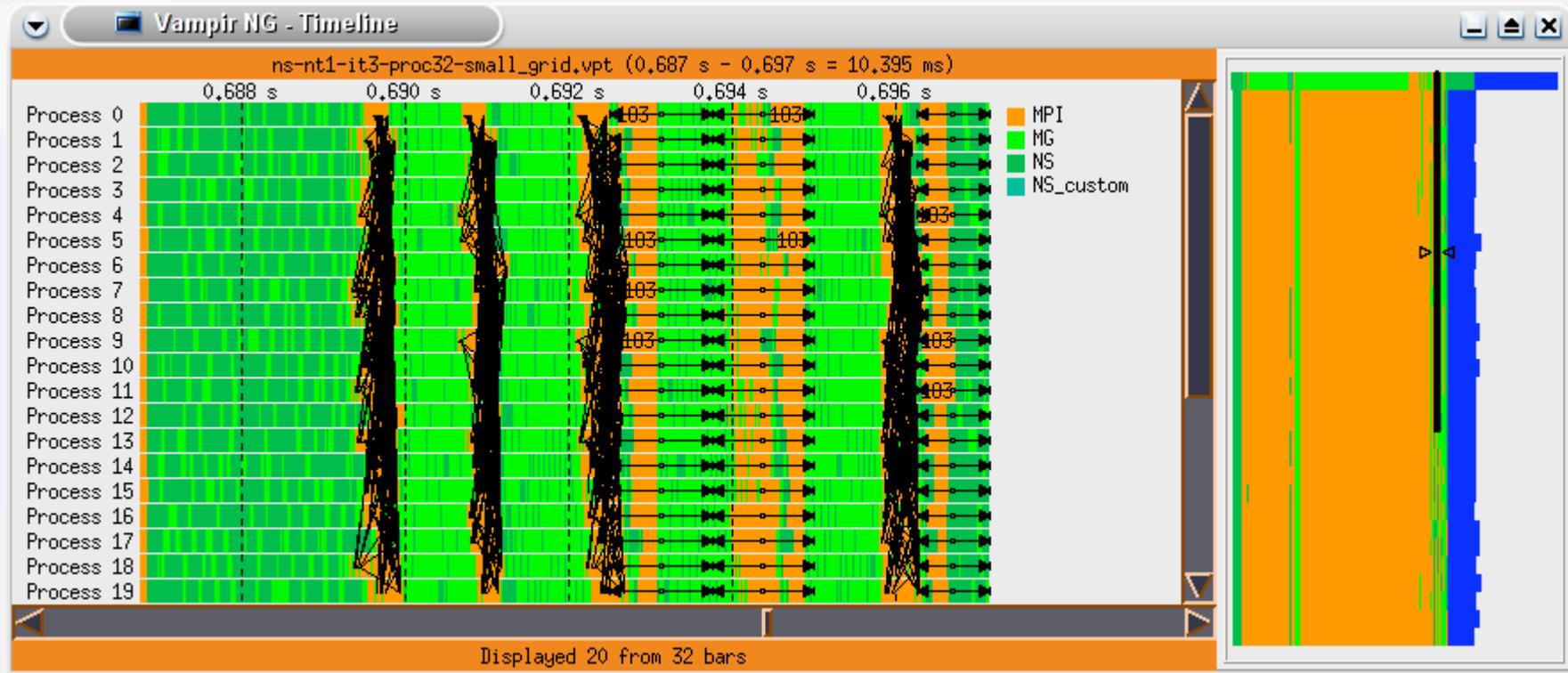
Global Timeline

Vampir Displays - Global Timeline with Thumbnail



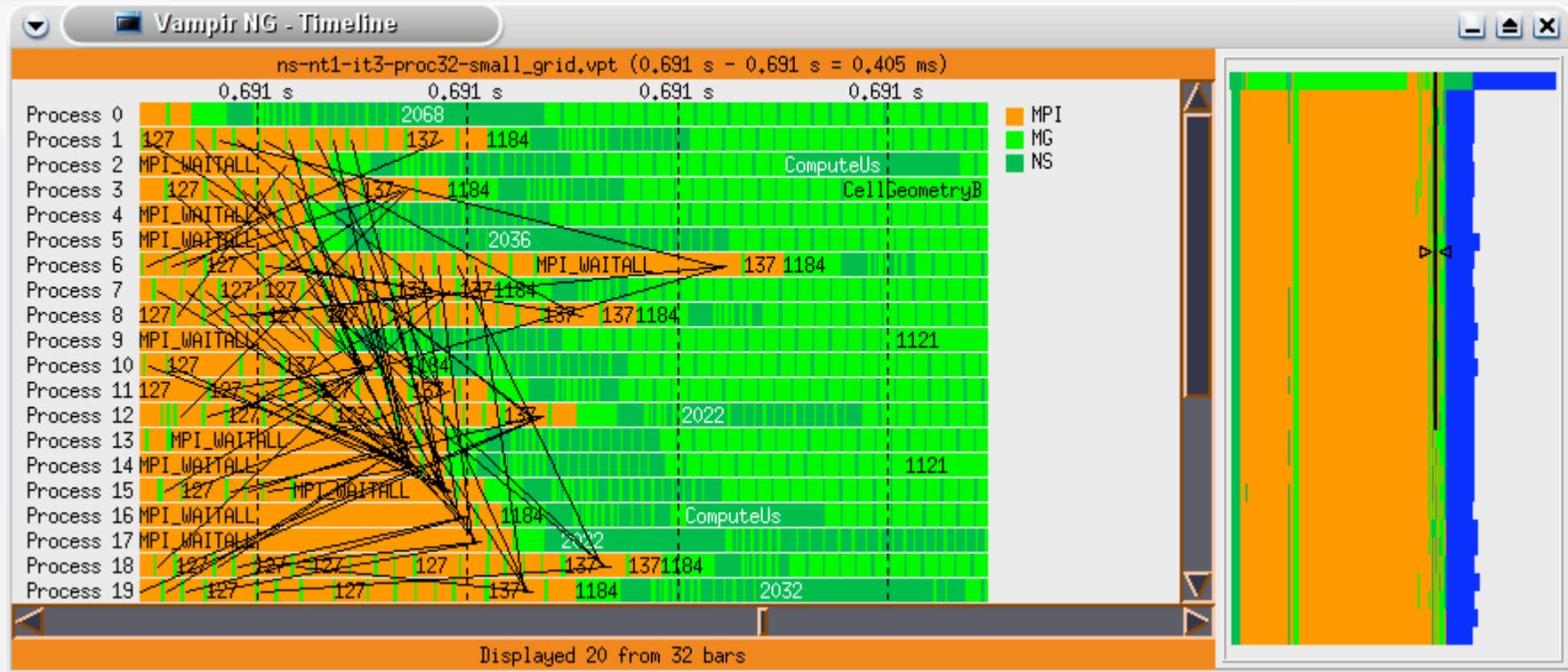
Global Timeline with Thumbnail View

Vampir Displays - Global Timeline with Thumbnail



Further Zoomed

Vampir Displays - Global Timeline with Thumbnail

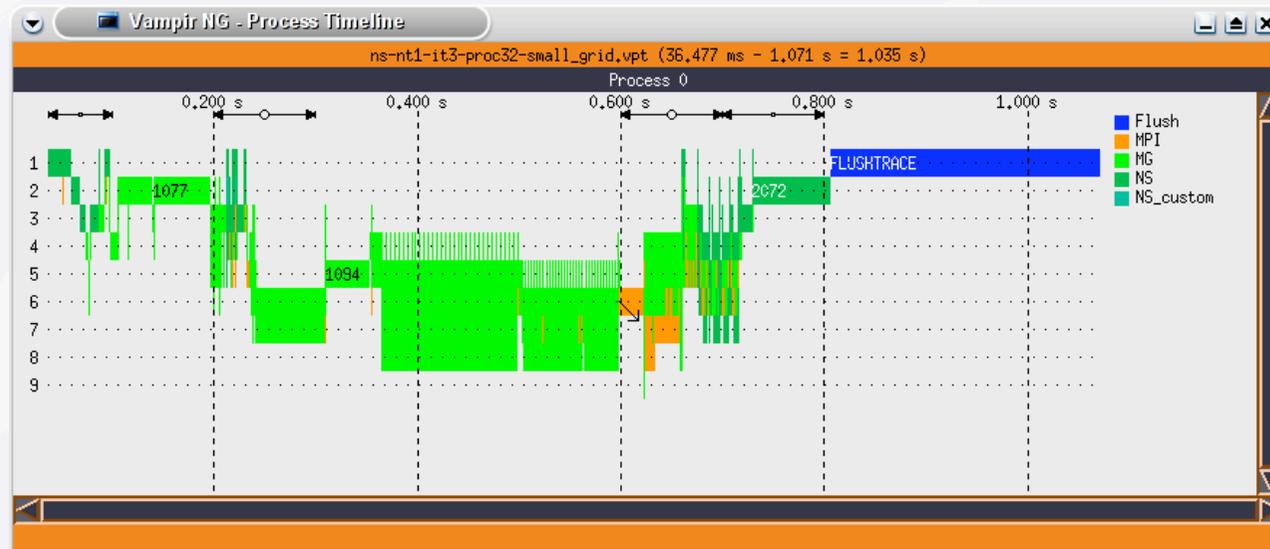


Even More Zoomed

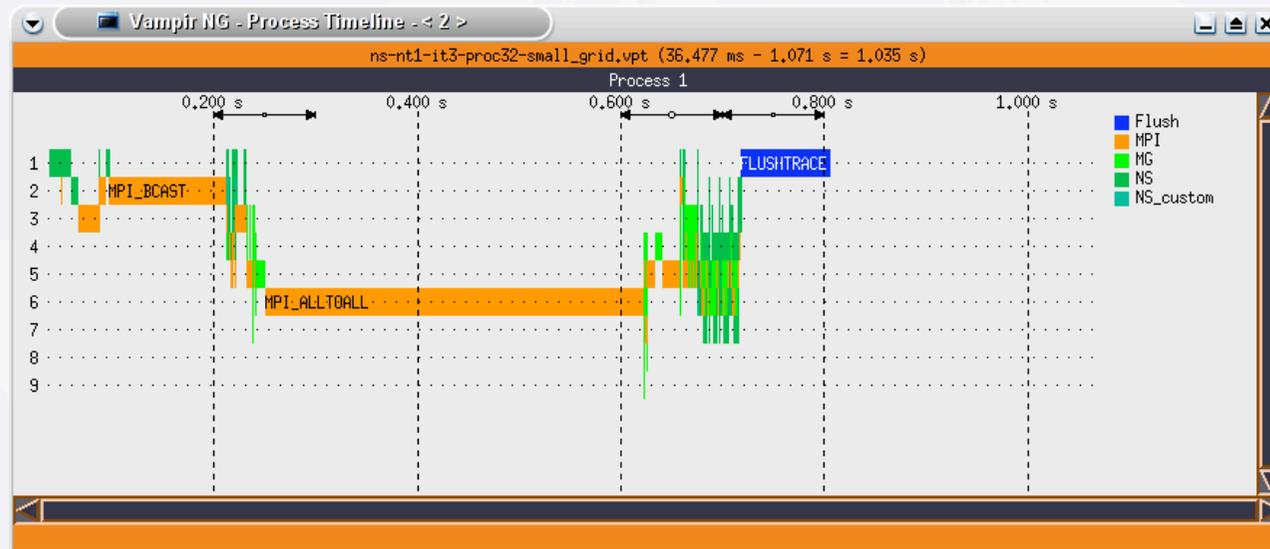
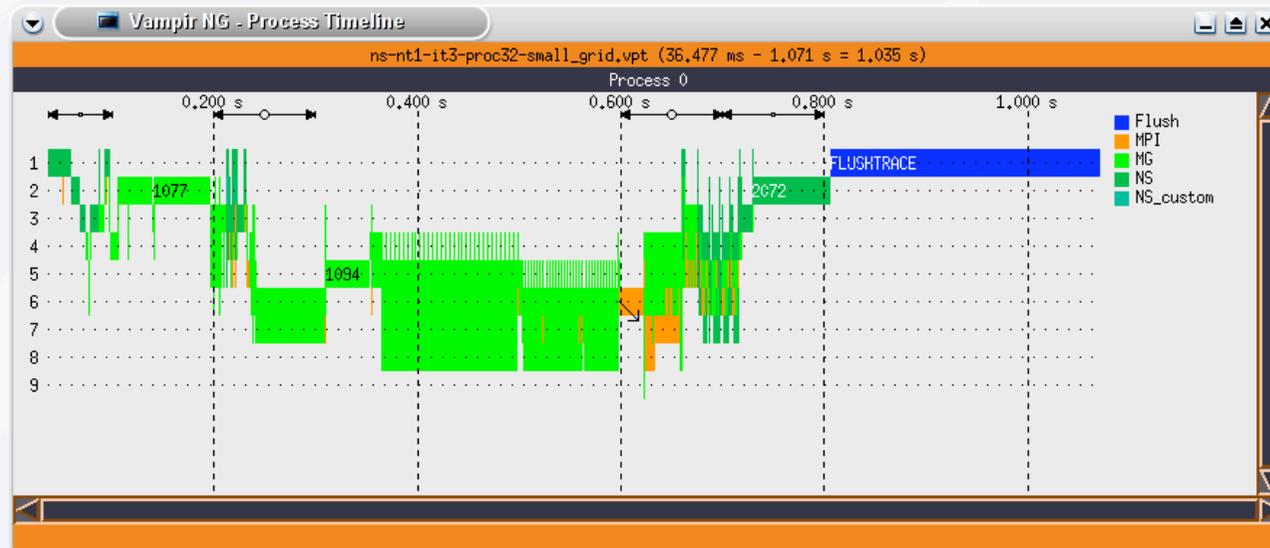
Vampir Displays - Global Timeline

- show temporal behavior for
 - processes
 - functions resp. function groups
 - messages
 - collective ops
 - I/O activities
- context information on click
- zoomable
- control actual time interval for other displays

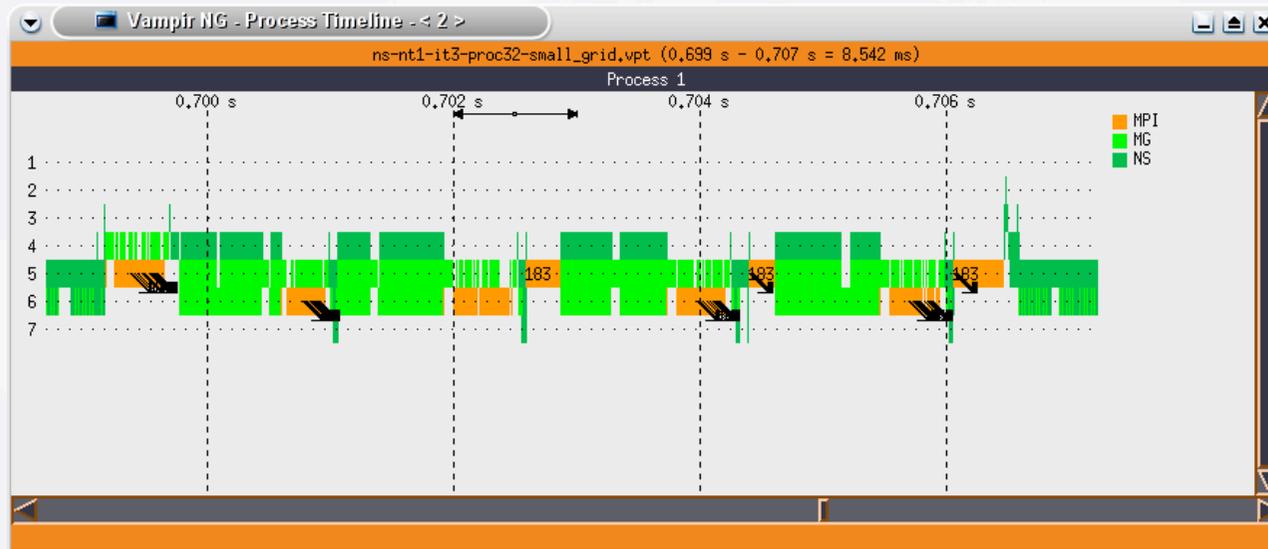
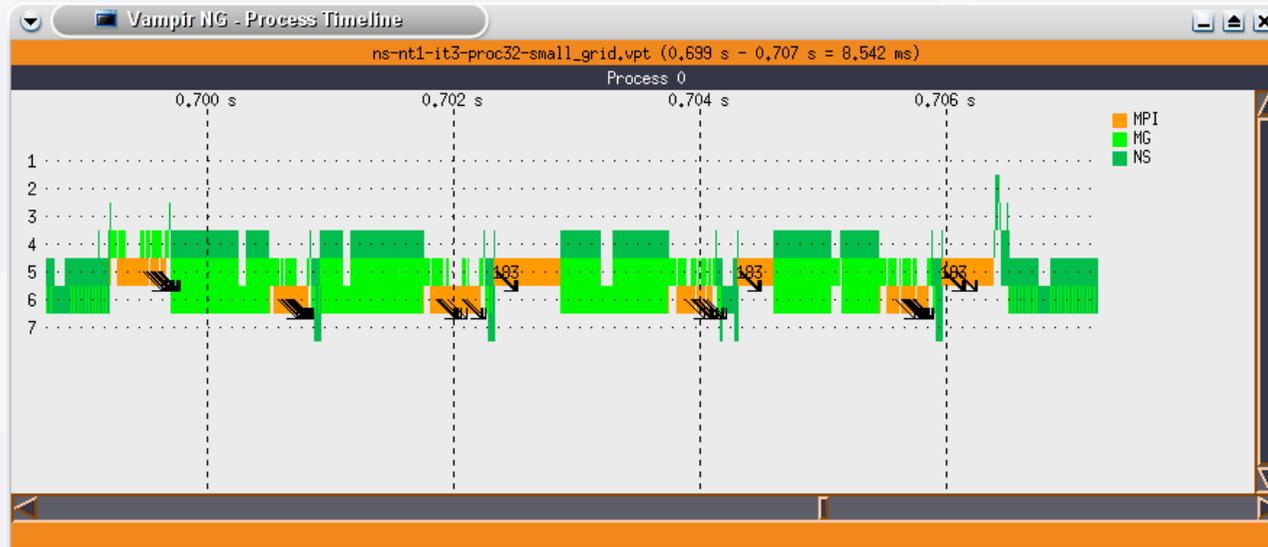
Vampir Displays - Process Timeline



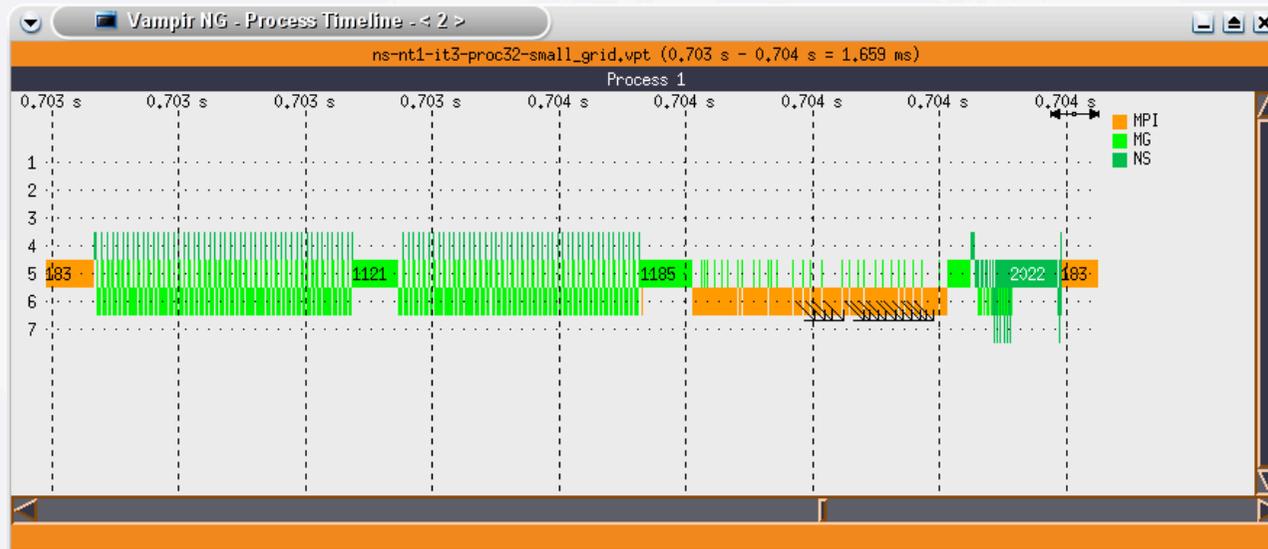
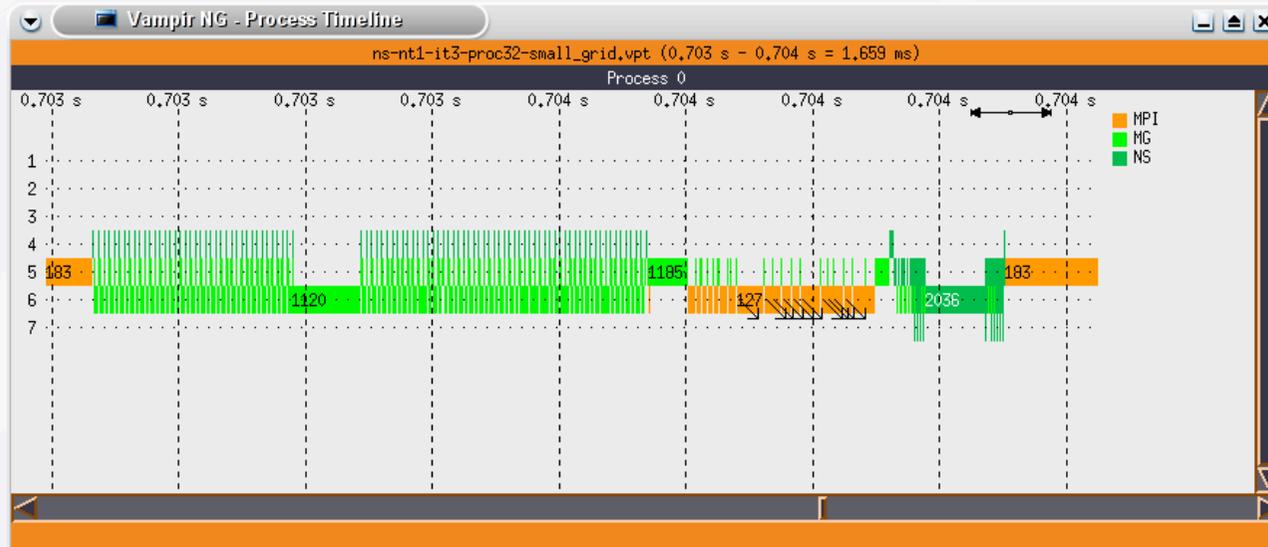
Vampir Displays – Aligned Process Timelines



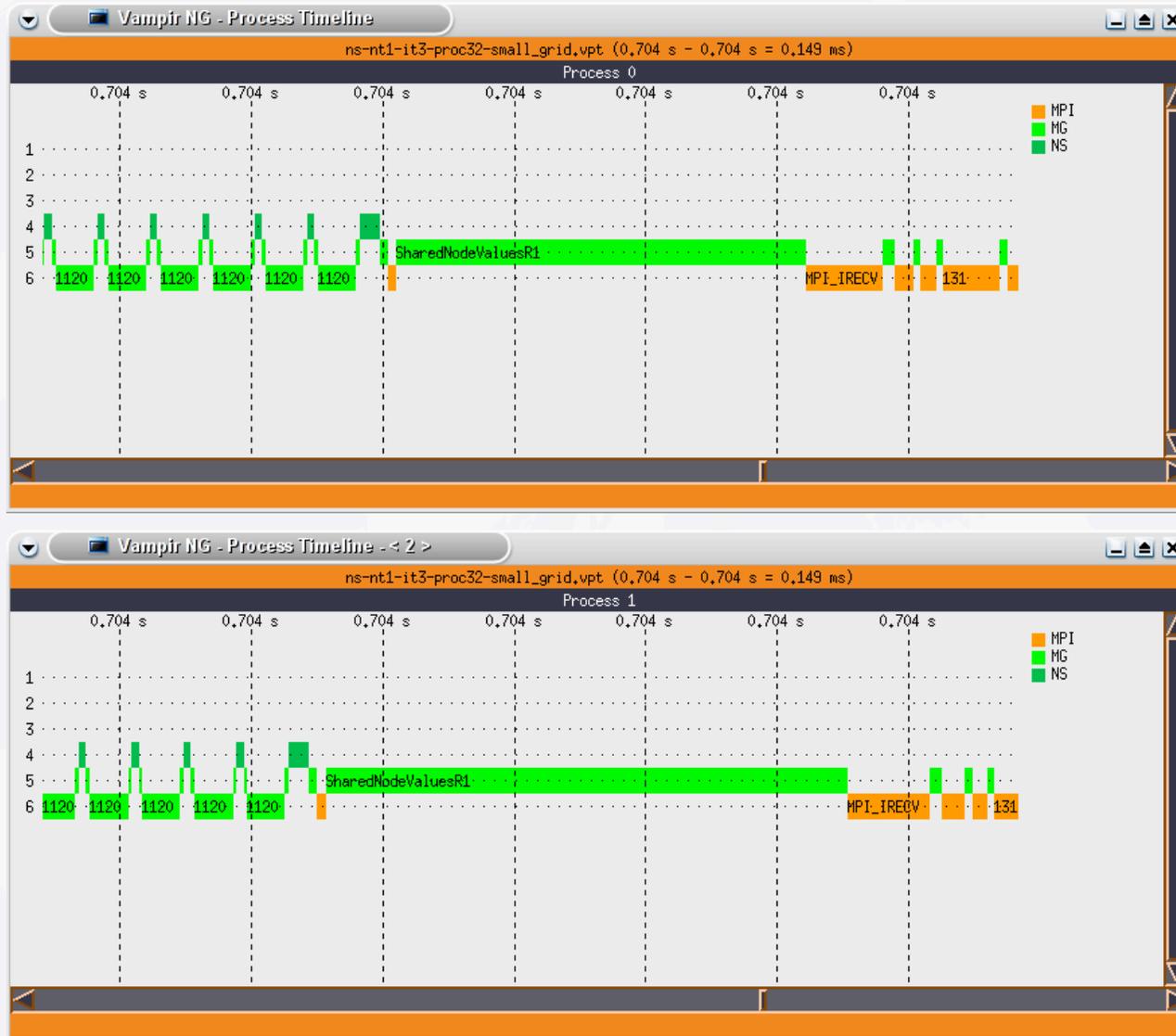
Vampir Displays - Process Timeline zoomed



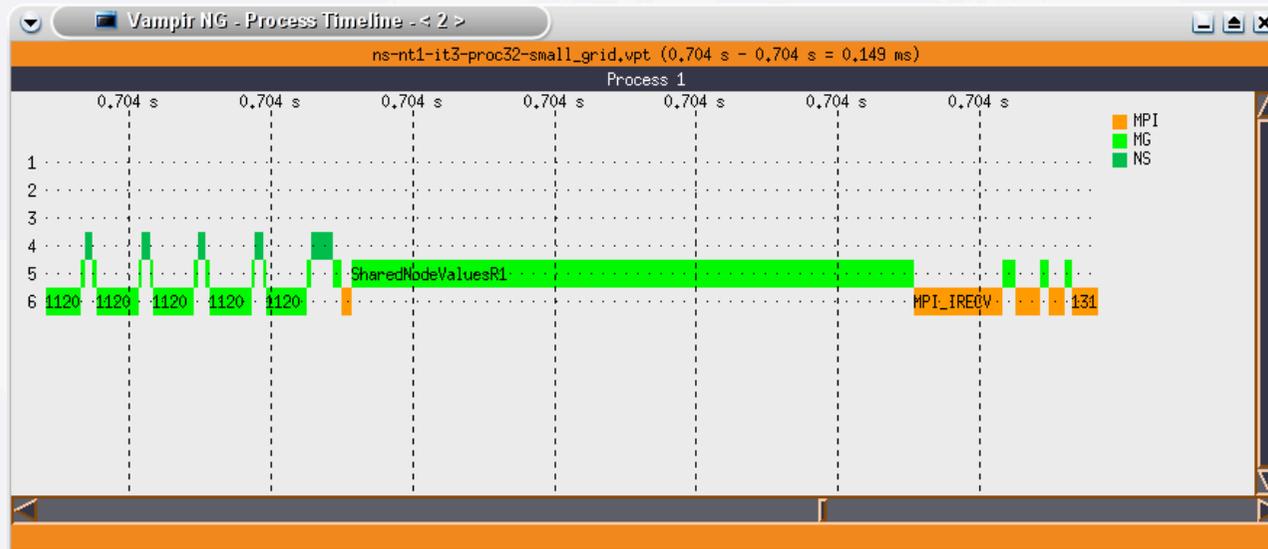
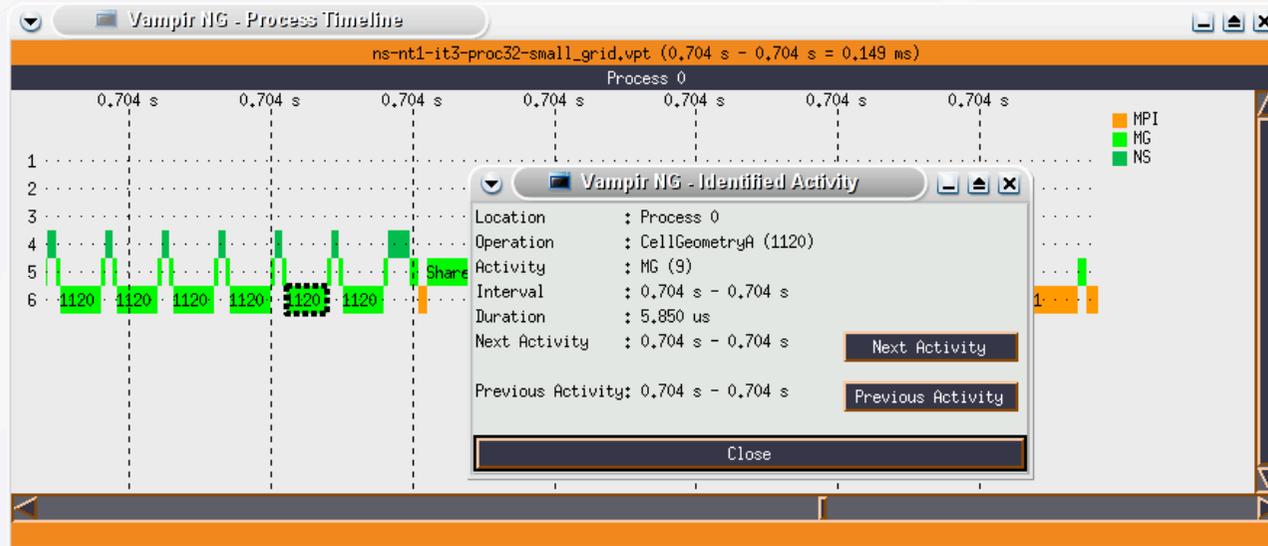
Vampir Displays - Process Timeline further zoomed



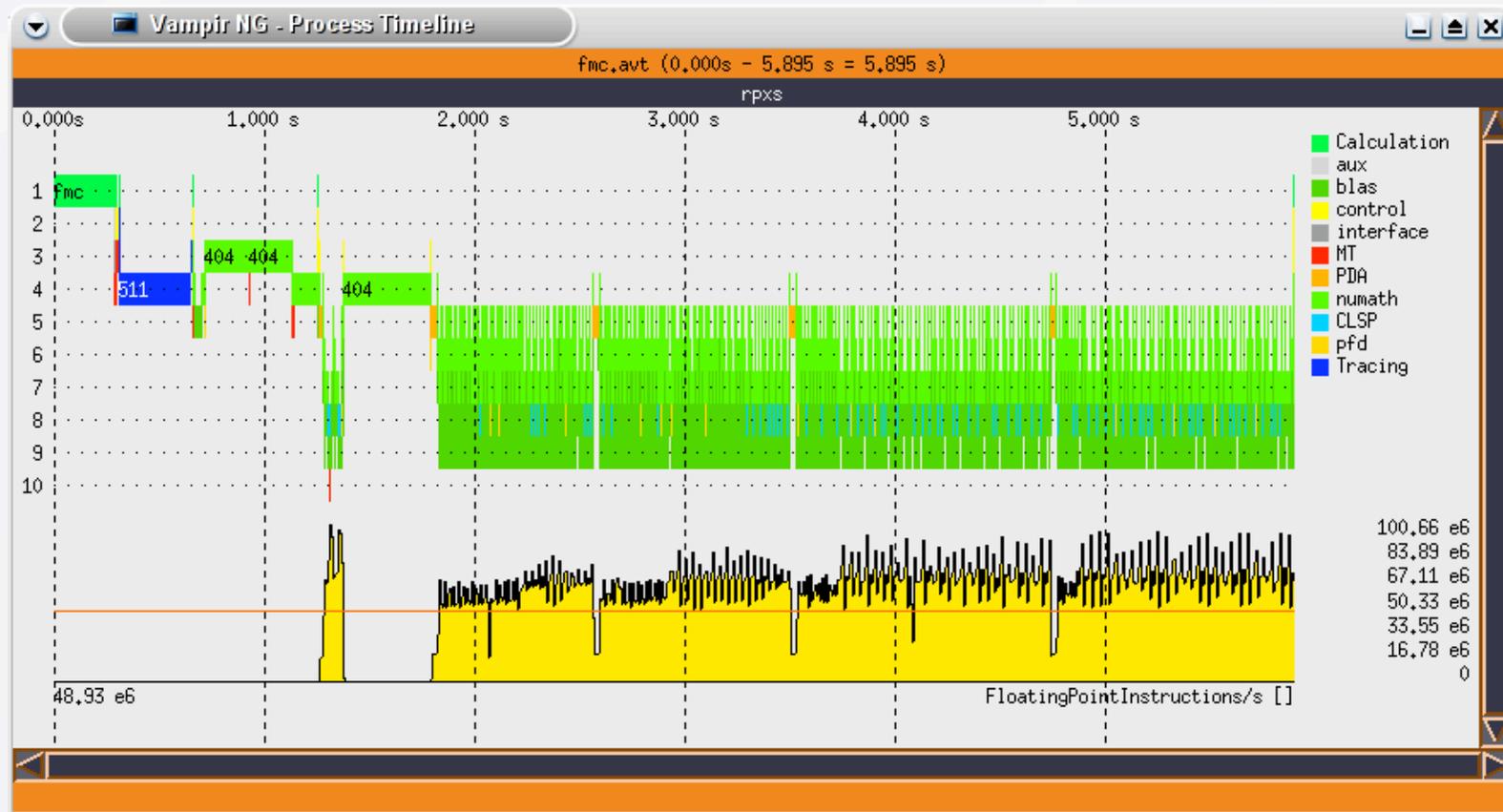
Vampir Displays - Process Timeline further zoomed II



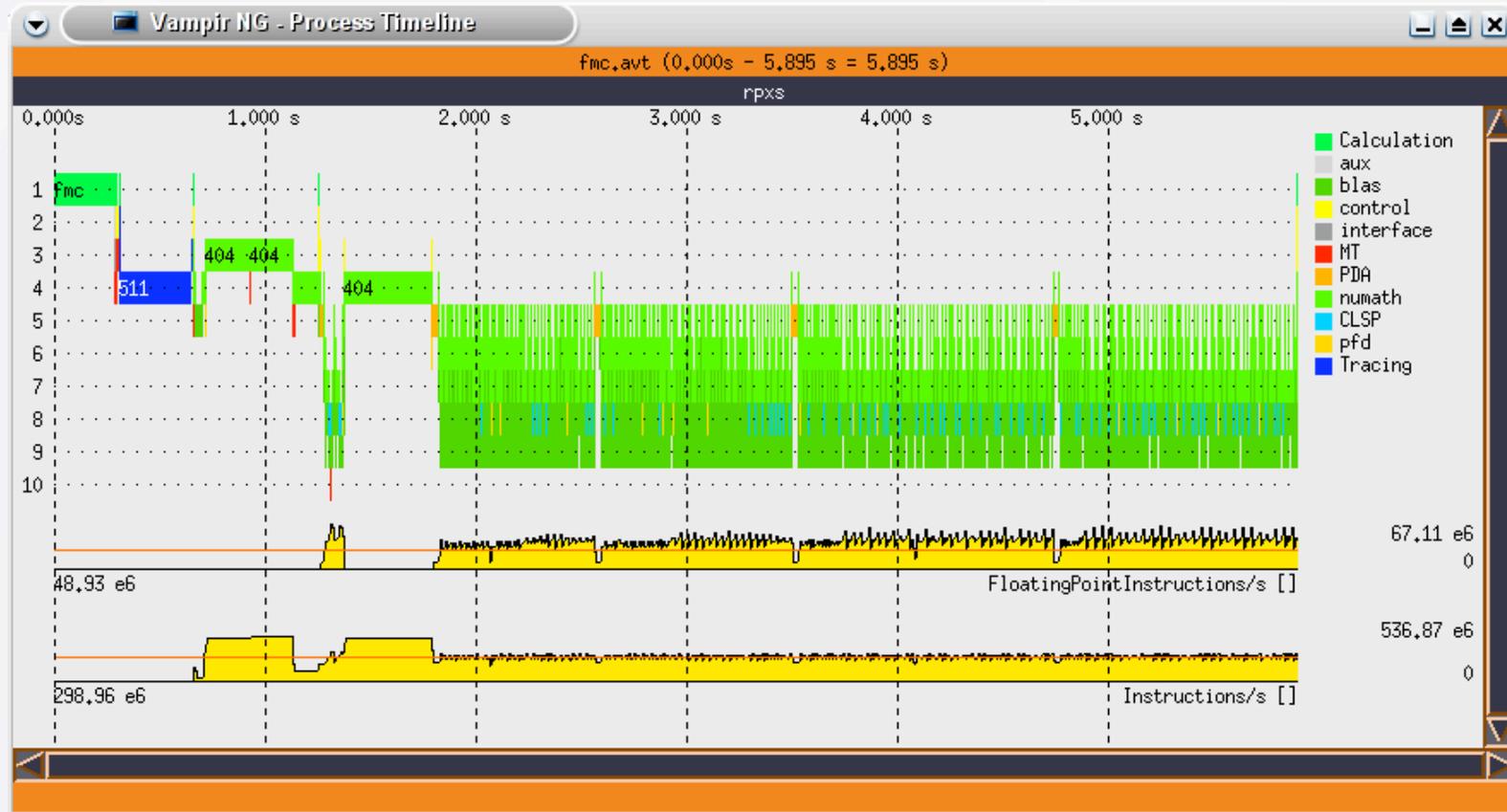
Vampir Displays - Process Timeline with context



Vampir Displays - Process Timeline with Flop/s



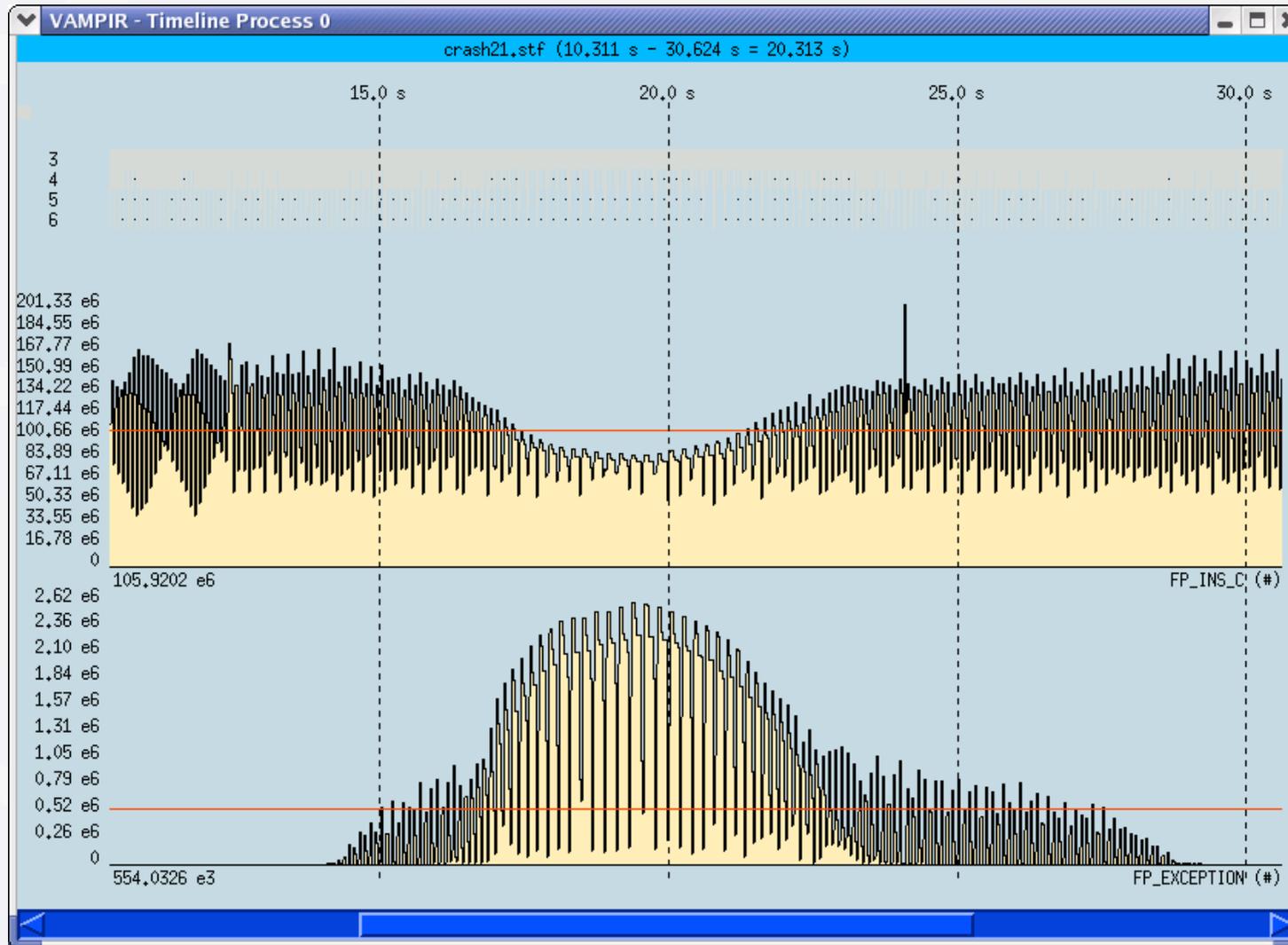
Vampir Displays - Process Timeline with two counters



Vampir Displays - Process Timeline zoomed



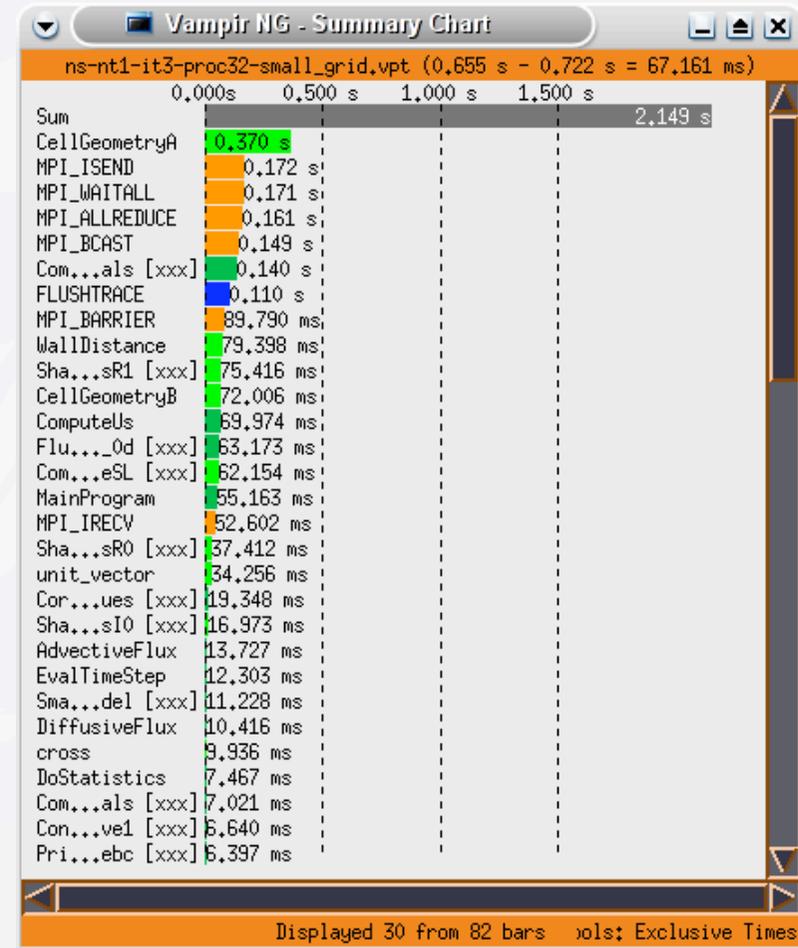
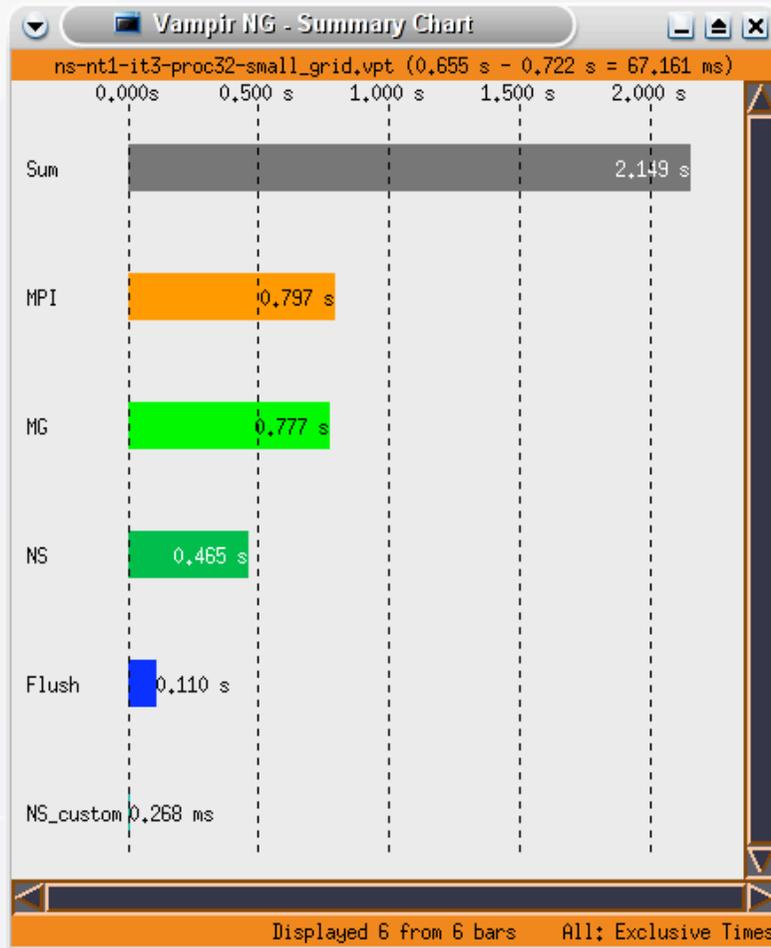
Vampir Displays - Counter Timeline



Vampir Displays - Process Timeline

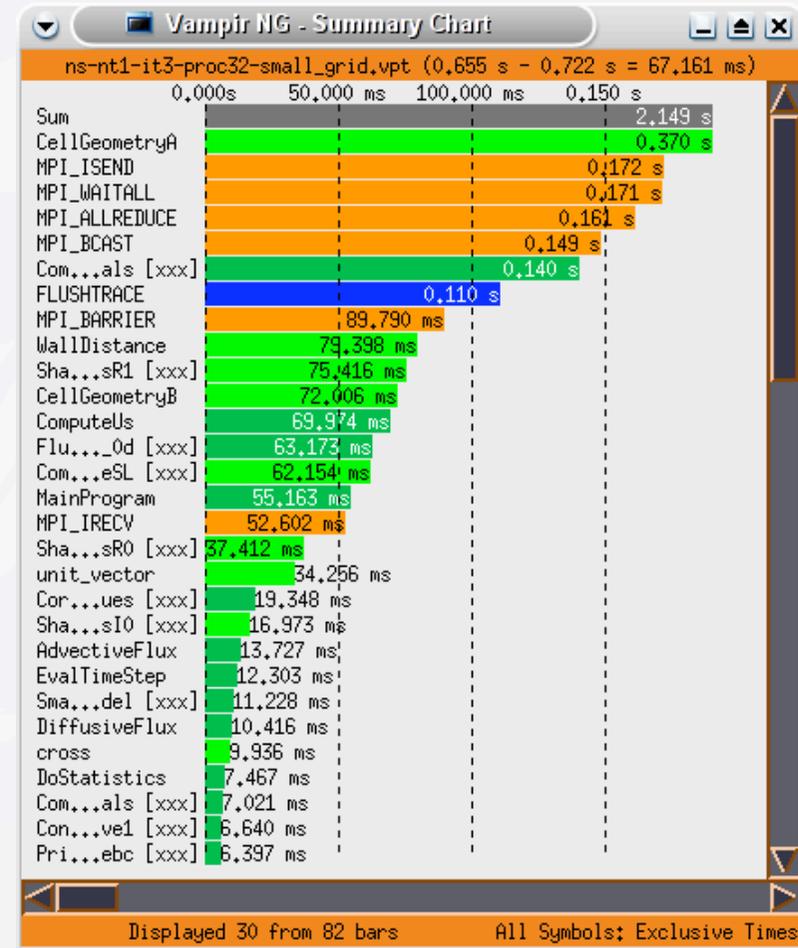
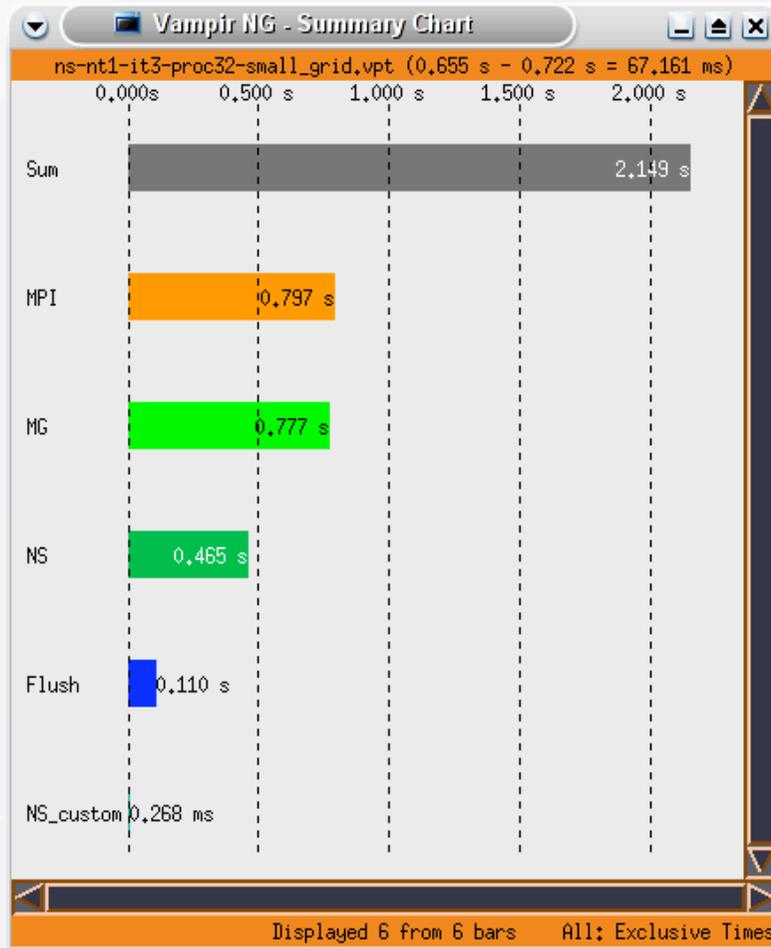
- timeline for single processes
- unfolded call stack in vertical dimension
- shows:
 - functions/function groups
 - messages, collective communication
 - I/O activities
 - performance counter values
- zoomable in time
- context display

Vampir Displays - Summary Chart



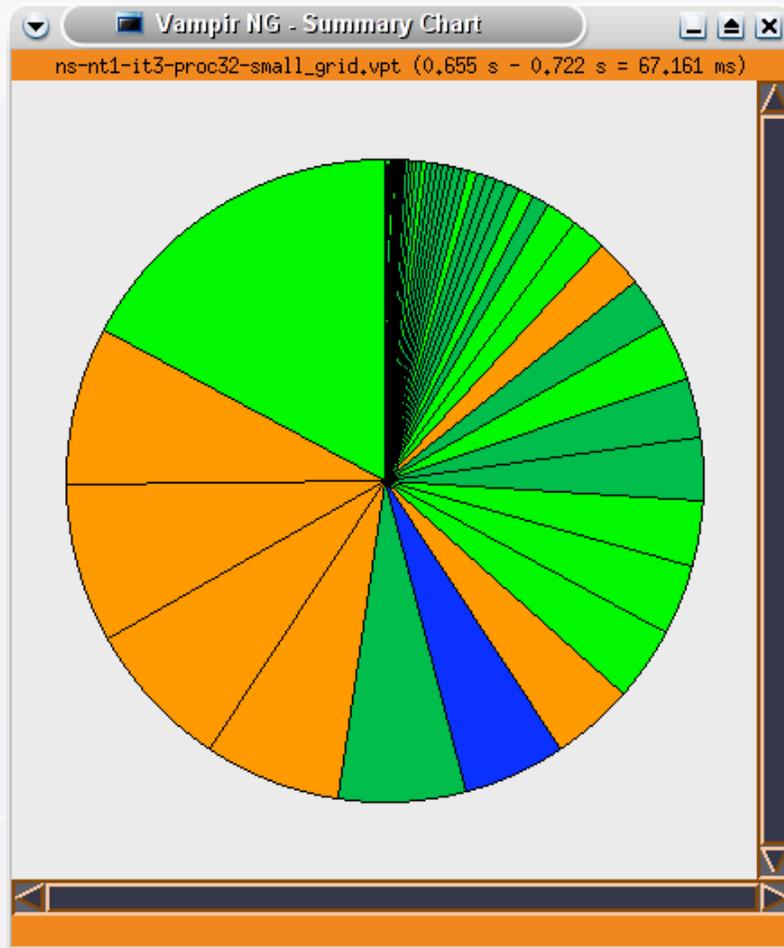
Grouped / Comprehensive Function Statistics

Vampir Displays - Summary Chart



Grouped / Comprehensive Function Statistics

Vampir Displays - Summary Chart



Name	Token	Value
CellGeometryA	[1120]	0,370 s
MPI_ISEND	[127]	0,172 s
MPI_WAITALL	[137]	0,171 s
MPI_ALLREDUCE	[183]	0,161 s
MPI_BCAST	[171]	0,149 s
Com...als [xxx]	[2018]	0,140 s
FLUSHTRACE	[115]	0,110 s
MPI_BARRIER	[170]	89,790 ms
WallDistance	[1192]	79,398 ms
Sha...sR1 [xxx]	[1185]	75,416 ms
CellGeometryB	[1121]	72,006 ms
ComputeUs	[2017]	69,974 ms
Flu..._0d [xxx]	[2032]	63,173 ms
Com...eSL [xxx]	[1130]	62,154 ms
MainProgram	[2065]	55,163 ms
MPI_Irecv	[131]	52,602 ms
Sha...sR0 [xxx]	[1184]	37,412 ms
unit_vector	[1212]	34,256 ms
Cor...ues [xxx]	[2022]	19,348 ms
Sha...sI0 [xxx]	[1183]	16,973 ms
AdvectiveFlux	[2008]	13,727 ms
EvalTimeStep	[2028]	12,303 ms
Sma...del [xxx]	[2082]	11,228 ms
DiffusiveFlux	[2023]	10,416 ms
cross	[1193]	9,936 ms

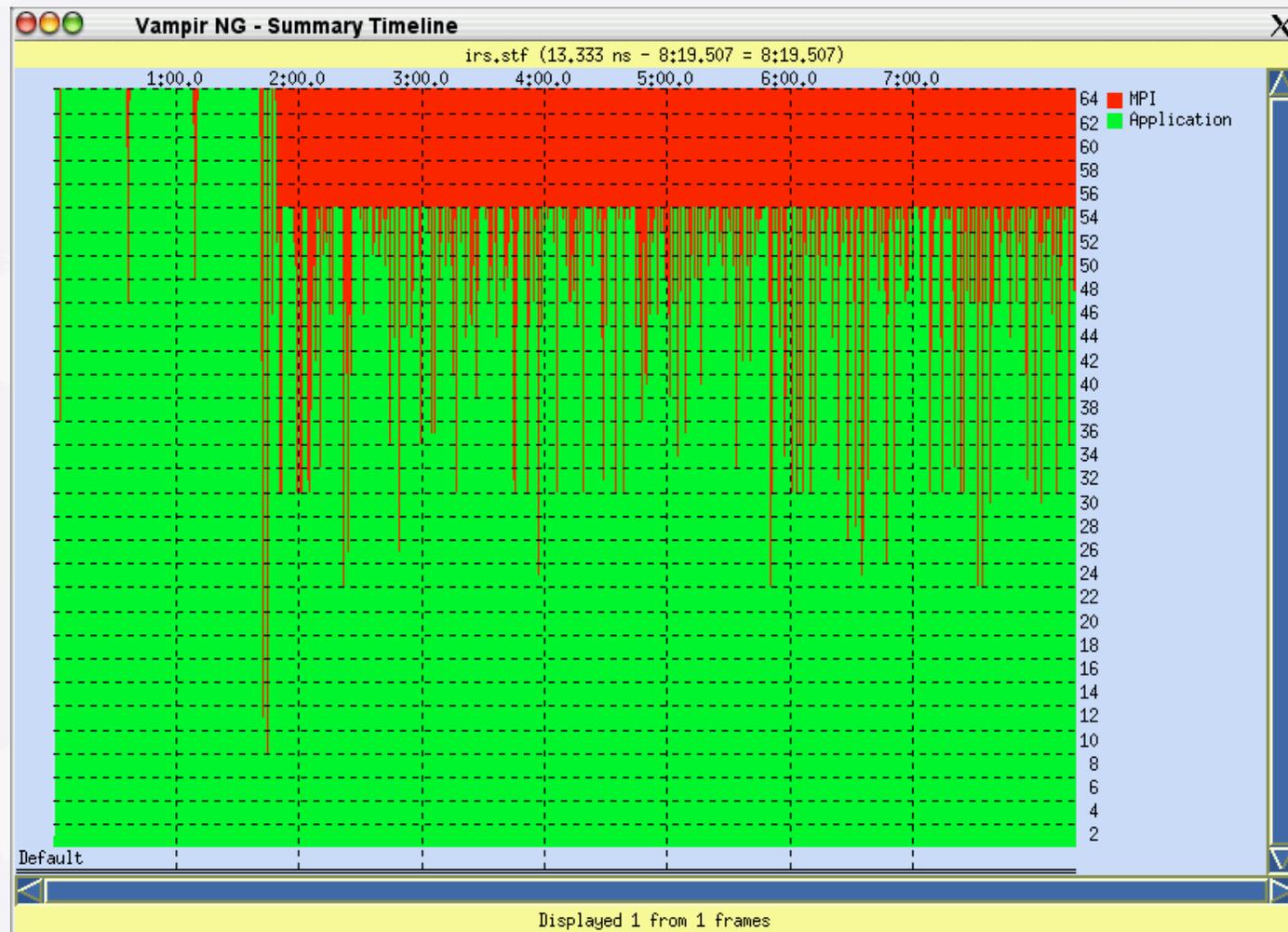
Sorted by Value Down

Alternative Representation

Vampir Displays - Summary Chart

- statistical overview over functions
resp. groups of functions
- as bar chart, pie chart or table
- global (all processes) or local (single process)
- exclusive/inclusive time, occurrences
- absolute or logarithmic scale
- absolute values or percentages
- zoomable
- respect the current time interval
(according to timeline)

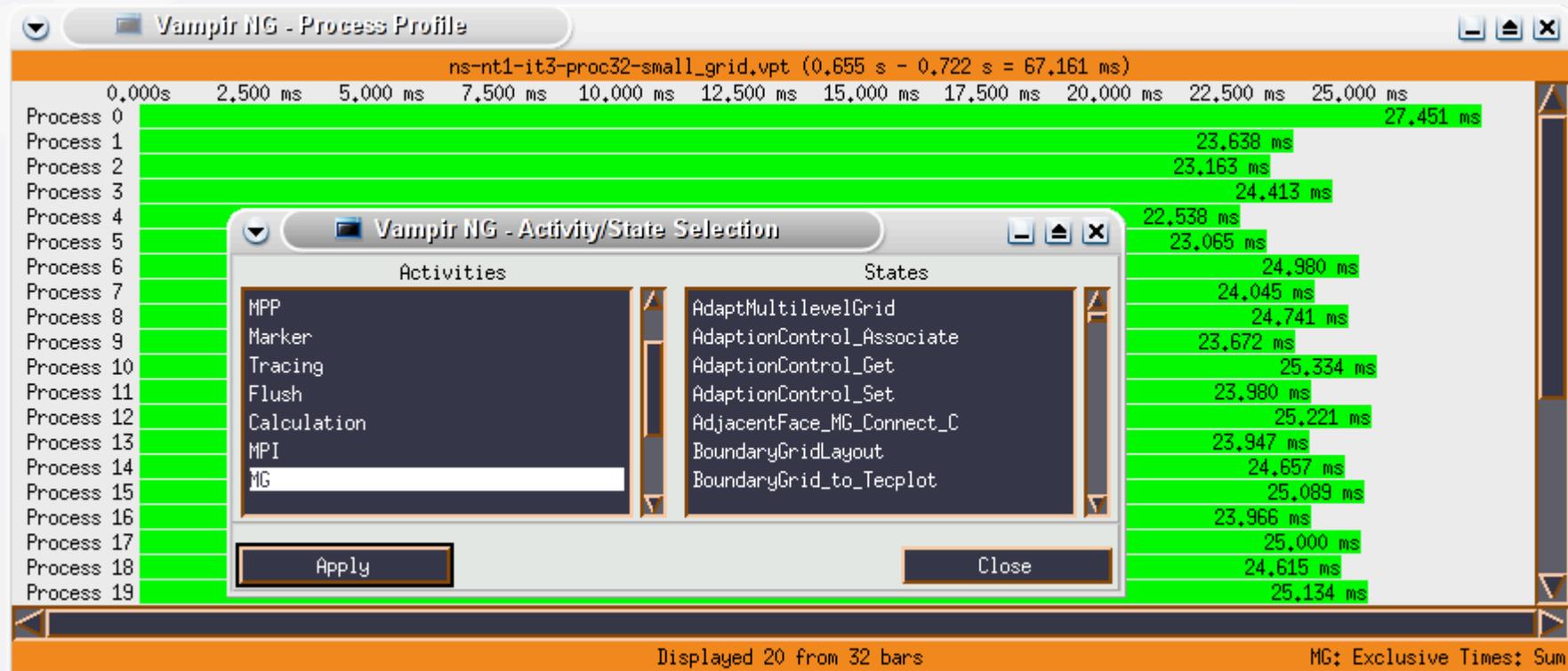
Vampir Displays: Summary Timeline



Vampir Displays - Summary Timeline

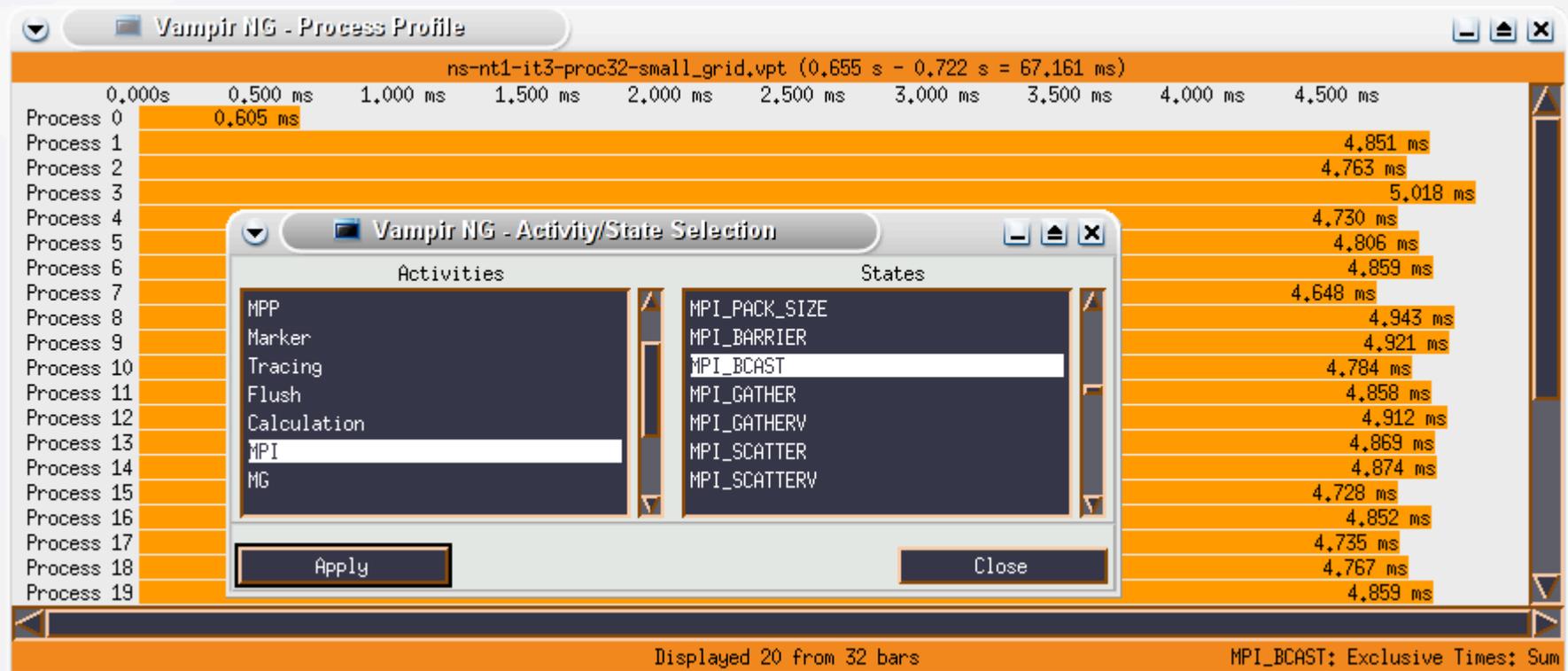
- Summary Chart over time
- same functionality
- zoomable

Vampir Displays - Process Profile



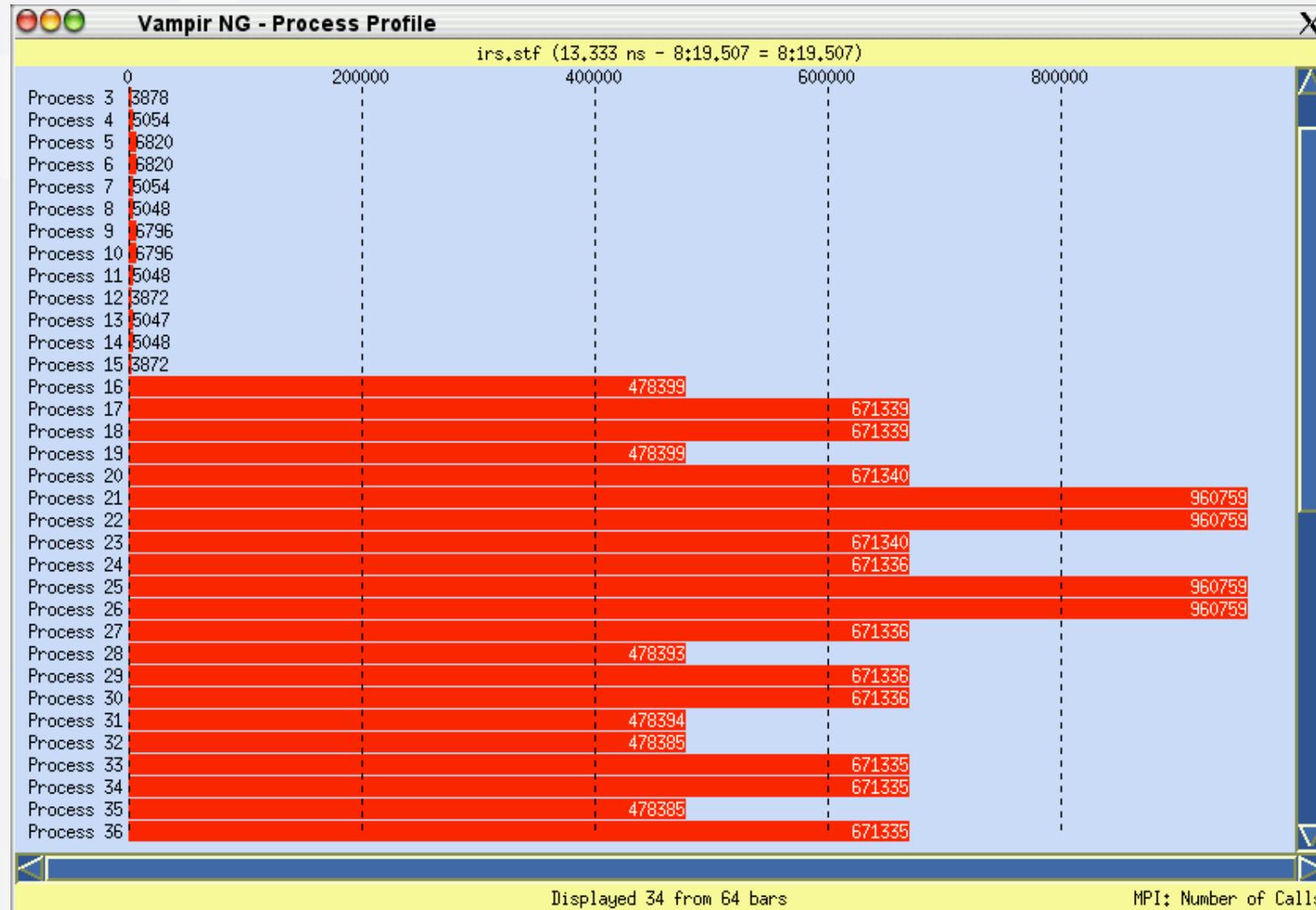
Process Profile

Vampir Displays - Counter Timeline



Process Profile

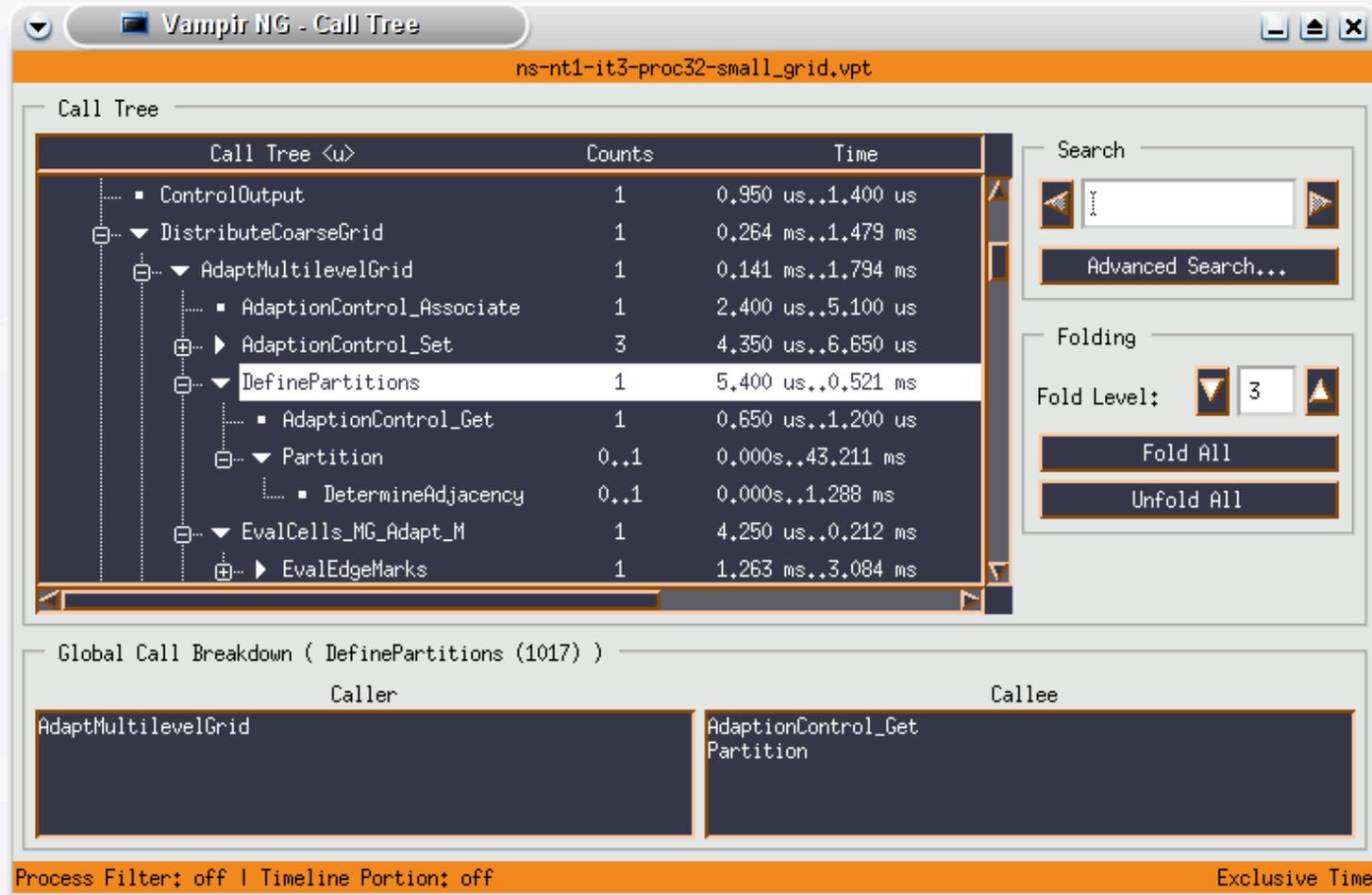
Vampir Display – Process Profile (MPI)



Vampir Displays - Process Profile

- profile for function/group of functions per process
- similar to results of regular profiling
- but now available for arbitrary time intervals

Vampir Displays - Call Tree

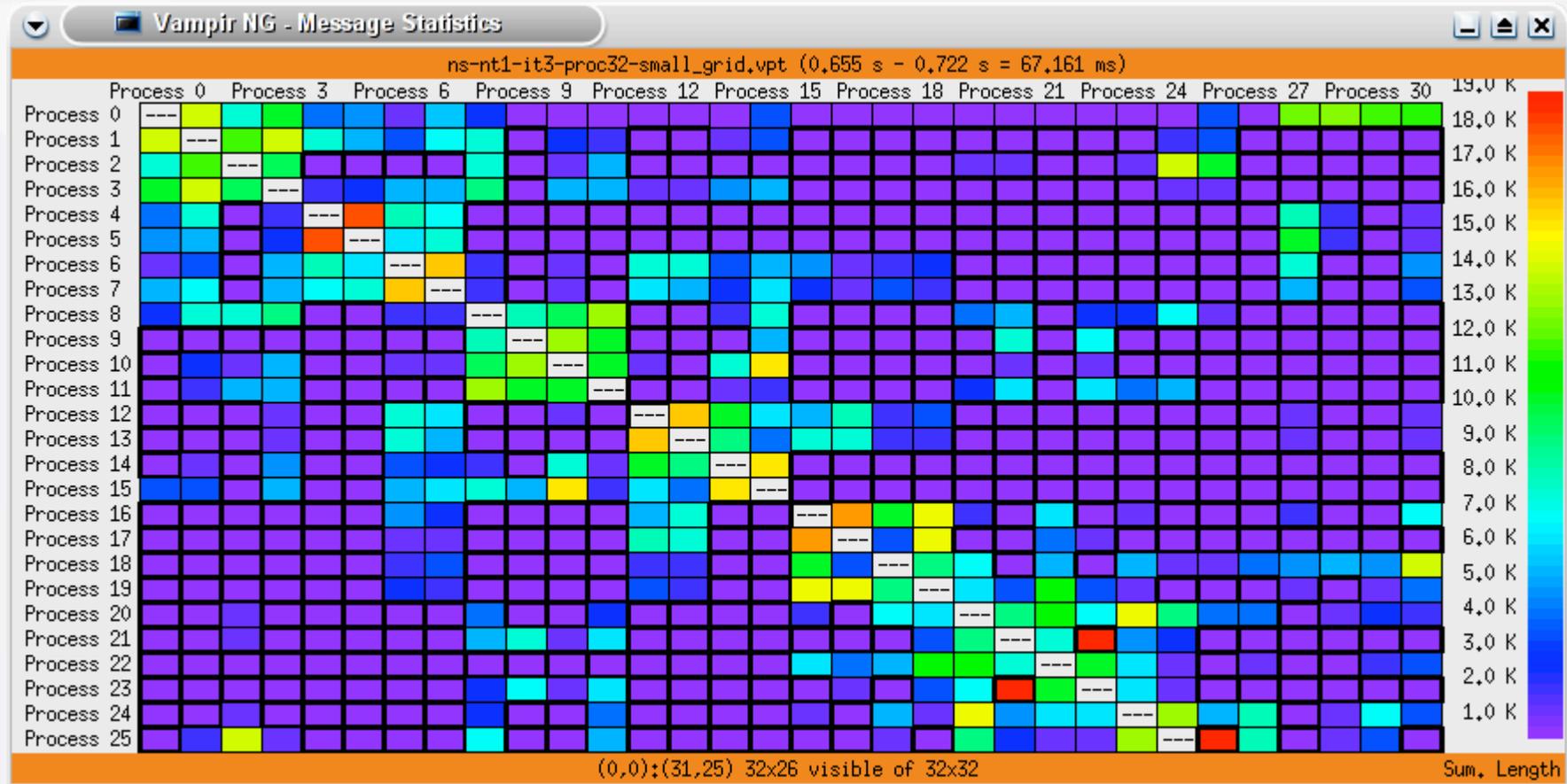


Call Tree

Vampir Displays - Call Tree

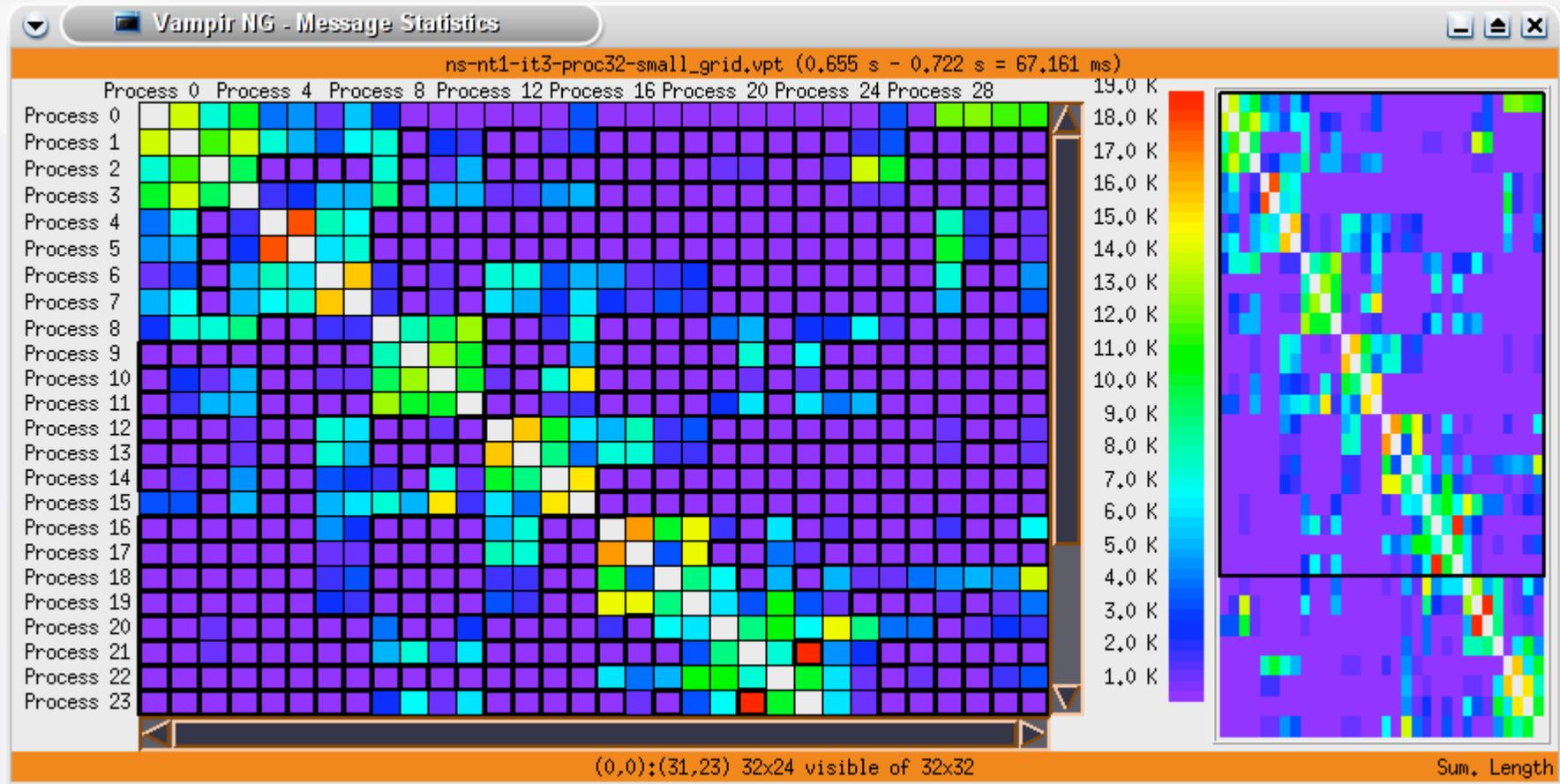
- shows function call hierarchy
- provide callers & callees
- call count
- min/max run time
- fold/unfold
- search

Vampir Displays - Message Statistics



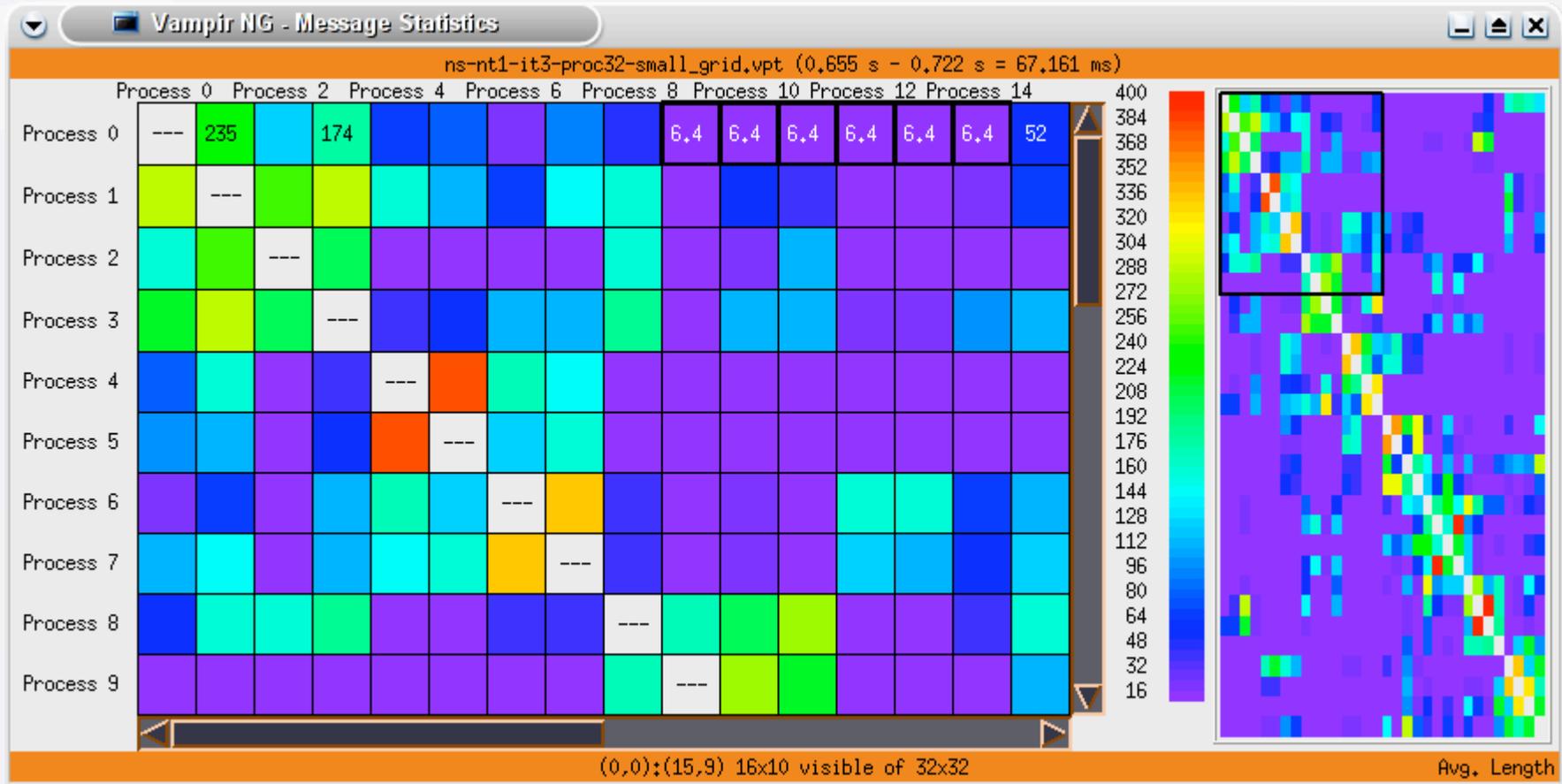
Message Statistics

Vampir Displays - Message Statistics



Message Statistics with Thumbnails

Vampir Displays - Message Statistics

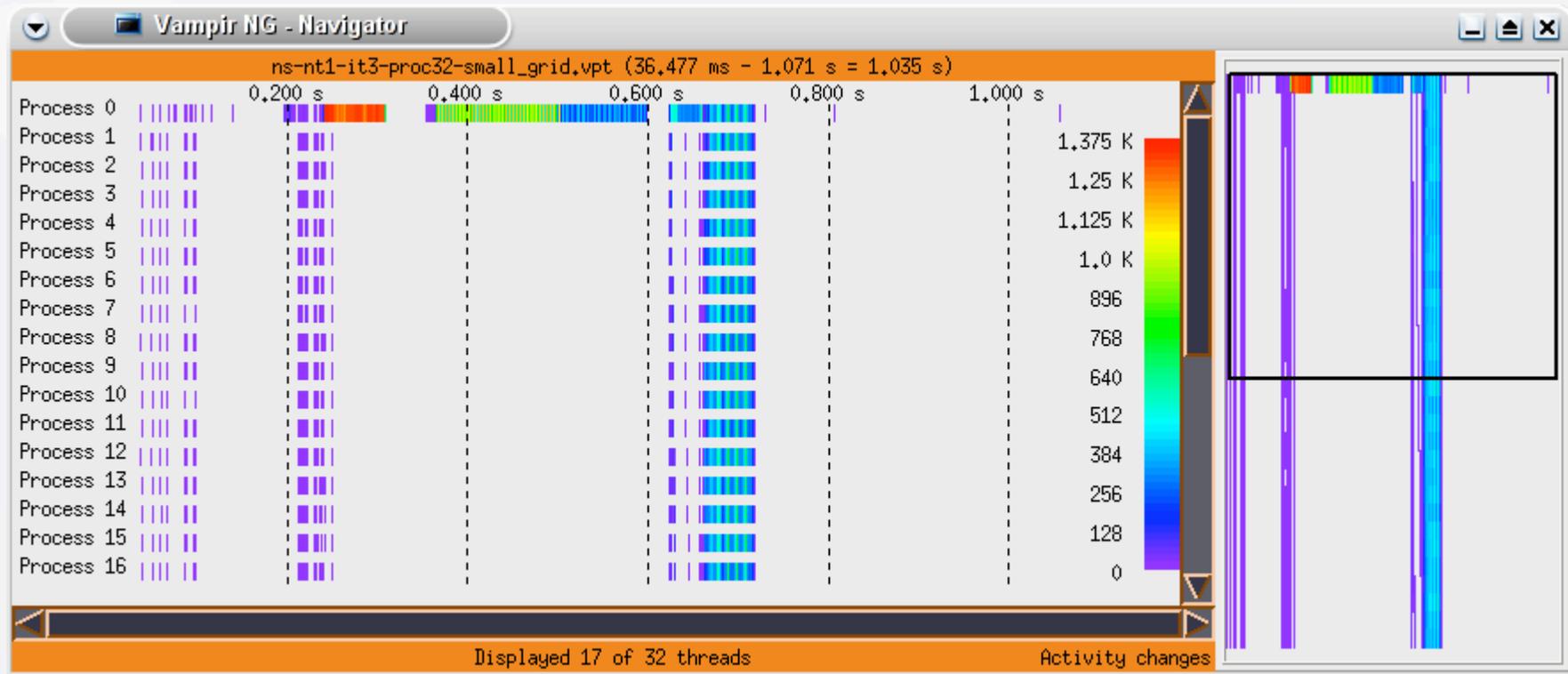


Zoomed Message Statistics

Vampir Displays - Message Statistics

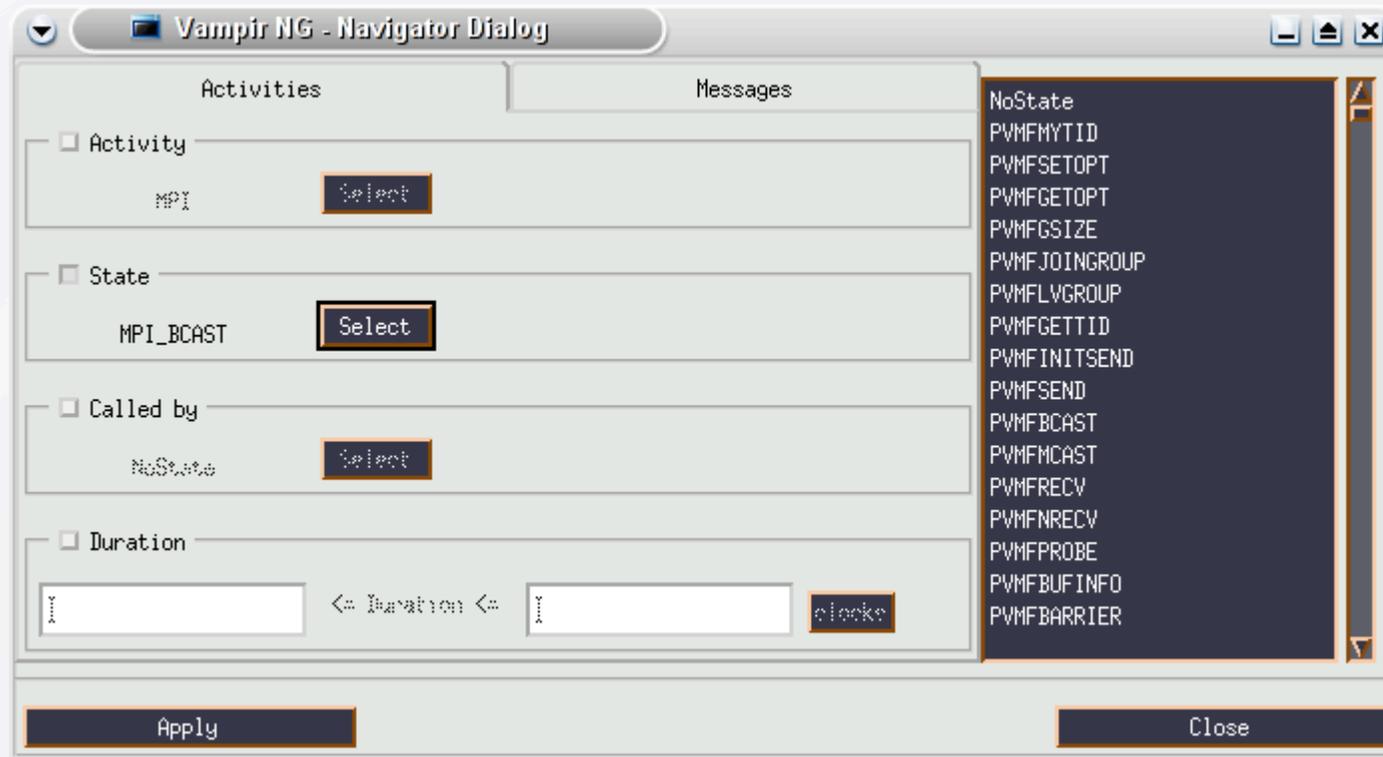
- sender-receiver matrix
- zoomable, for > 1000 processes
- show message properties:
 - Length, rate
 - duration
 - count
- optionally as:
 - min, max, avg, sum
- (message length histogram)

Vampir Displays - Navigator



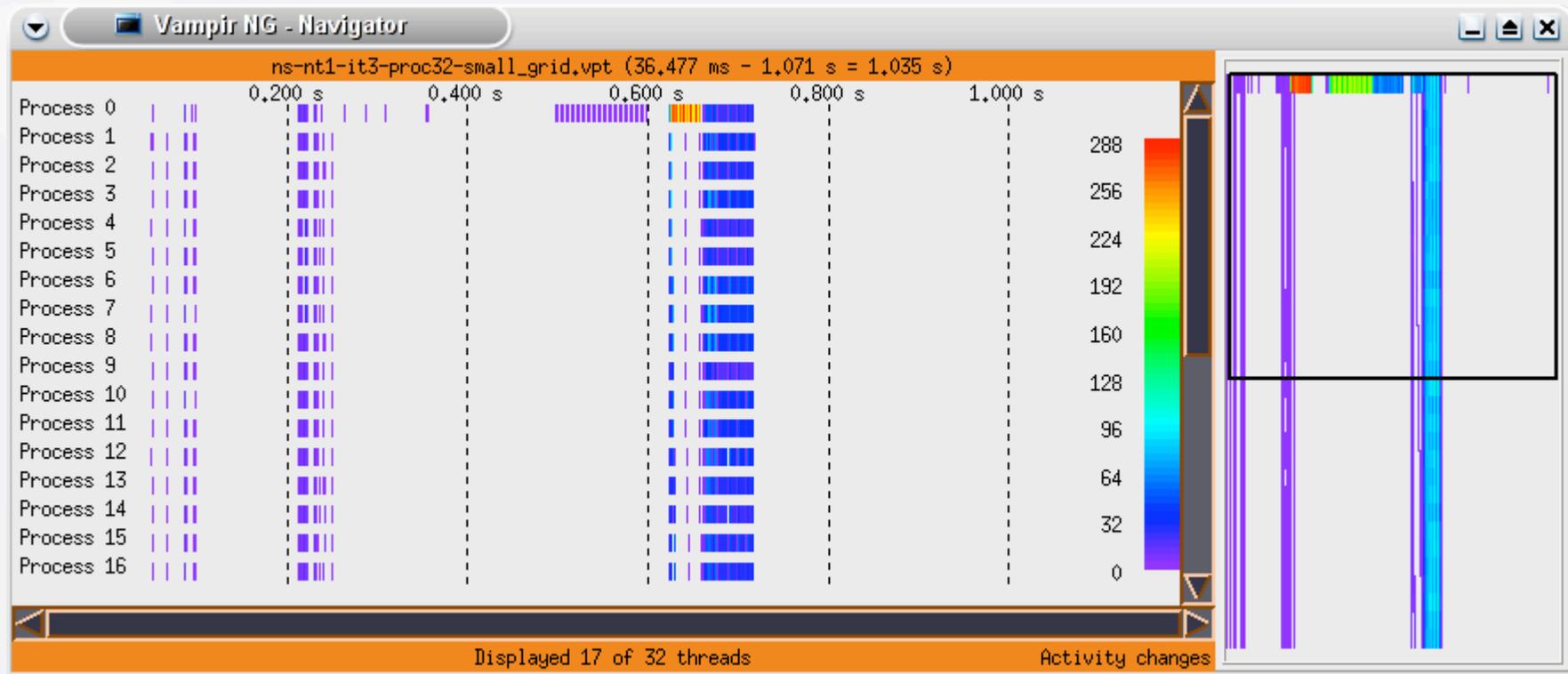
Event Navigator

Vampir Displays - Navigator



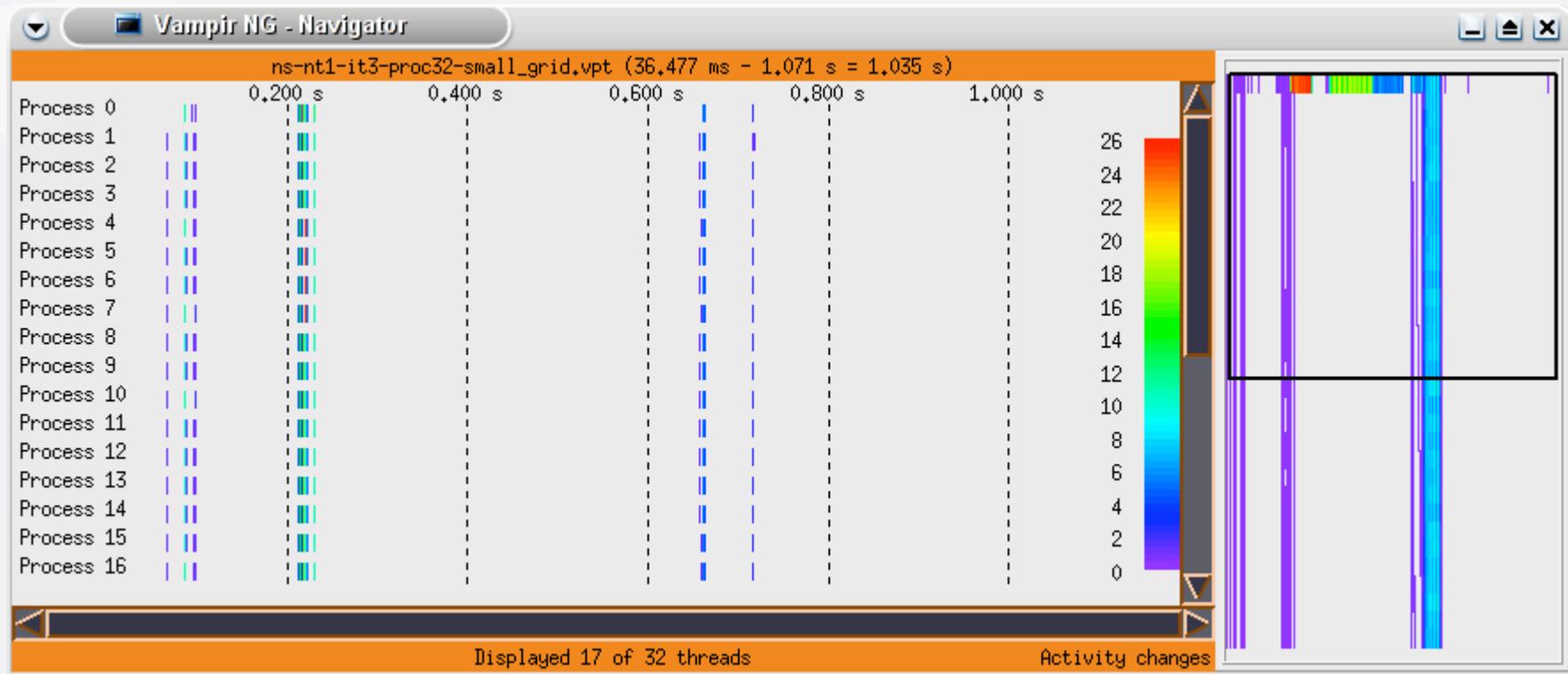
Navigator's Function Selection Dialog

Vampir Displays - Navigator



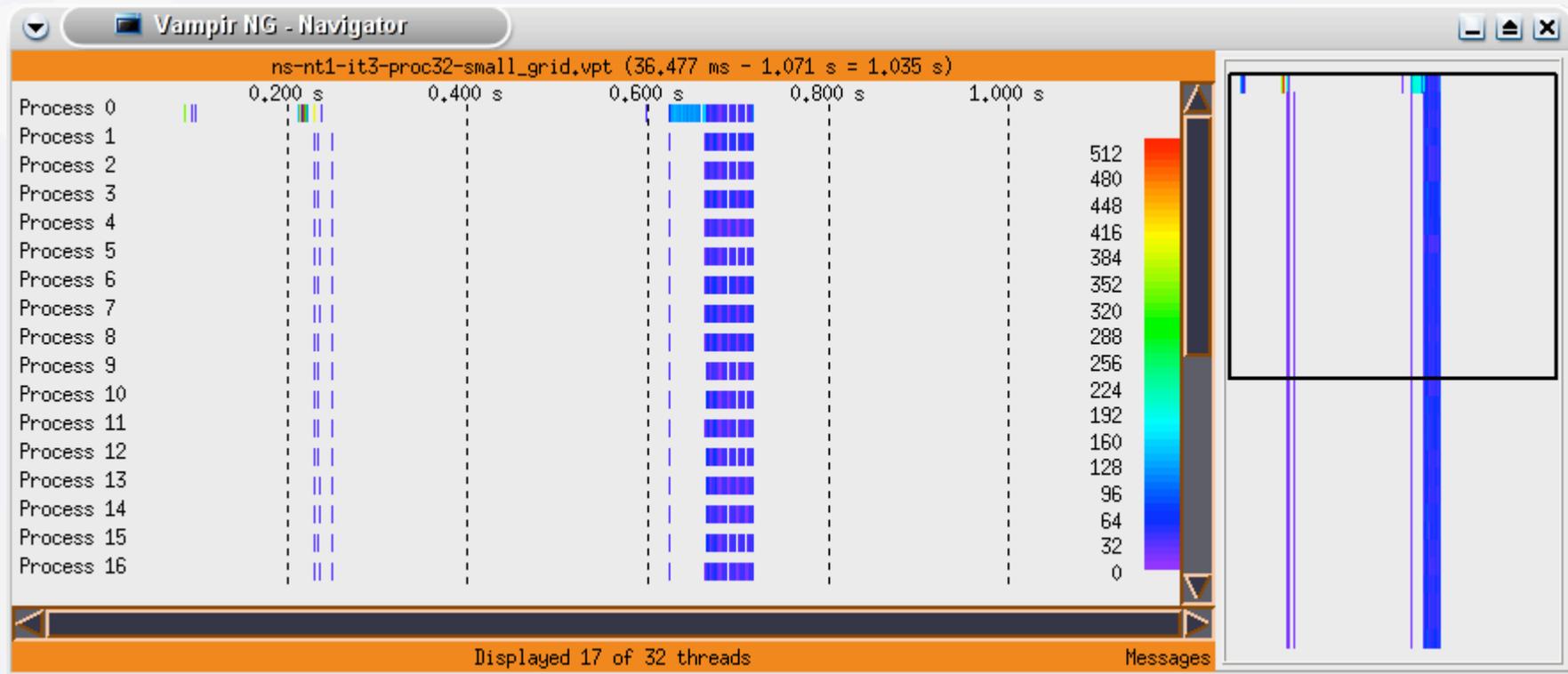
Event Navigator for Function Group (MPI)

Vampir Displays - Navigator



Event Navigator for Function (MPI_Bcast)

Vampir Displays - Navigator



Event Navigator for Messages

Vampir Displays - Navigator

The screenshot shows a window titled "Vampir NG - Navigator Dialog" with two tabs: "Activities" and "Messages". The "Messages" tab is active. The dialog contains several filter sections, each with a checkbox and a label:

- Length: Two text input fields with a "<= Length <=" label between them.
- Sent by: A "Process" label followed by a "Select" button.
- Received by: A "Process" label followed by a "Select" button.
- Duration: Two text input fields with a "<= Duration <=" label between them, and a "clock" icon to the right.
- Data Rate: Two text input fields with a "<= Data Rate <=" label between them.
- Tag/Comm.: "Comm: 1" followed by a "Select" button, and "Tag: 0" followed by a "Select" button.

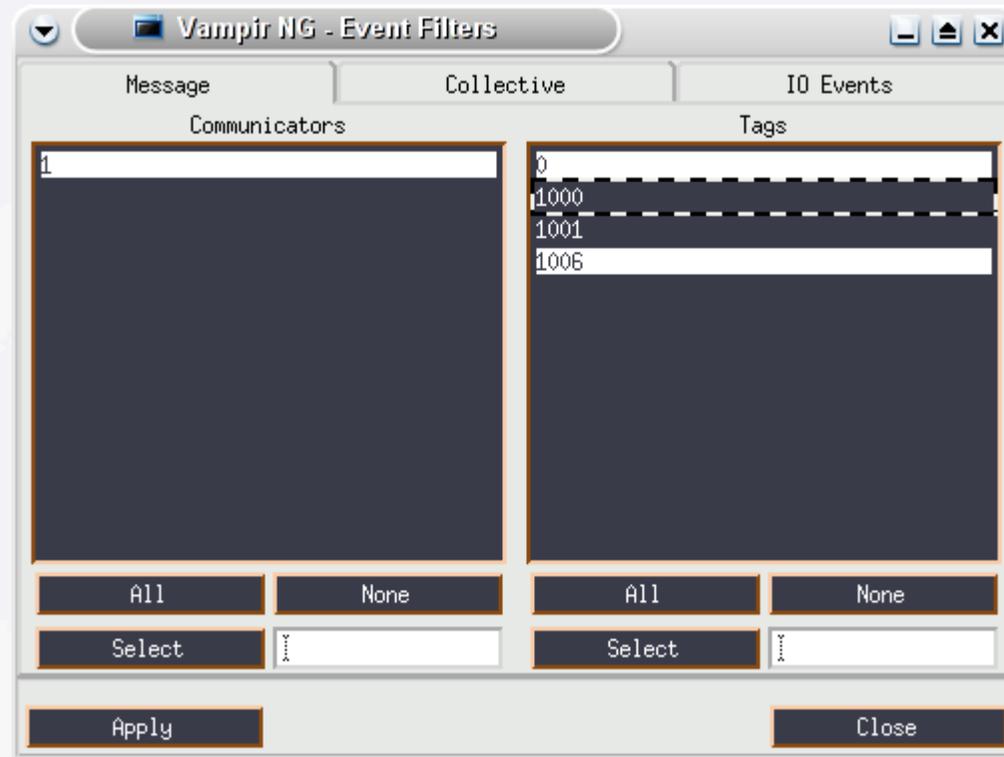
At the bottom of the dialog are two buttons: "Apply" and "Close".

Navigator's Message Selection Dialog

Vampir Displays - Navigator

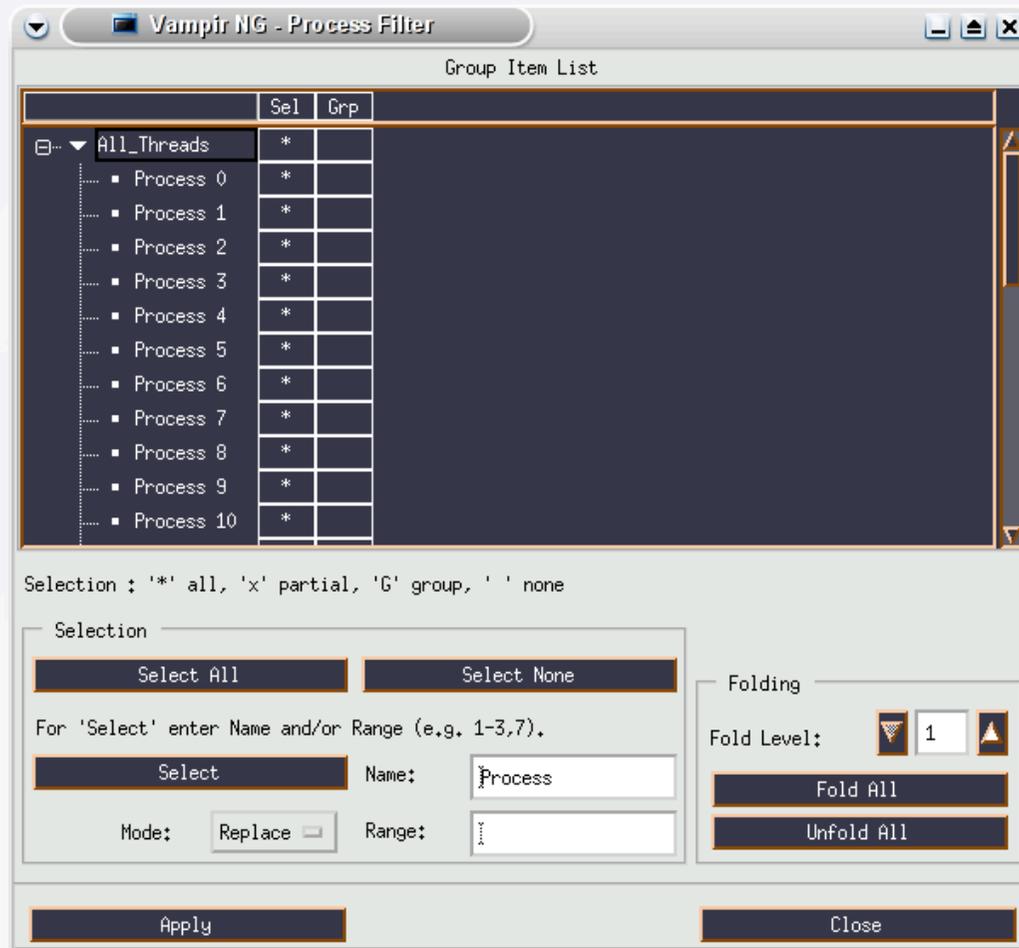
- find occurrences of certain kinds of events
- show frequency by color scale
- filter functions by
 - functions, function groups, caller
 - duration interval
- filter messages by
 - sender, receiver processes
 - message length, duration or speed
 - MPI communicator or tag
- zoomable in time

Vampir Displays - Global Filters



Event Filter Dialog

Vampir Displays - Global Filters

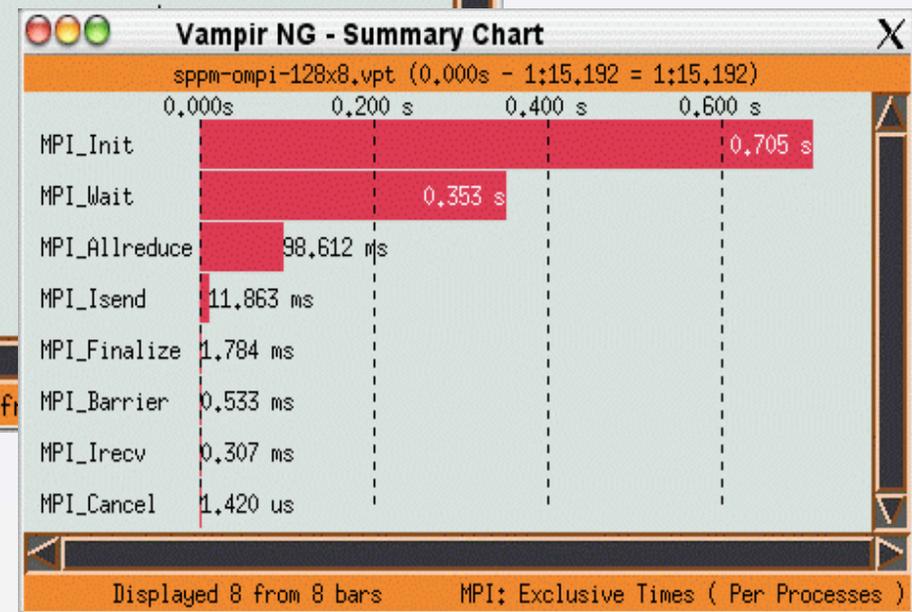
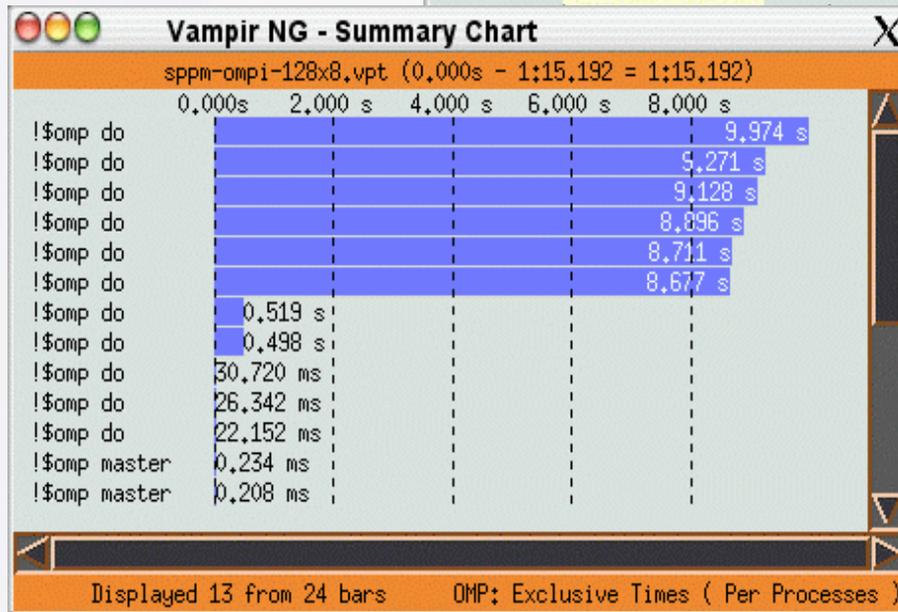
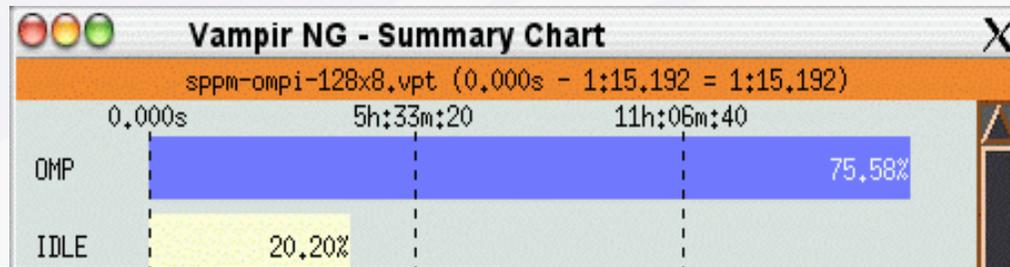


Process Filter Dialog

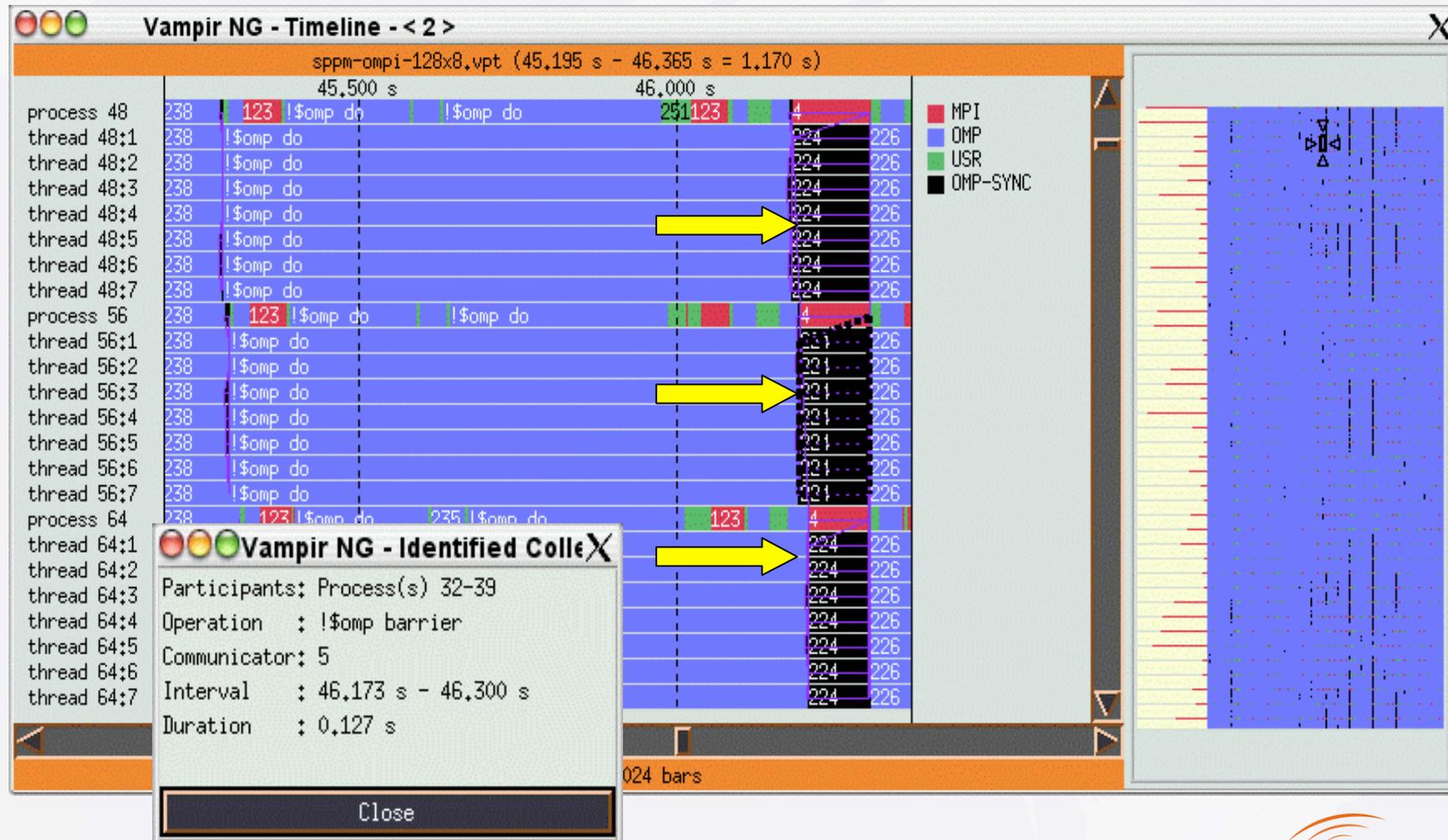
Vampir Displays - Global Filters

- ignore certain items globally
- processes/threads or groups of them
- messages
 - by communicator or by tag
- collective operations
 - by communicator or by type
- I/O events
 - by communicator
 - by file
 - by read/write access

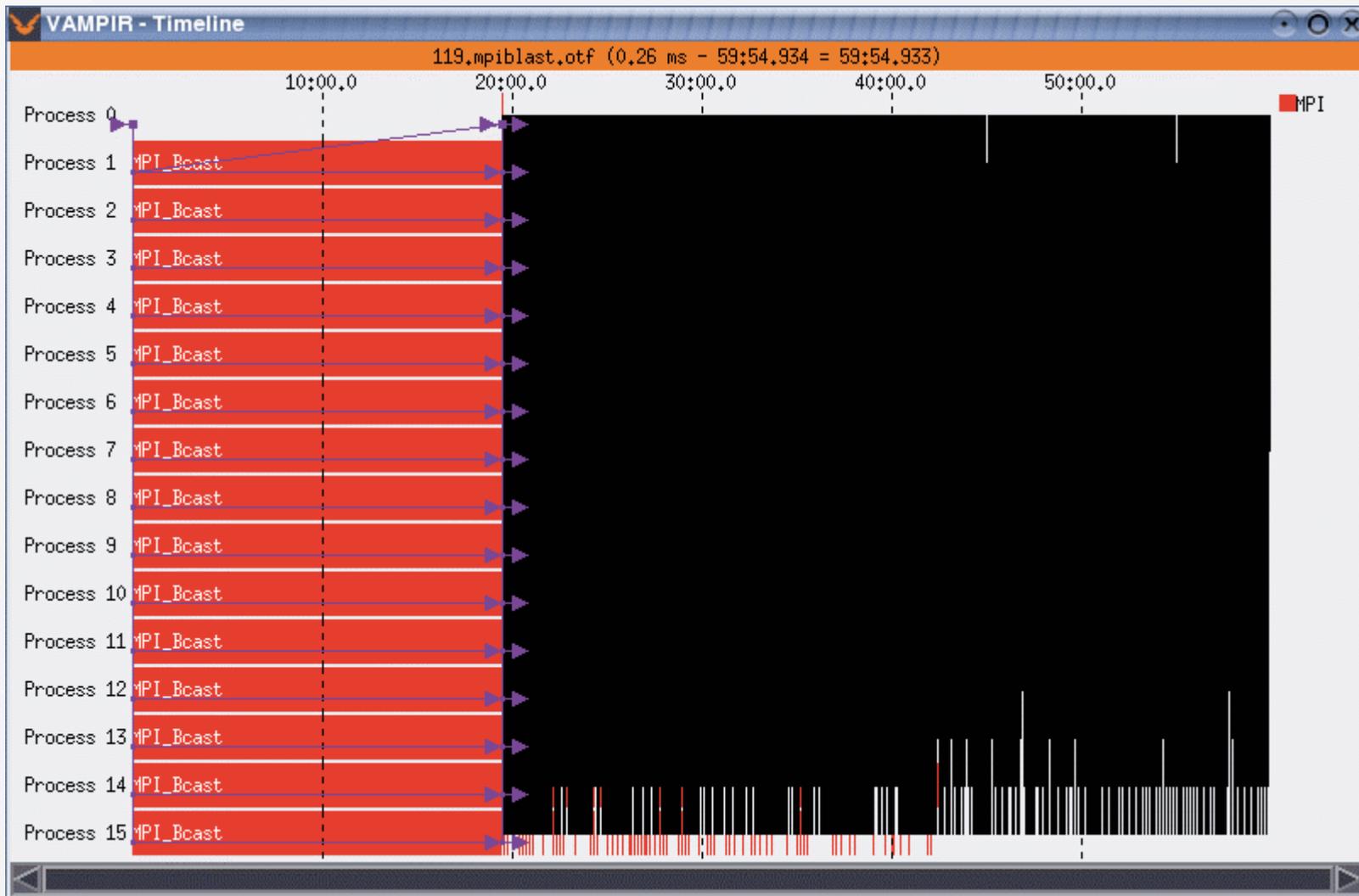
Vampir Displays - OpenMP and MPI Profiles



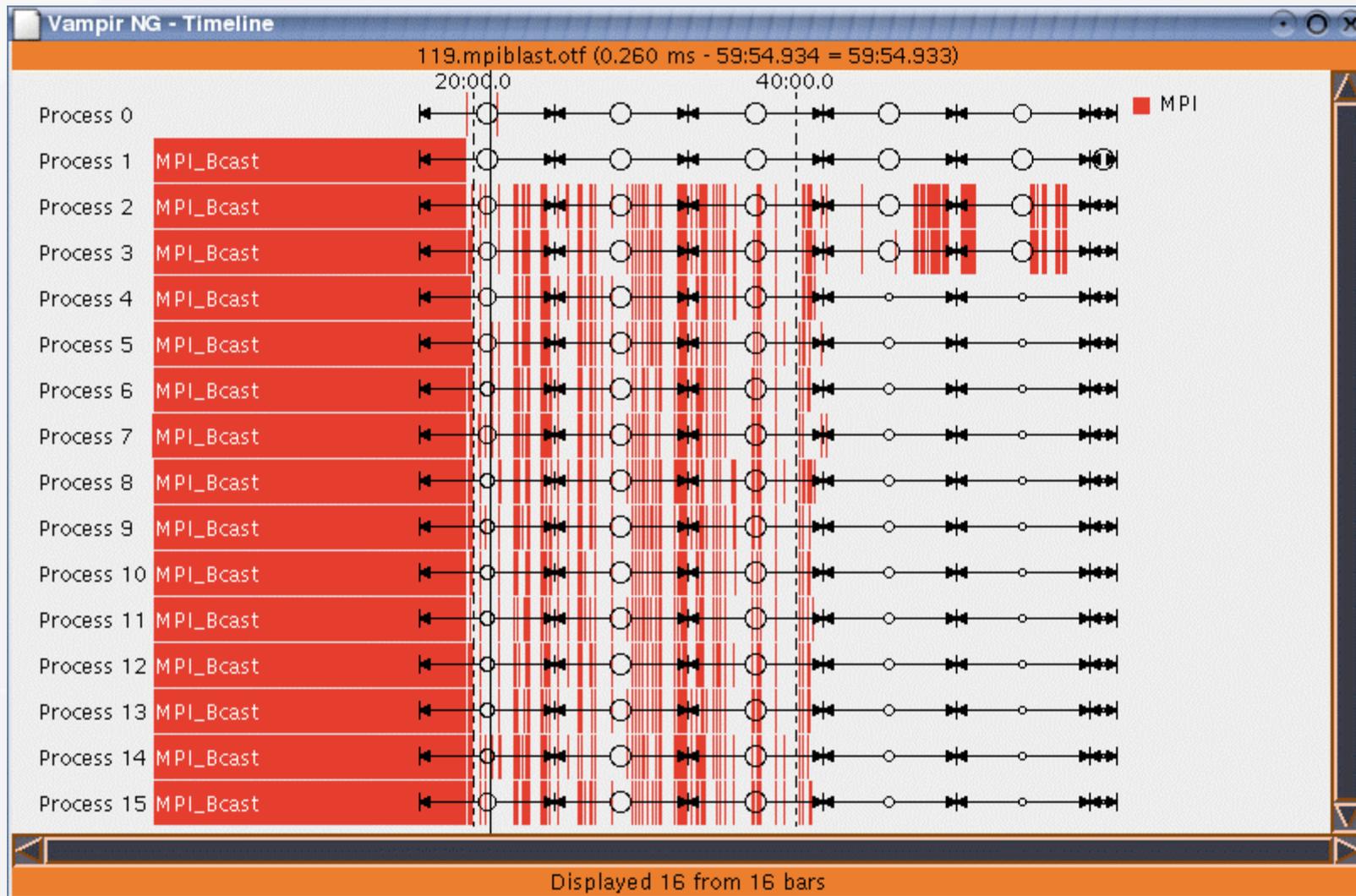
Vampir Displays - OpenMP Barrier Synchronization



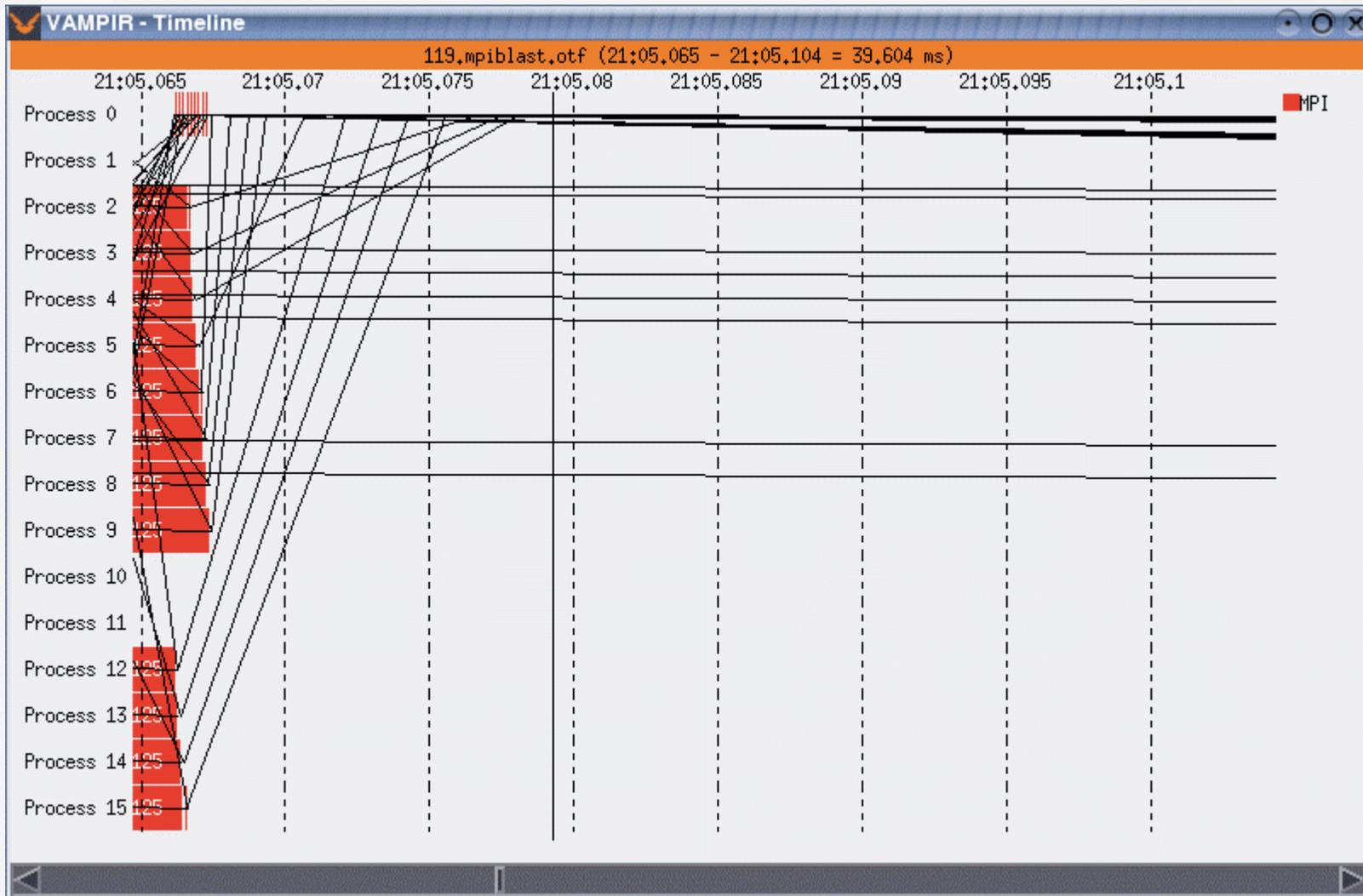
Vampir (old)



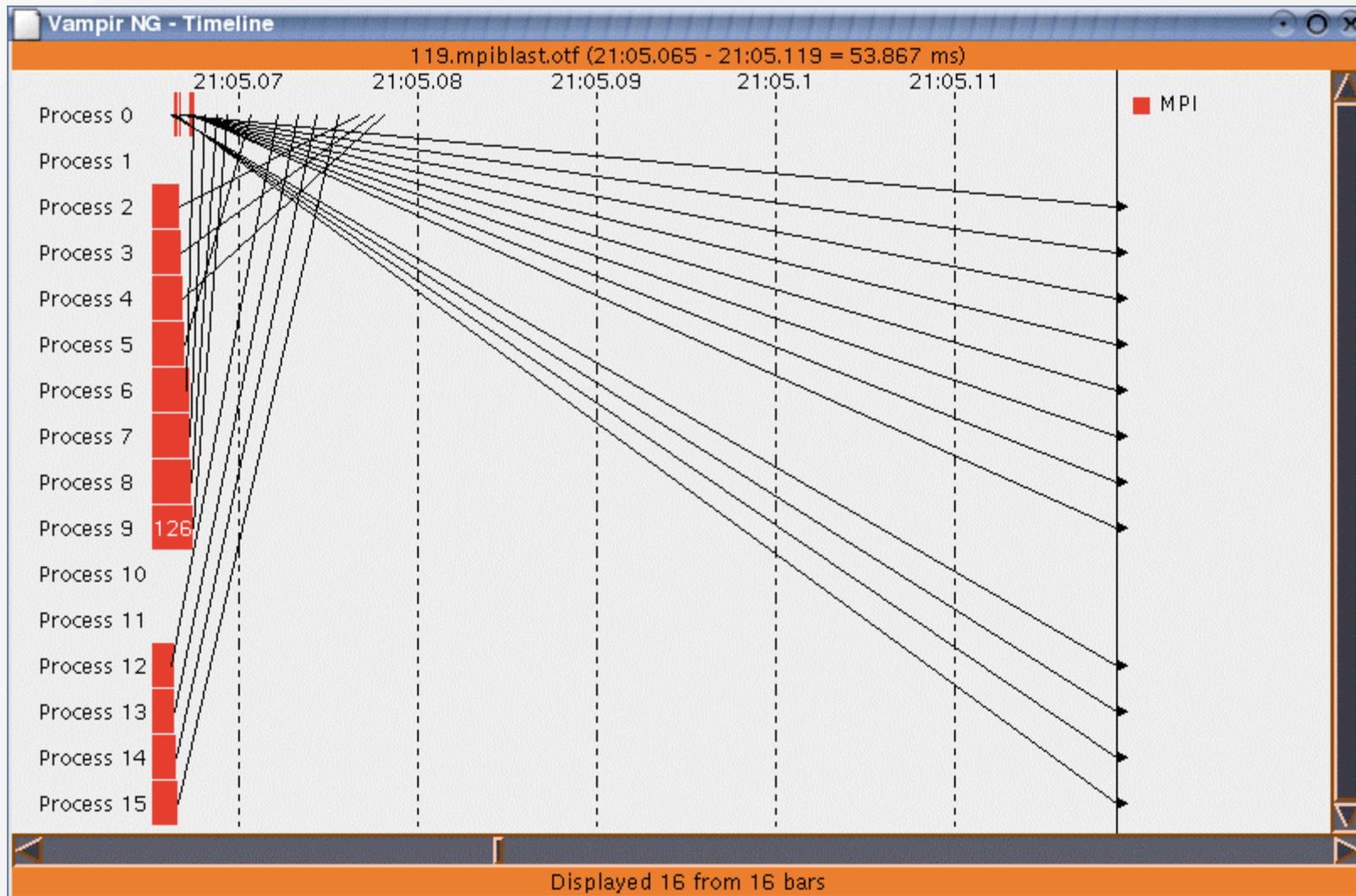
Vampir (new)



Vampir (old)



Vampir (new)



Summary

Vampir/Vampir NG offers

- Comprehensive performance analysis environment
- Demonstrated scalability to 5000 processors
- Perfect solution for large centers with remote users and a variety of platforms

Vampir Future Developments

- I/O Analysis
- One sided communication

Vampir Literature

- H. Brunst, H.-Ch. Hoppe, W. E. Nagel, and M. Winkler: **Performance Optimization for Large Scale Computing: The Scalable VAMPIR Approach.** Proc. of ICCS 2001, San Francisco, CA, USA, May 28-30 2001, Springer LNCS 2074, pp 751-760, 2001
- H. Brunst, D. Kranzlmüller, W. E. Nagel: **Tools for Scalable Parallel Program Analysis - Vampir NG and DeWiz.** Distributed and Parallel Systems, Cluster and Grid Computing, Springer, Kluwer Intl. Series in Engineering and Computer Science, Vol 777, Budapest, Hungary, 2004
- A. Knüpfer, H. Brunst, and W. E. Nagel: **High Performance Trace Visualization.** In Proc. of 13th Euromicro Conference on Parallel, Distributed and Network-based Processing, Lugano, Switzerland, February 9-11 2005
- Vampir User Guide

Performance-Analyse mit Vampir

Am Beispiel eines komplexen Modellsystems

Zellescher Weg 12

Willers-Bau A106

Tel. +49 351 - 463 31945

Matthias Lieber

Matthias.Lieber@tu-dresden.de

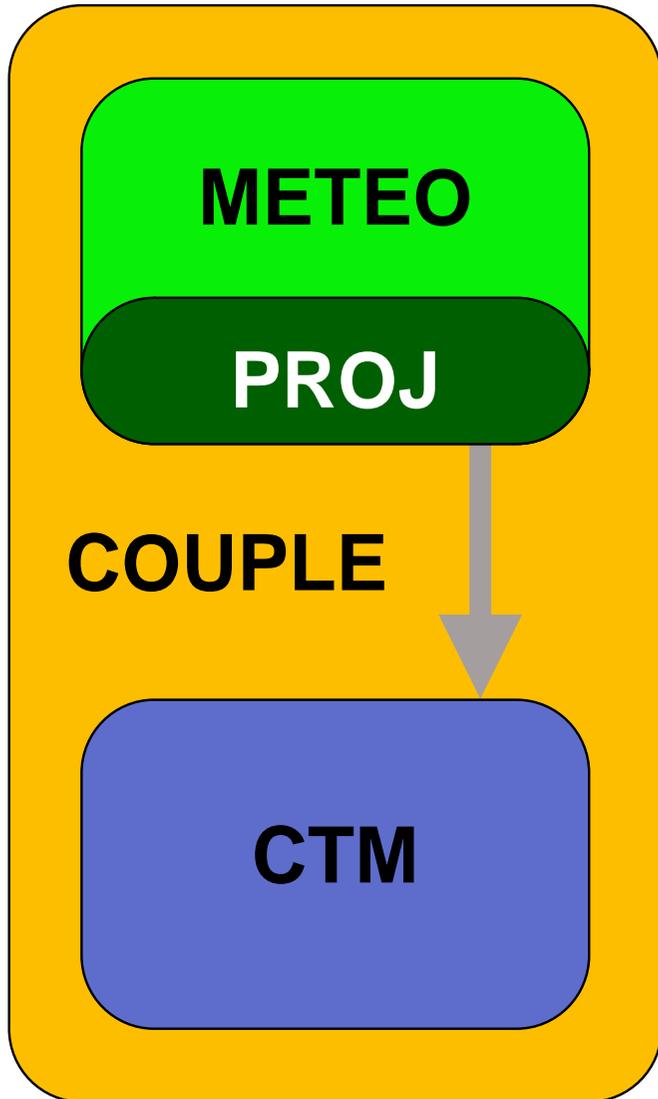
Inhalt

- kurze Einführung zum Modellsystem
- Analyse der Tracefiles mit Vampir NG

Modellsystem

- Simulation von Umwandlungs- und Transportprozessen von Schadstoffen in der Atmosphäre
- Anwendungen: Ozonbelastung, Feinstaub-Bildung, Saharastaub-Ausbreitung
- entwickelt am Leibniz-Institut für Troposphärenforschung Leipzig (IfT)
- Kopplung zweier unabhängiger Simulationsmodelle
 - Wettervorhersagemodell des DWD (**METEO**)
 - Chemie-Transport-Modell (**CTM**)

Kopplung



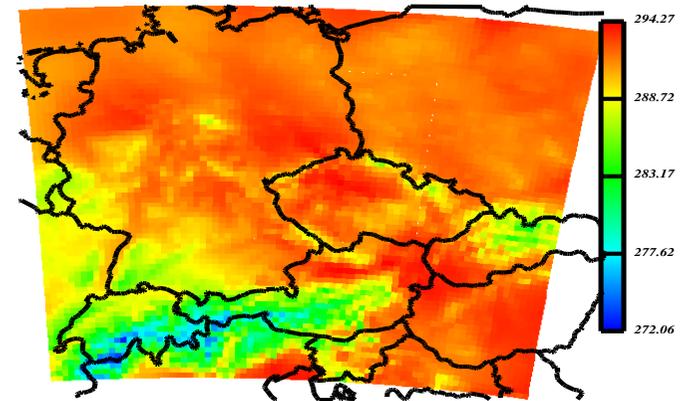
Temp.
Druck
Wind

...

Feinstaub
Ozon
SO₂

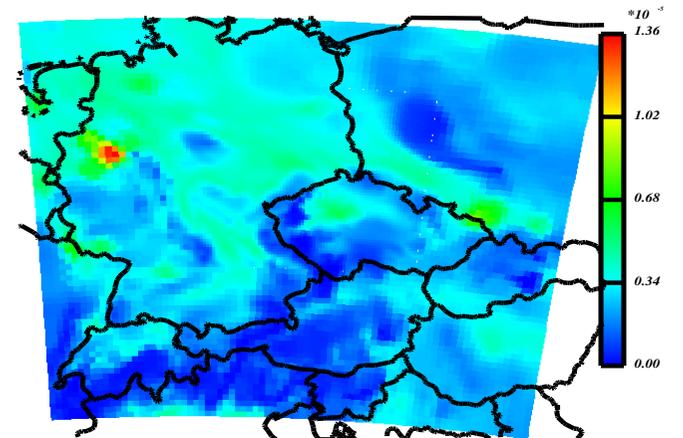
...

TEMP [K] euro.avs



Sat Aug 10 00:00:00 2002

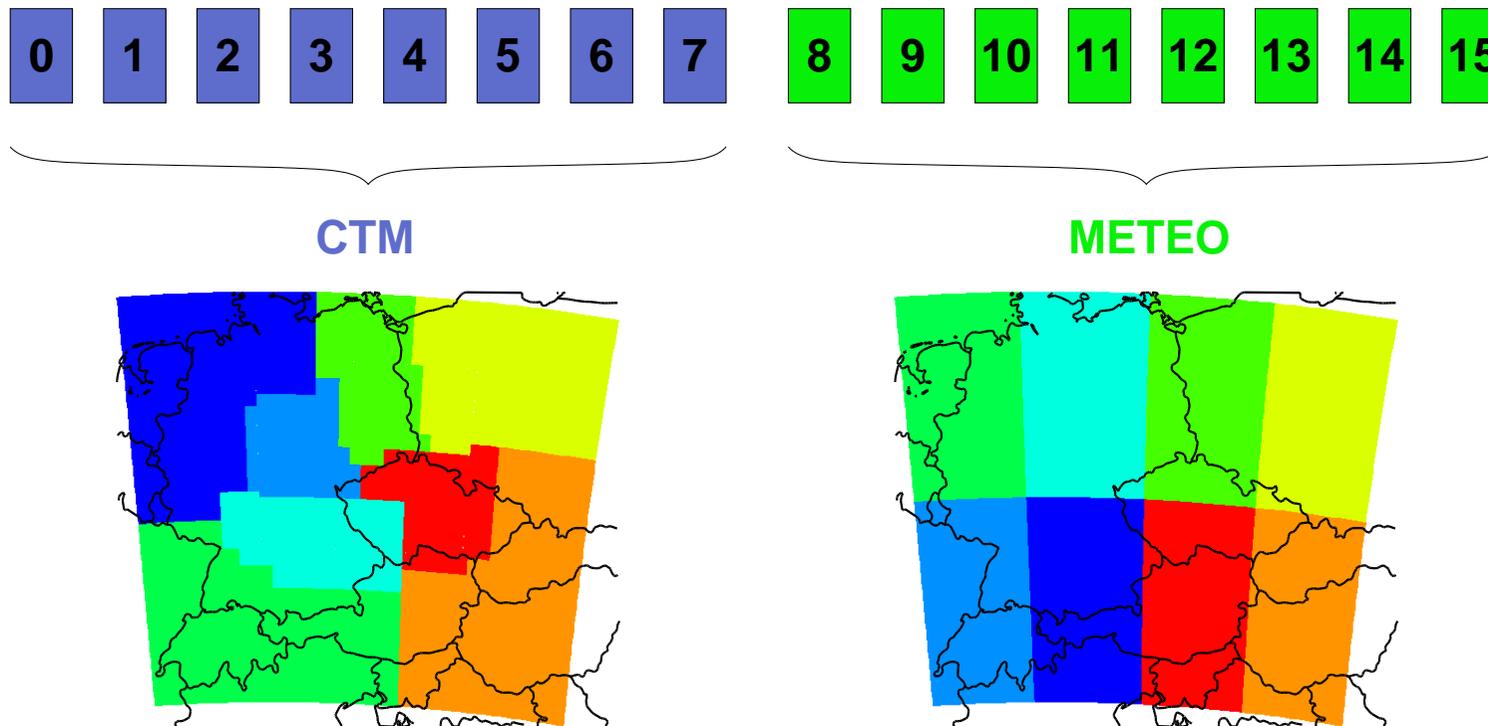
PPM10 [g/m3] euro.avs



Sat Aug 10 06:00:00 2002

Parallelisierung

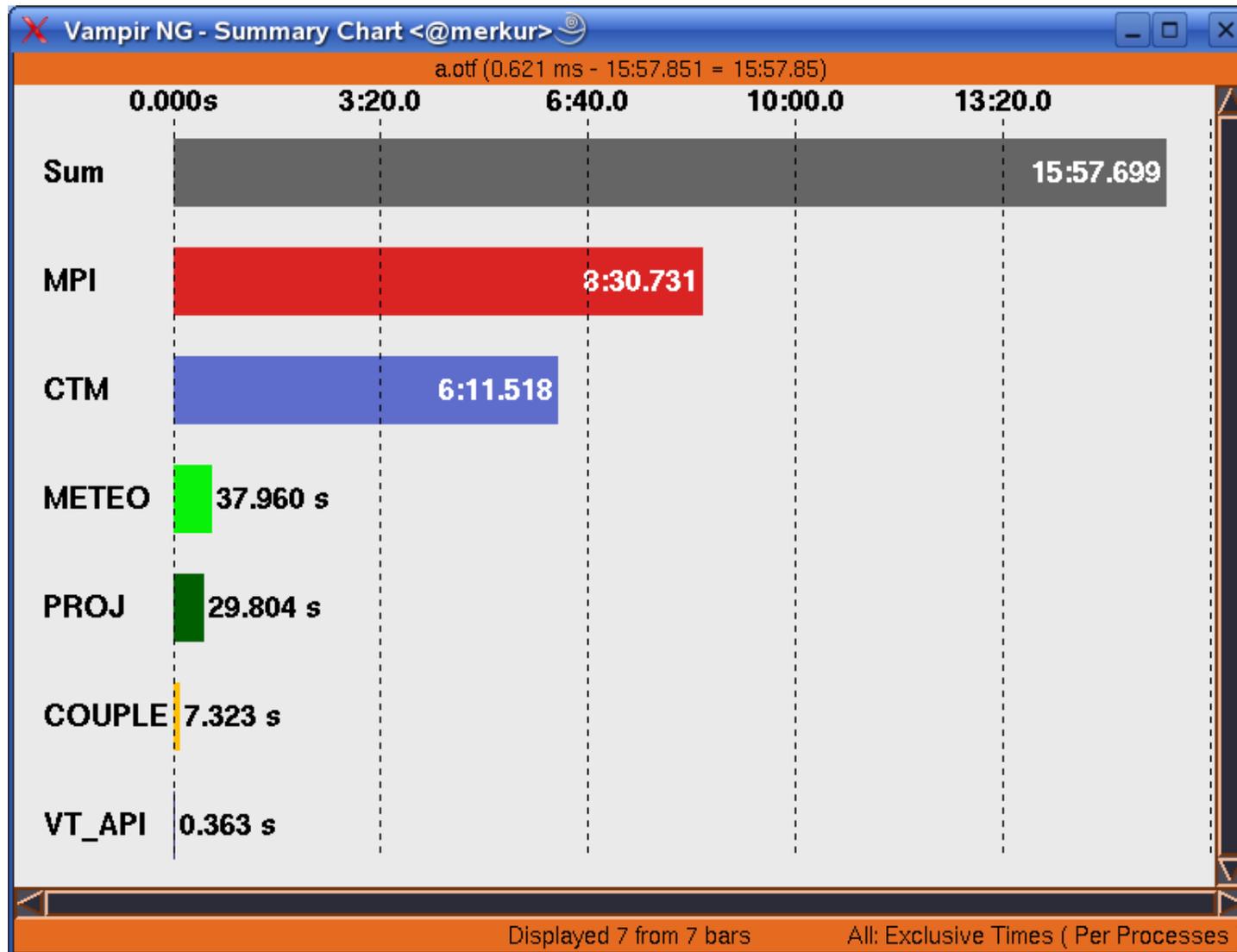
- Aufteilung der CPUs in **CTM** und **METEO**
- innerhalb der Modelle durch Aufteilung des Modell-Gebietes unter den Prozessoren



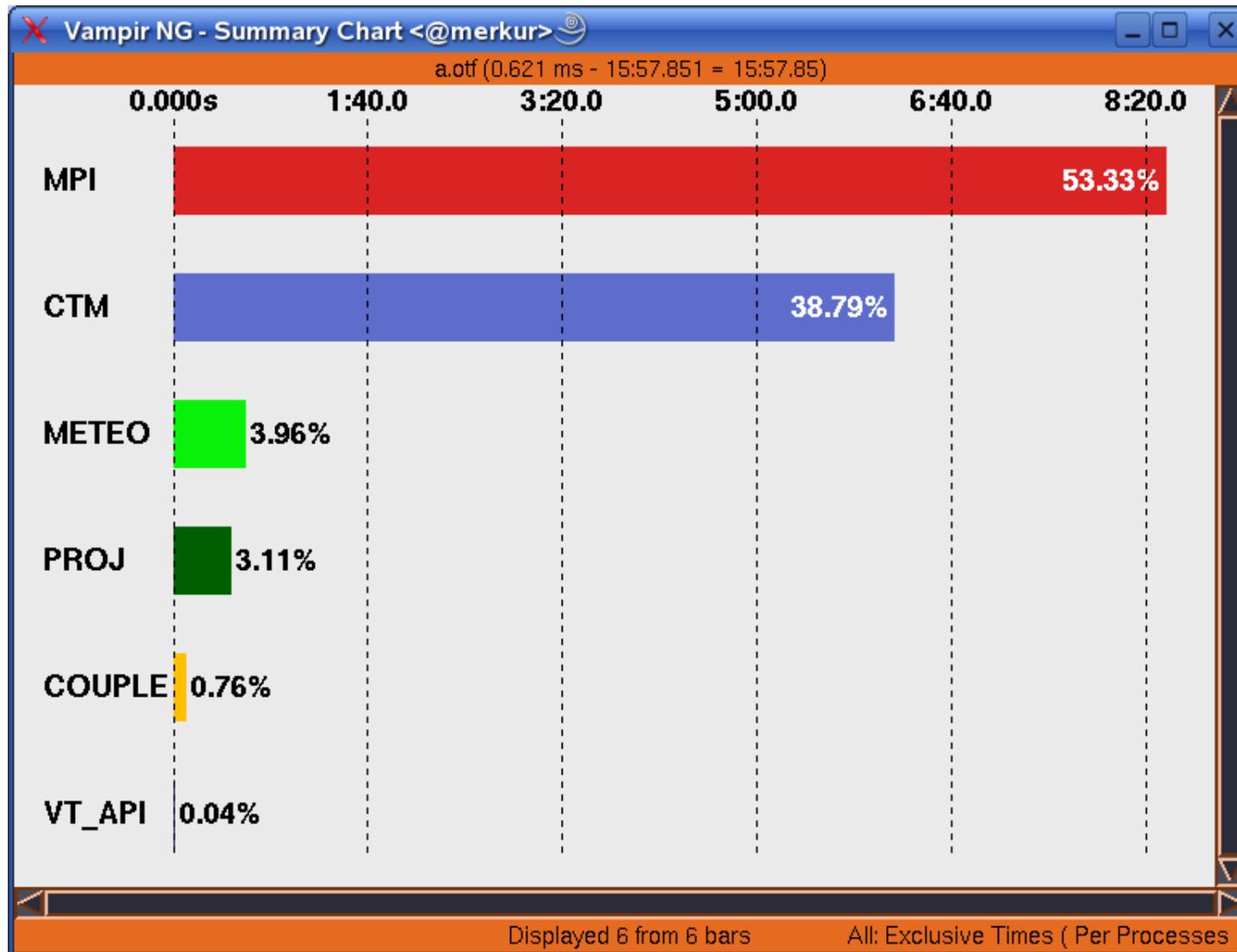
Konfiguration der Performance-Analyse

- Gruppen (*activities*):
 - Chemie-Transport-Modell (**CTM**)
 - Wettervorhersagemodell des DWD (**METEO**)
 - Projektion – notwendige Korrektur von Wind-Daten (**PROJ**)
 - Routinen zur Kopplung der beiden Modelle (**COUPLE**)
 - Kommunikation und Synchronisation via MPI (**MPI**)
- 16 CPUs auf merkur – 8 für **CTM** und 8 für **METEO**
- Vorhersagezeit: 6h
- Filter für Funktionen notwendig um Größe der Tracefiles zu reduzieren
- Größe der Tracefiles: ca. 300 MB (ohne Filter: ca. 4 GB)

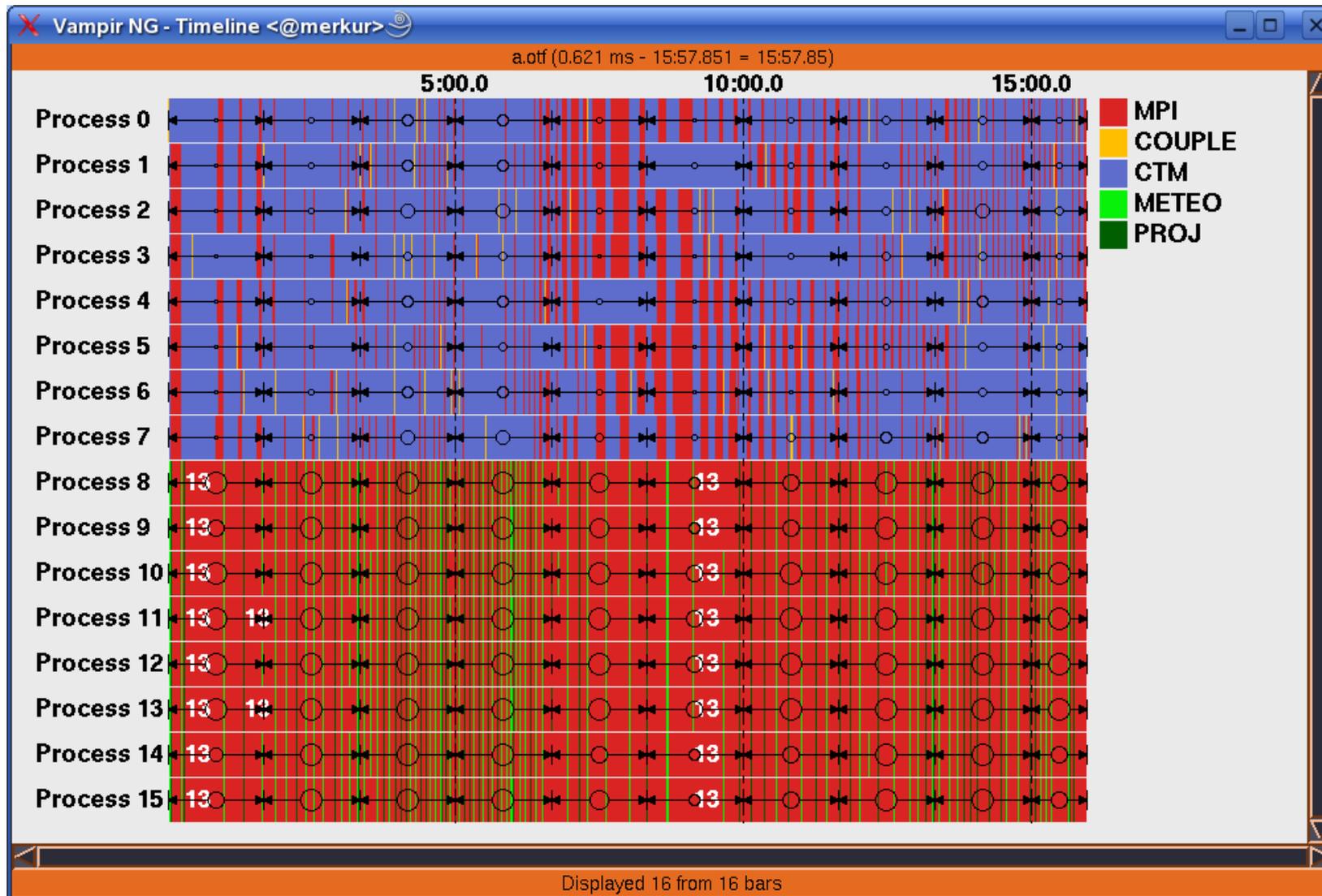
Summary Chart - Gesamtlaufzeit



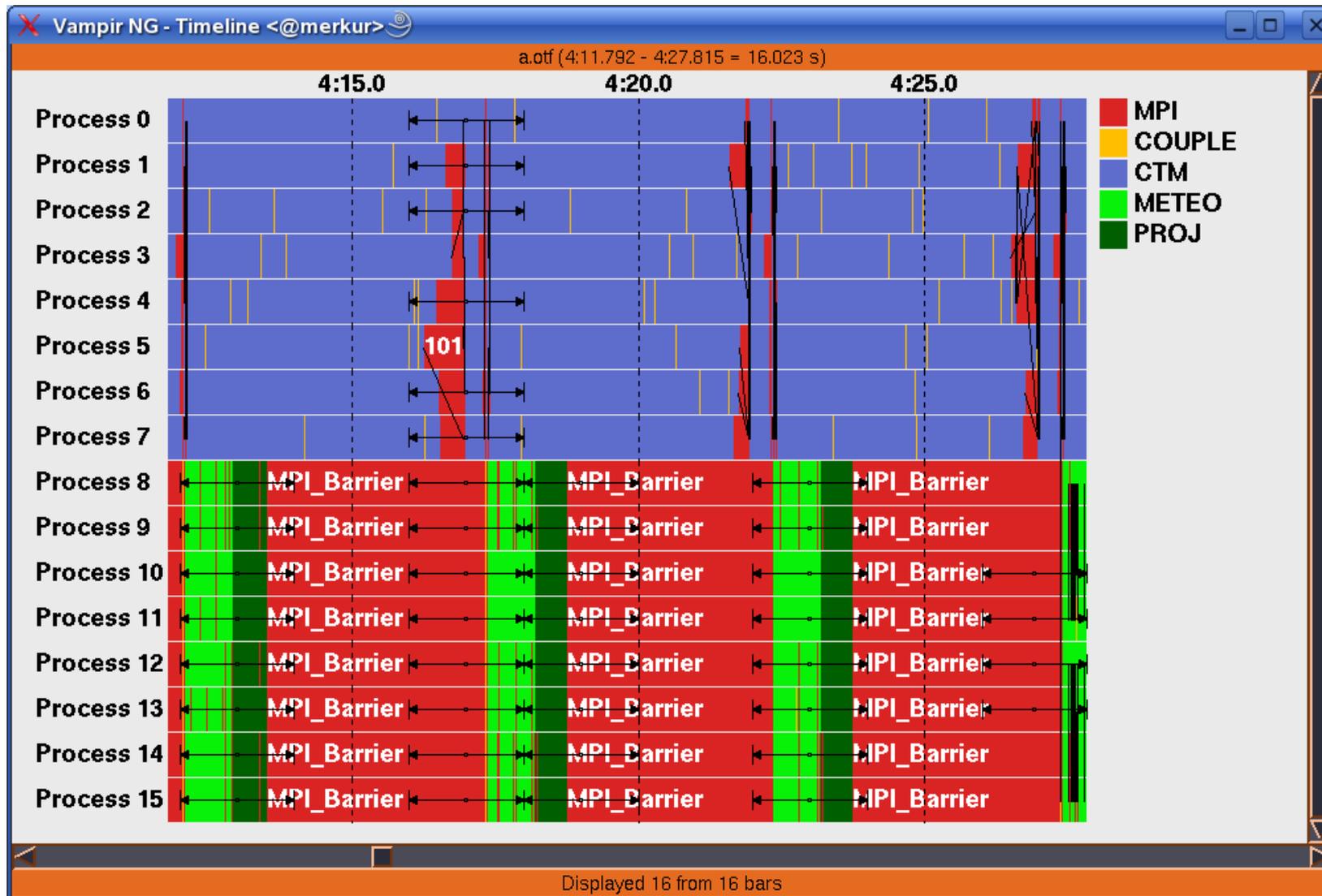
Summary Chart - Anteile der Gruppen



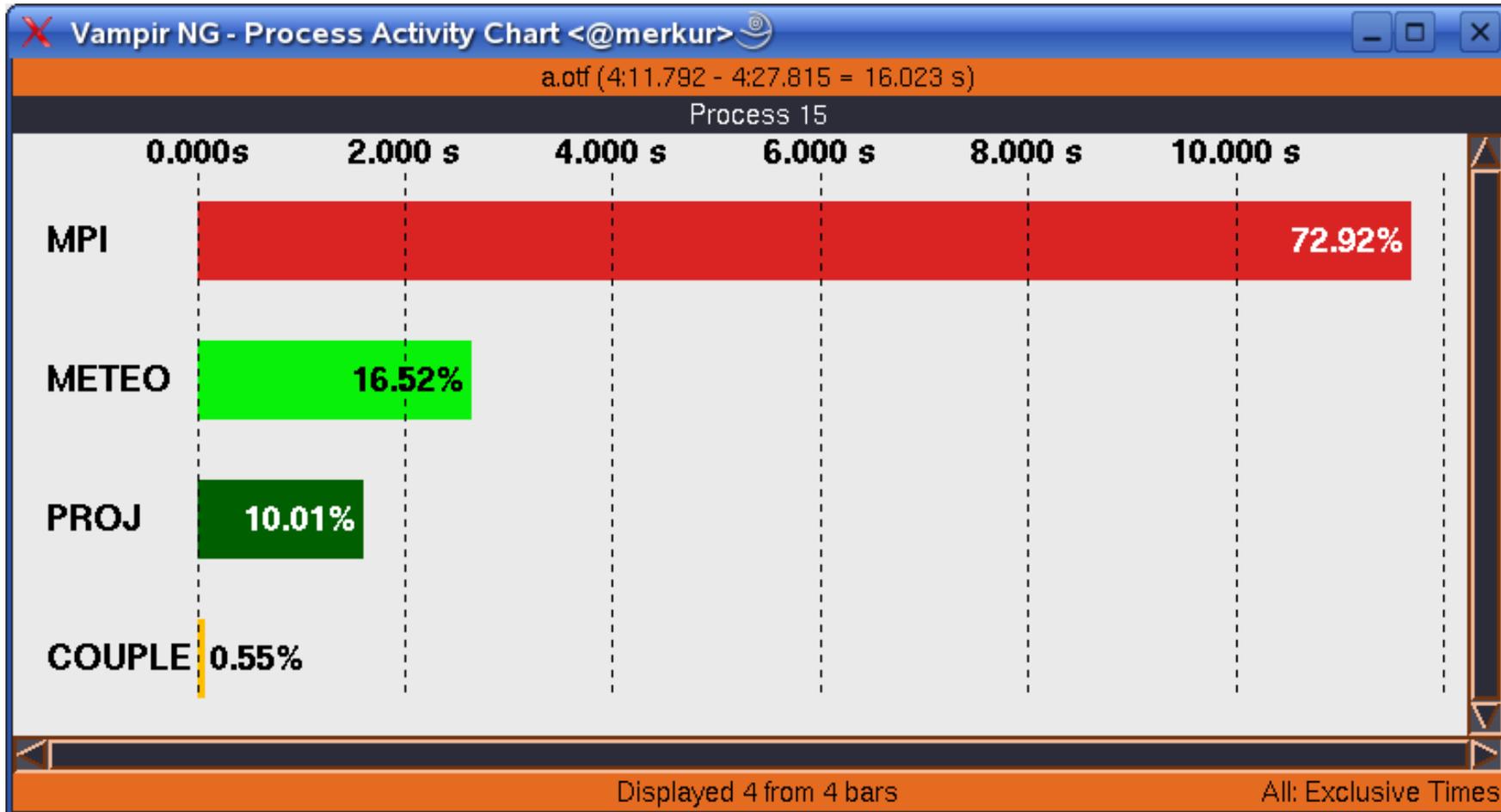
Timeline



Timeline - Ausschnitt



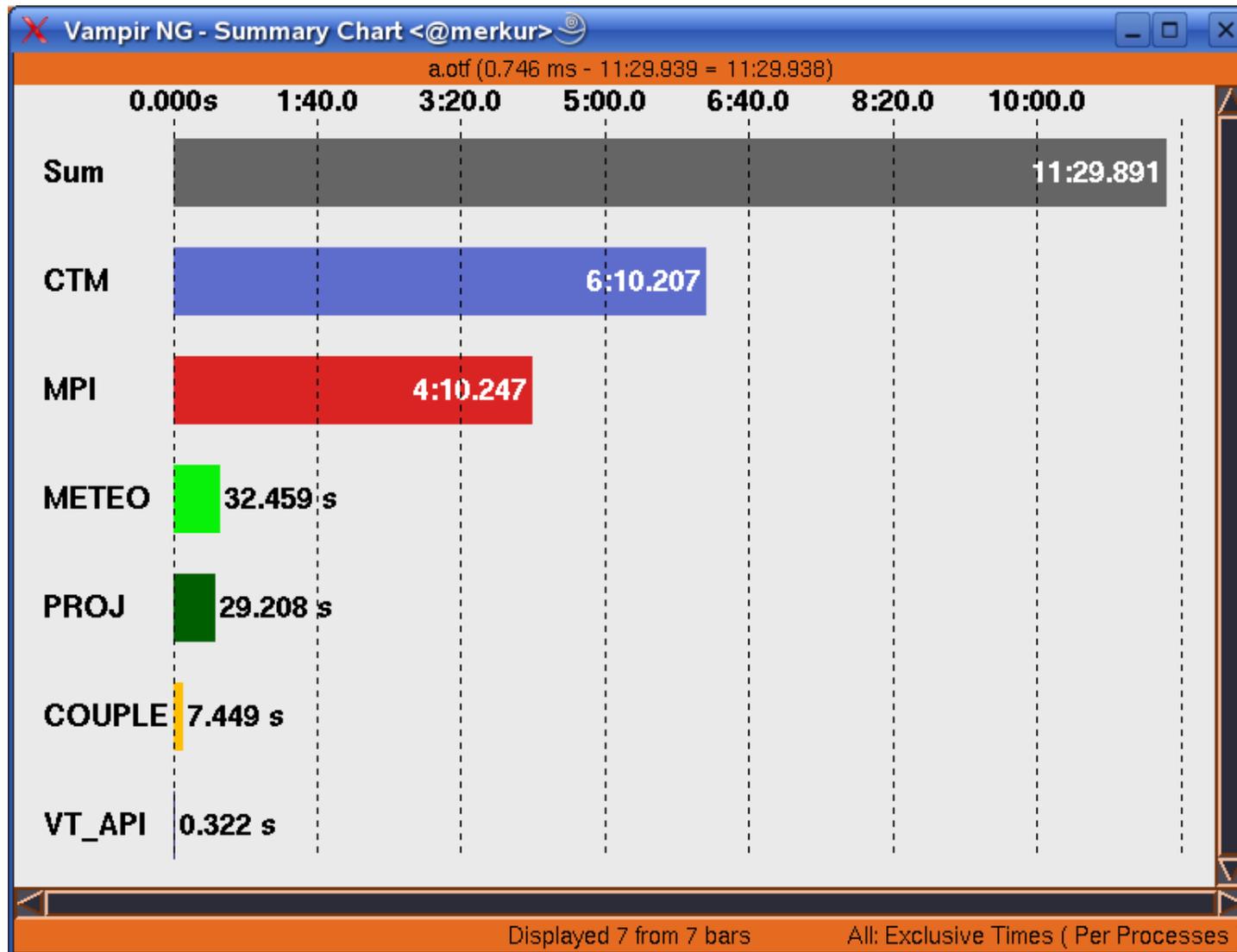
Process Activity Chart - Prozess 15 (METEO)



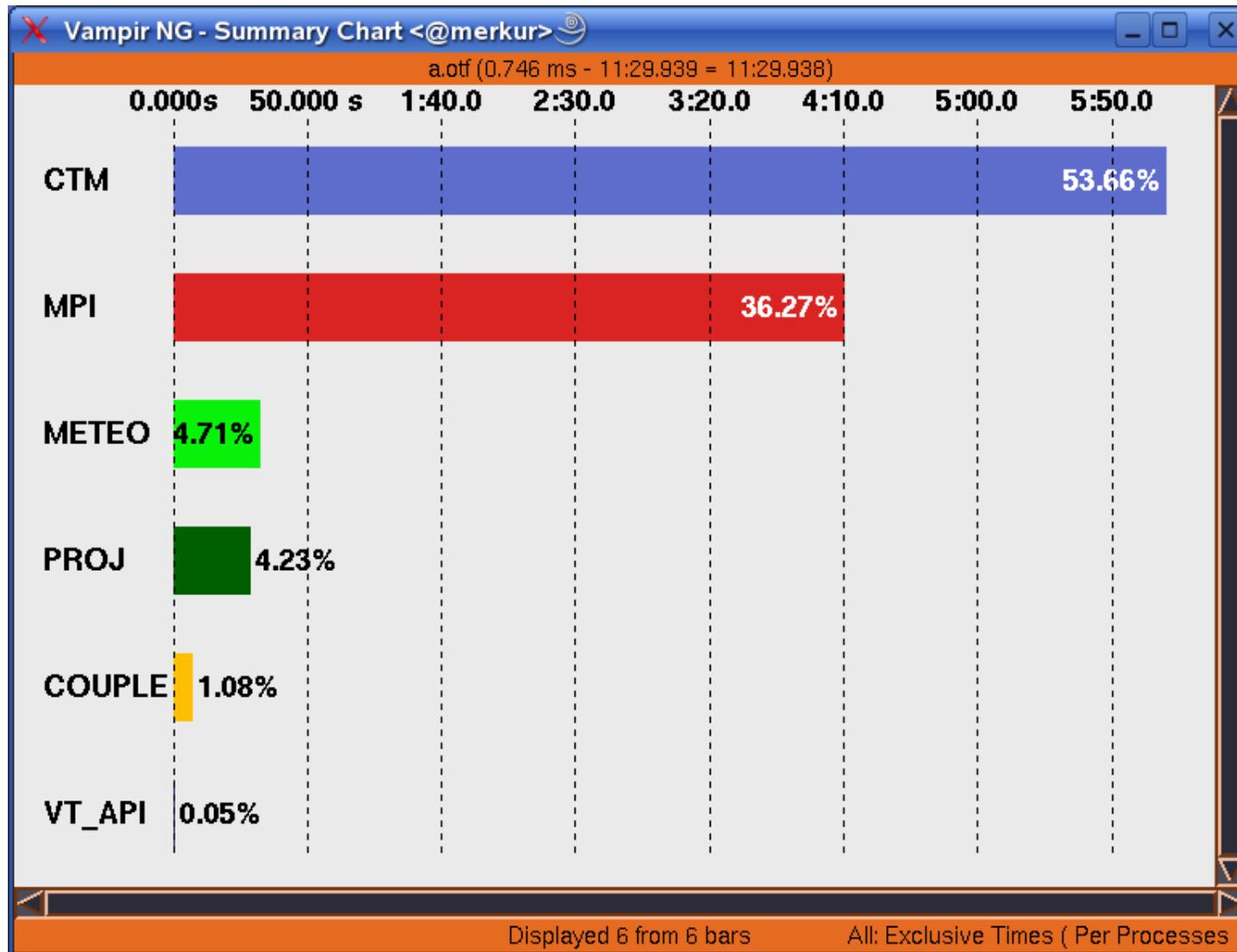
Diskussion

- CPUs von **METEO** haben die meiste Zeit gewartet
 - Load-Imbalance durch Kopplung der Modelle
 - Verbesserung: Das Verhältnis der **CTM**- zu **METEO**-CPUs optimieren!
- weniger **METEO**-CPUs, mehr **CTM**-CPUs
- z. B.: 4 CPUs für **METEO** und 12 CPUs für **CTM**
- Reduzierung der Laufzeit?

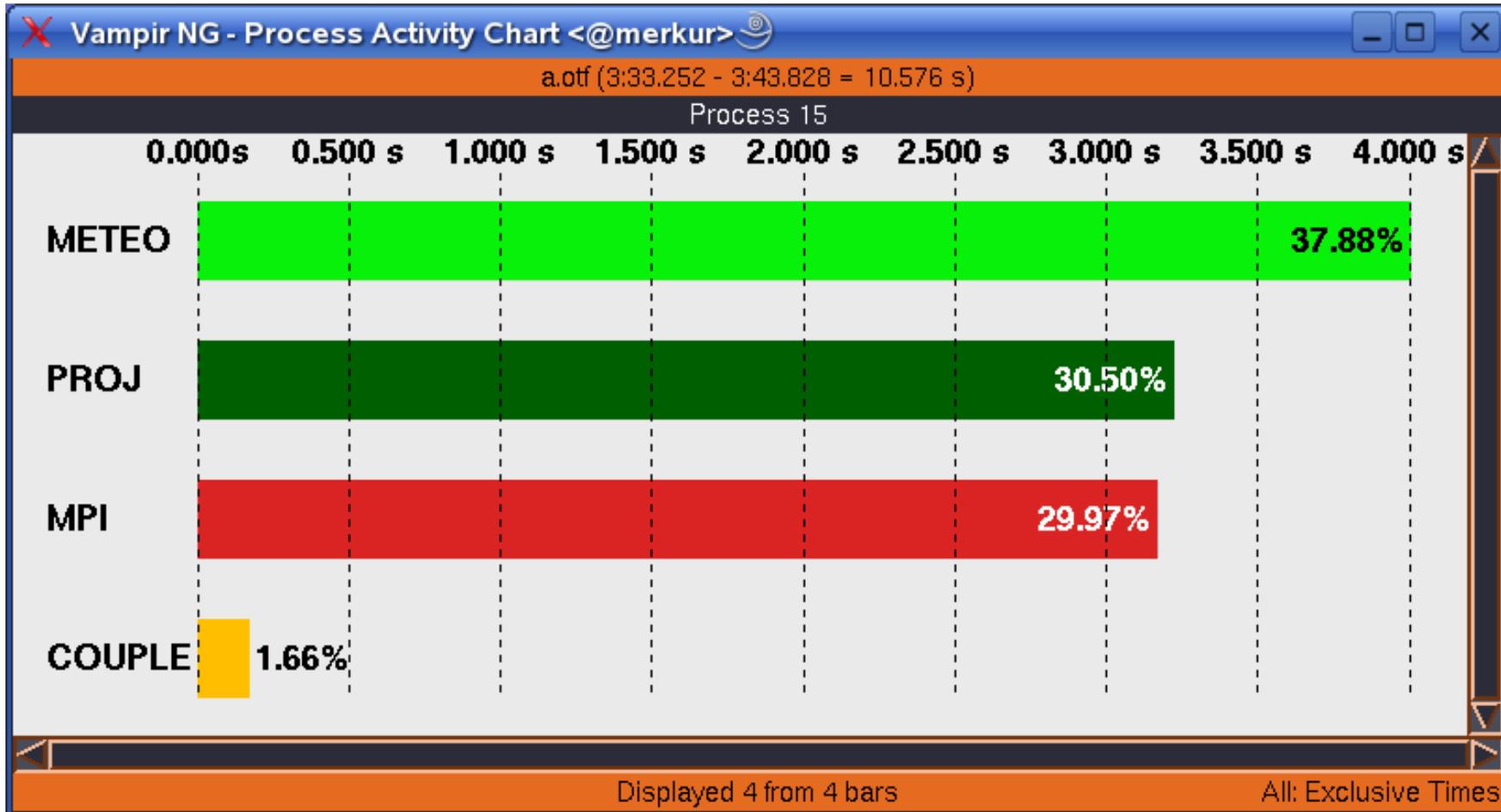
Summary Chart - Gesamtlaufzeit



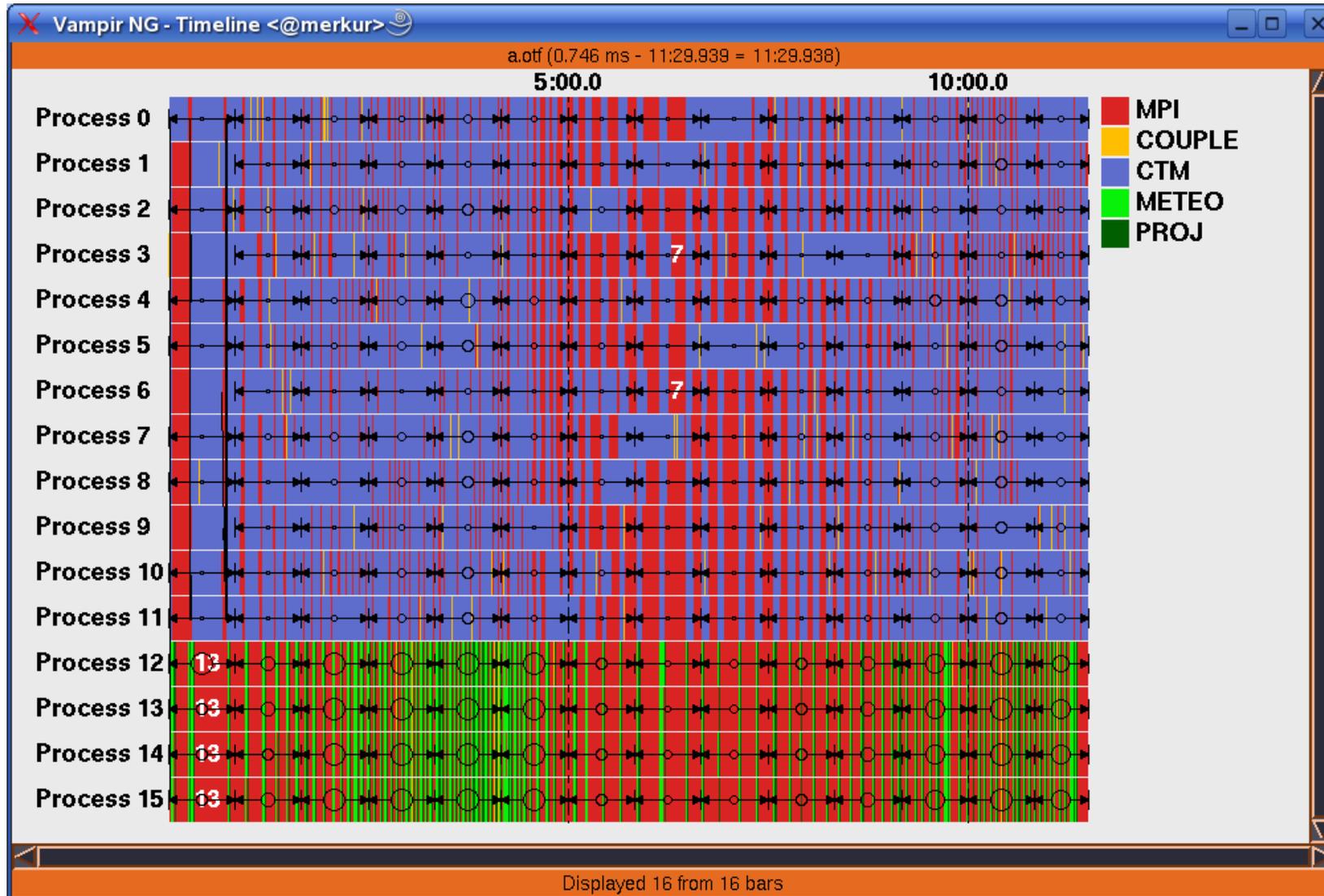
Summary Chart - Anteile der Gruppen



Process Activity Chart - Prozess 15 (METEO)



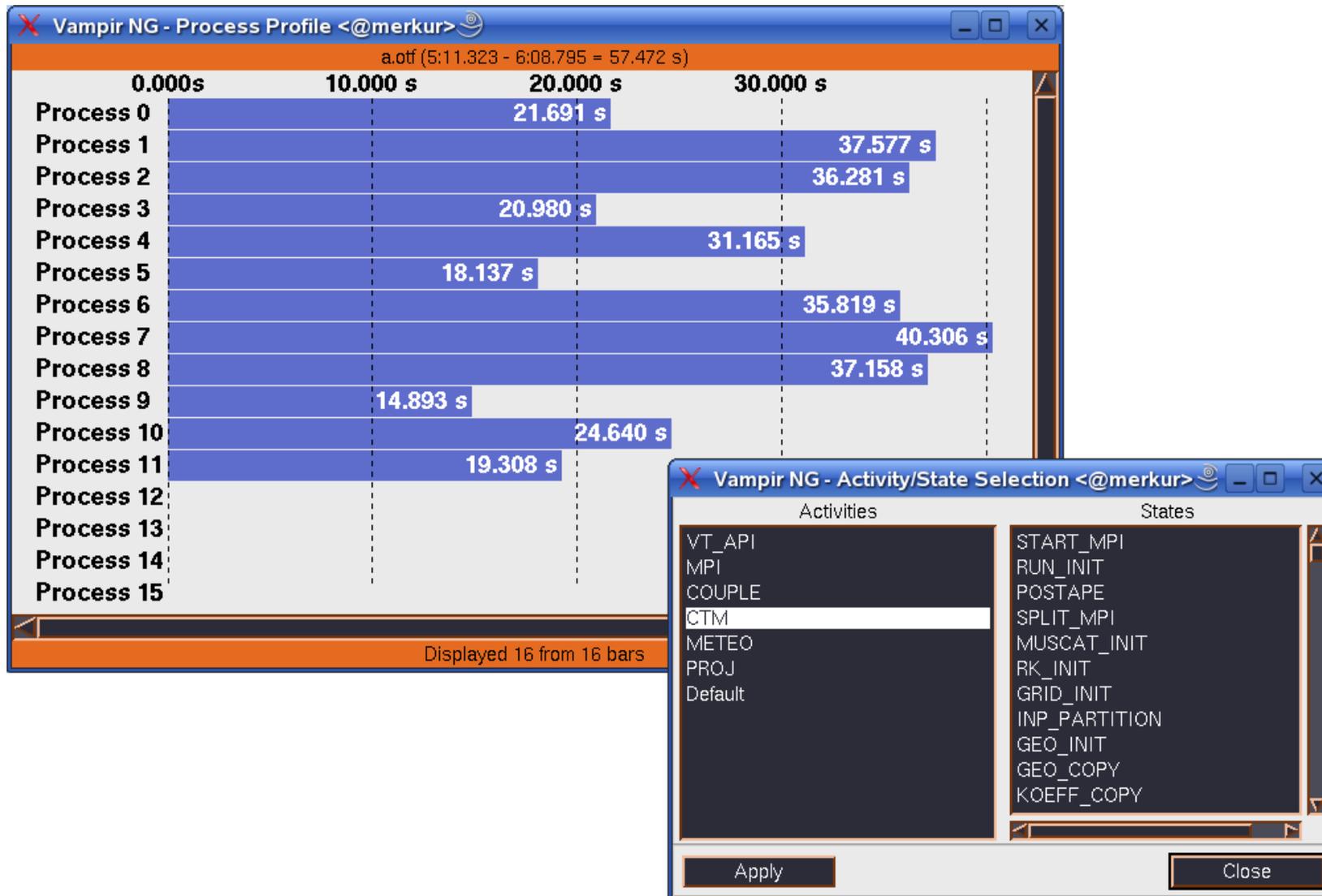
Timeline



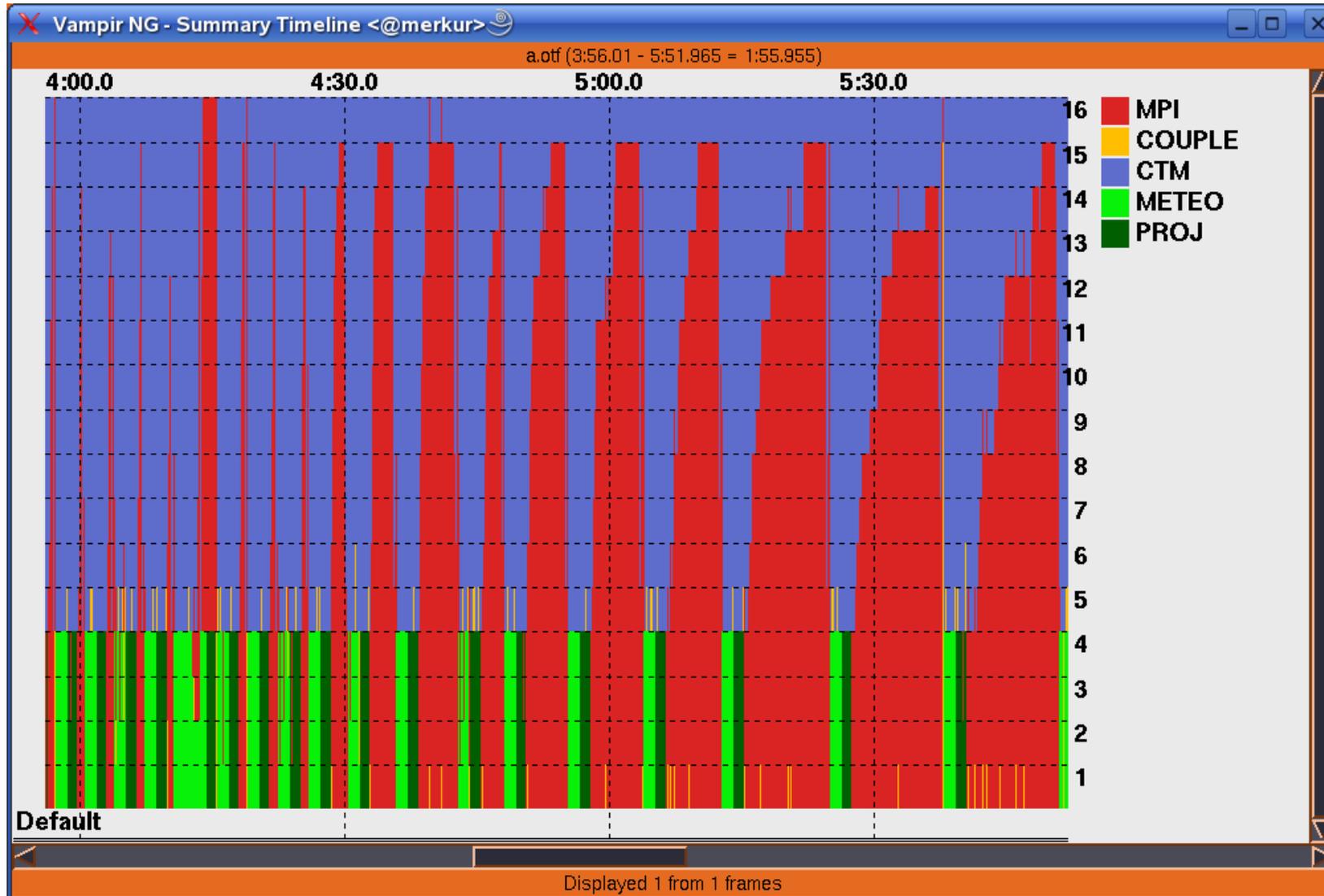
Timeline - Ausschnitt



Process Profile - CTM

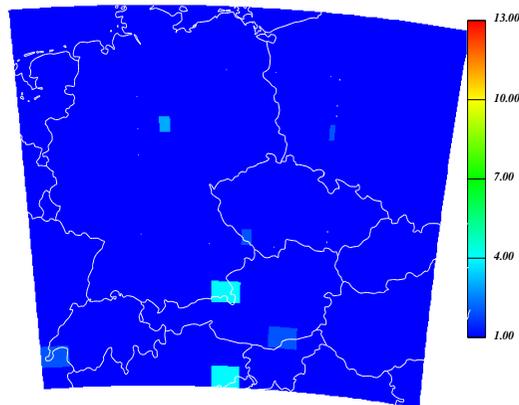


Summary Timeline

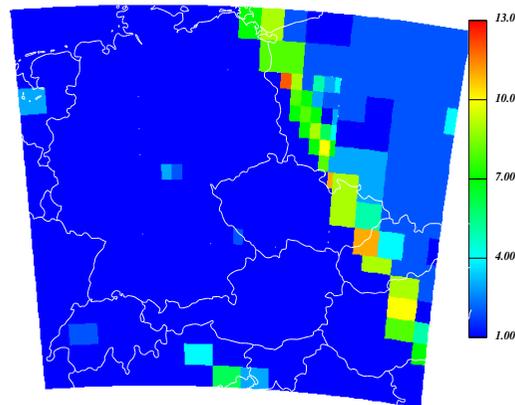


Ursachen der Load-Imbalance im CTM

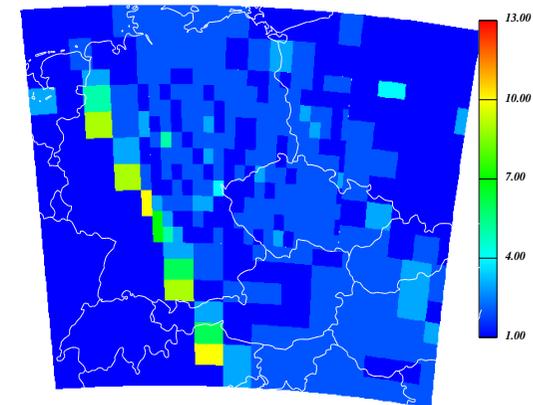
Starke räumliche und zeitliche Schwankungen des Rechenaufwands, in diesem Fall zu Sonnenaufgang:



Sat Aug 10 02:56:00 2002



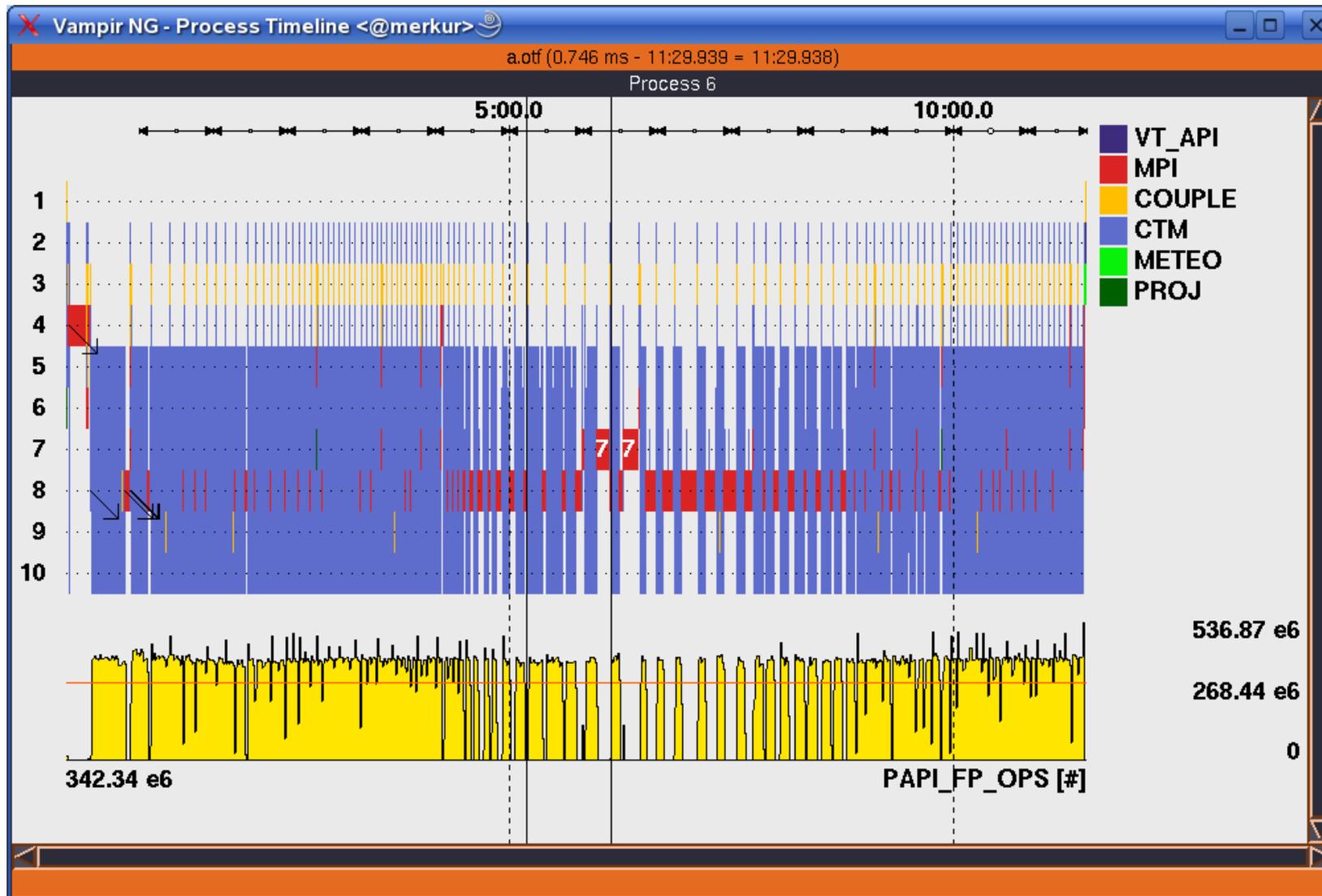
Sat Aug 10 03:35:59 2002



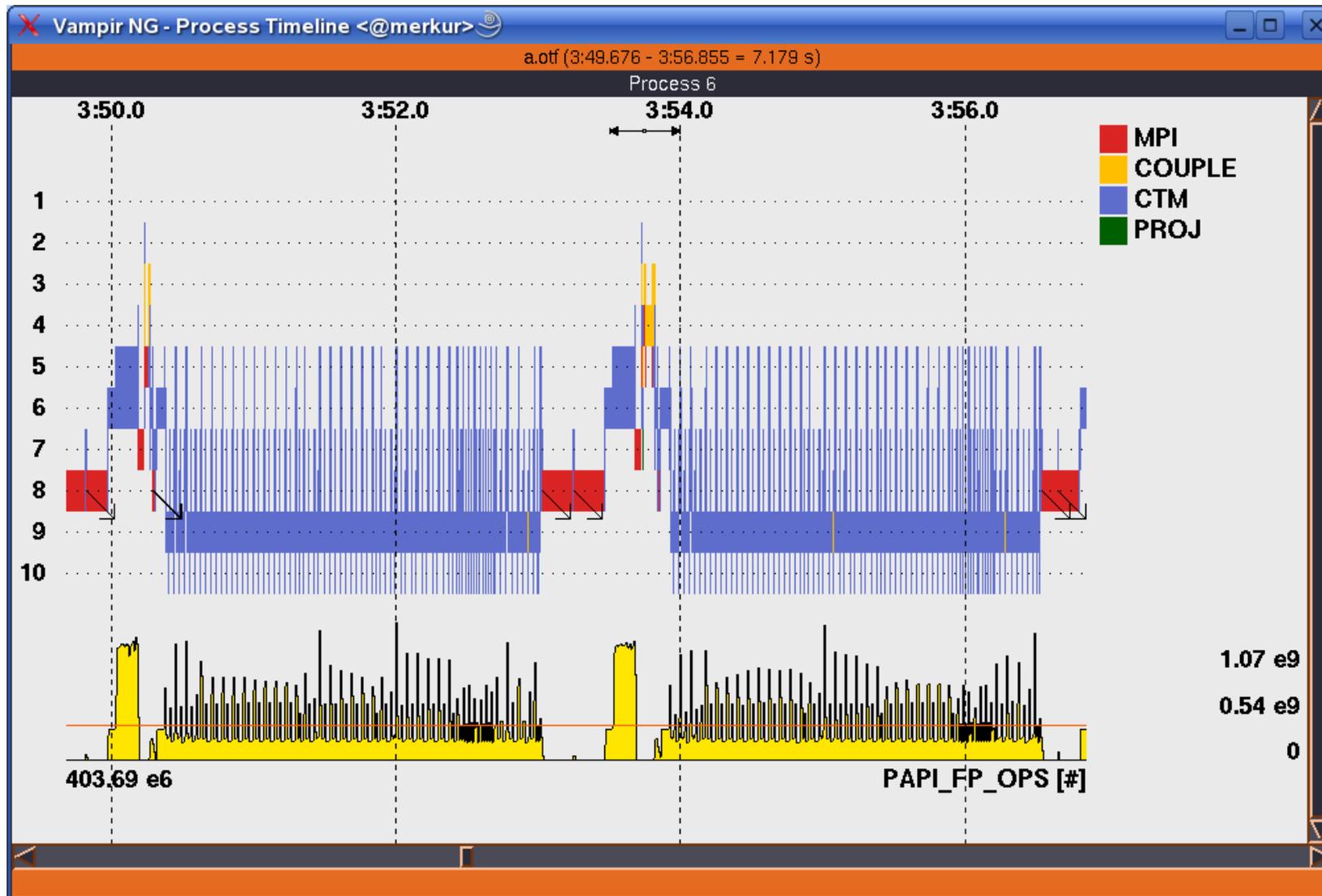
Sat Aug 10 04:09:20 2002

Räumliche Verteilung des Rechenaufwands zu 3 Zeitpunkten

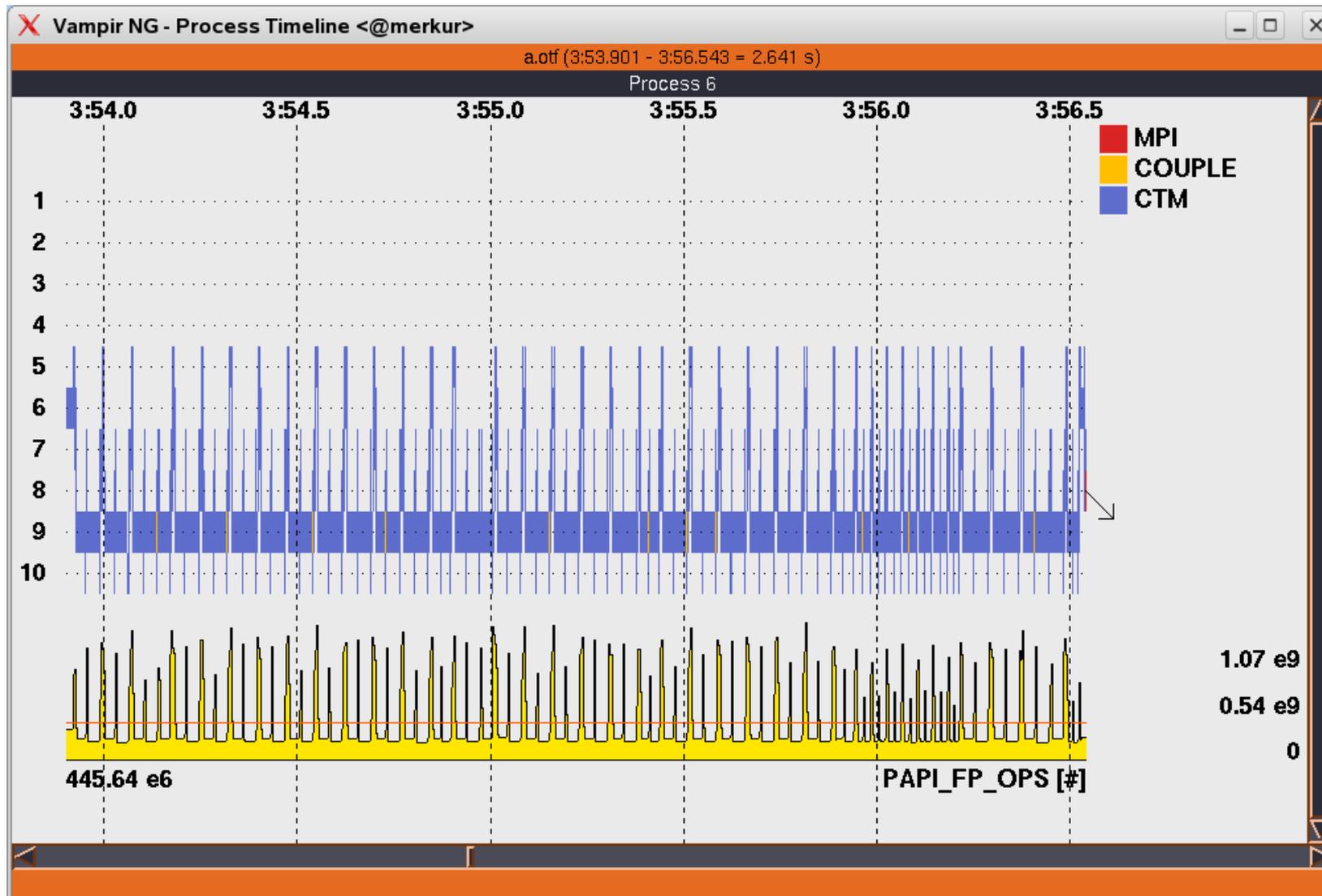
Process Timeline - Prozess 6 (CTM)



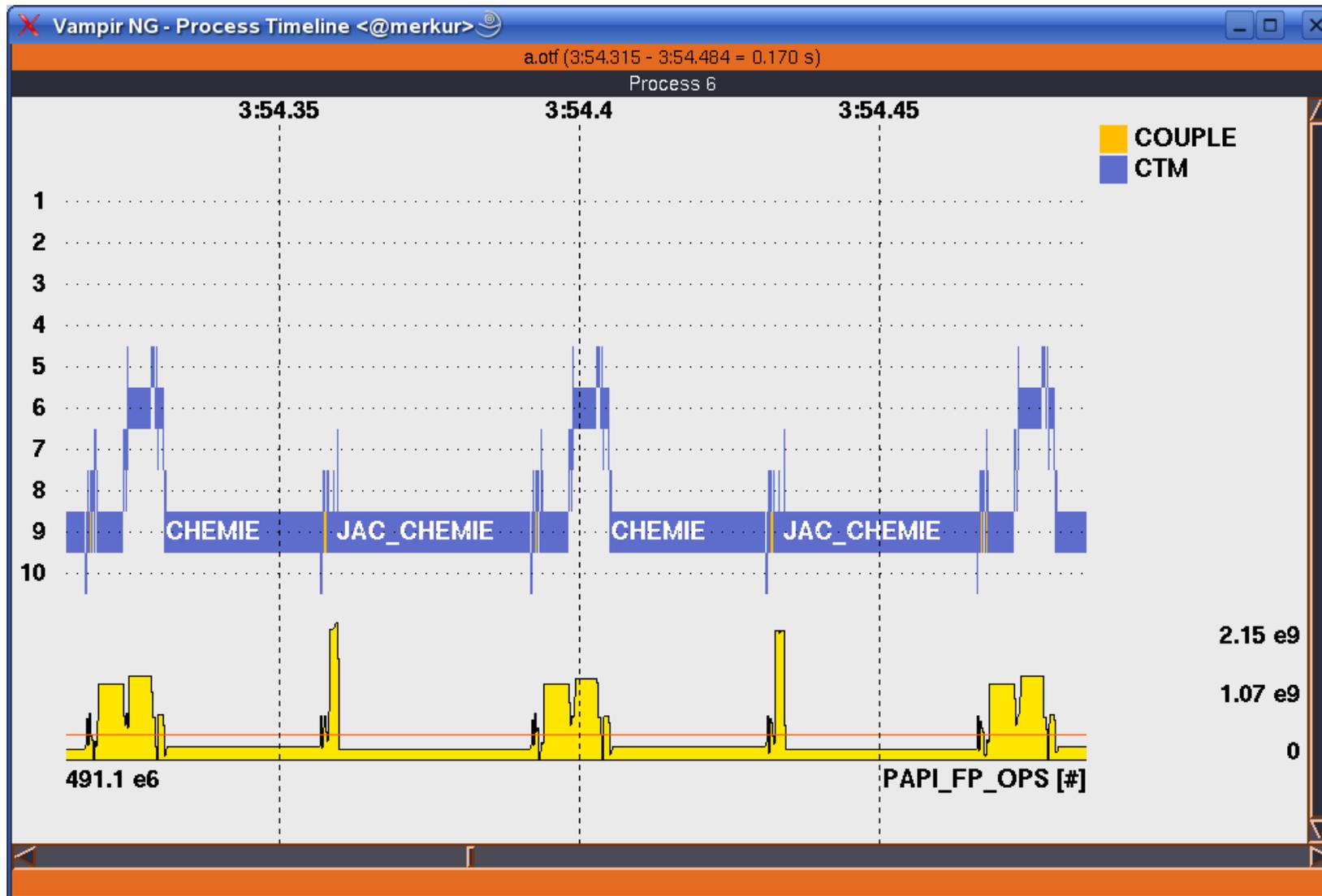
Process Timeline - Prozess 6 (CTM) - Zoom #1



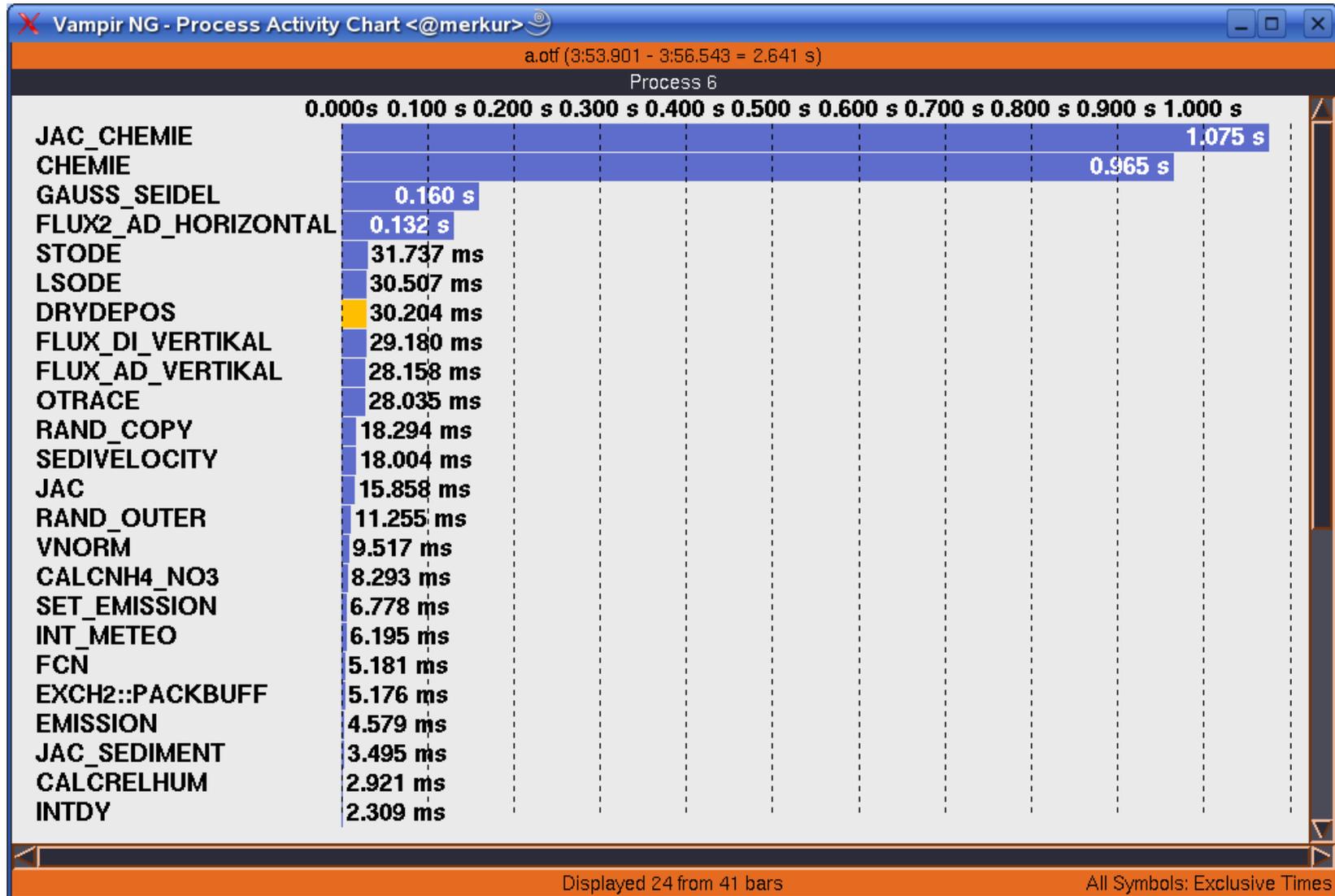
Process Timeline - Prozess 6 (CTM) - Zoom #2



Process Timeline - Prozess 6 (CTM) - Zoom #3



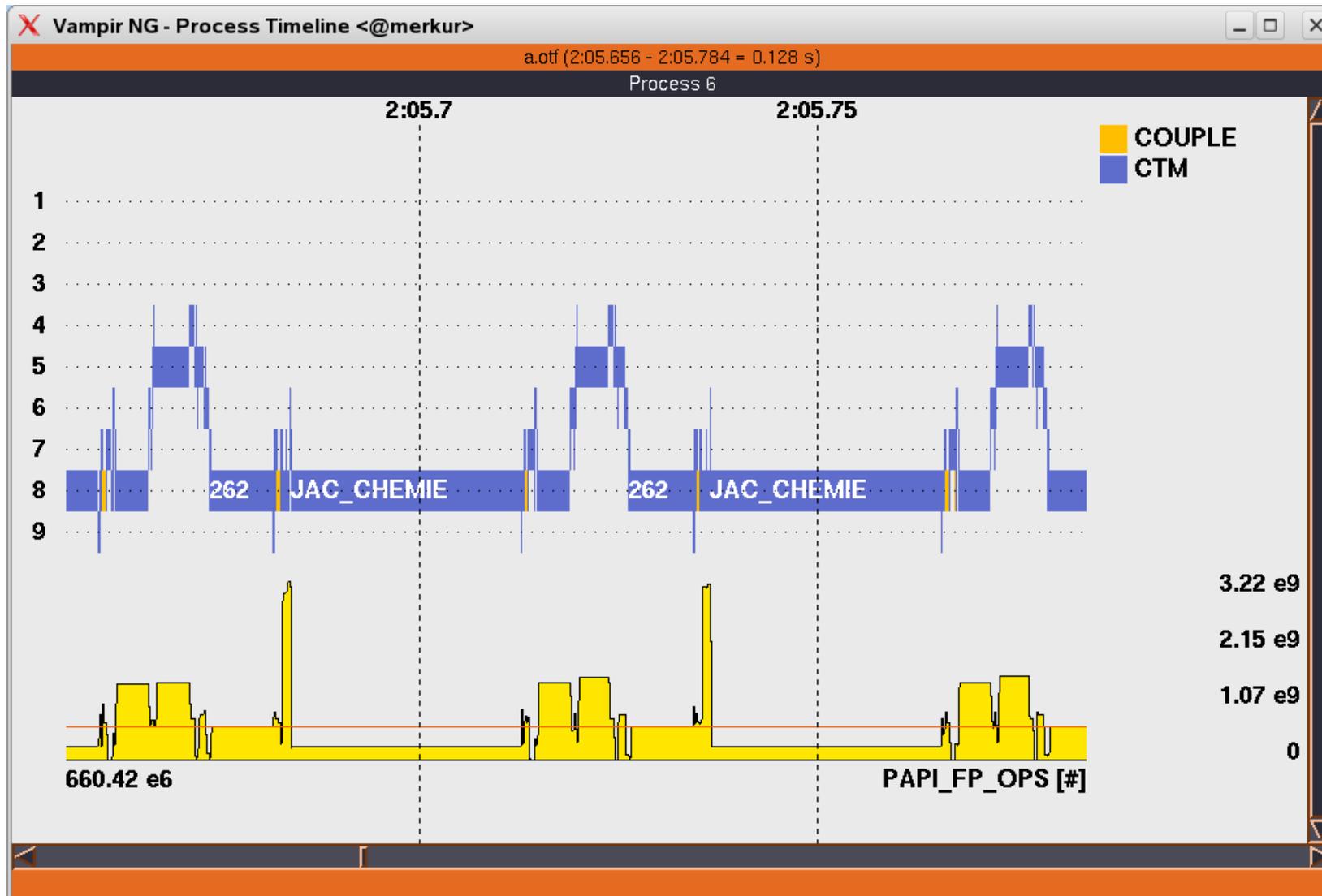
Process Activity Chart - Prozess 6 (CTM)



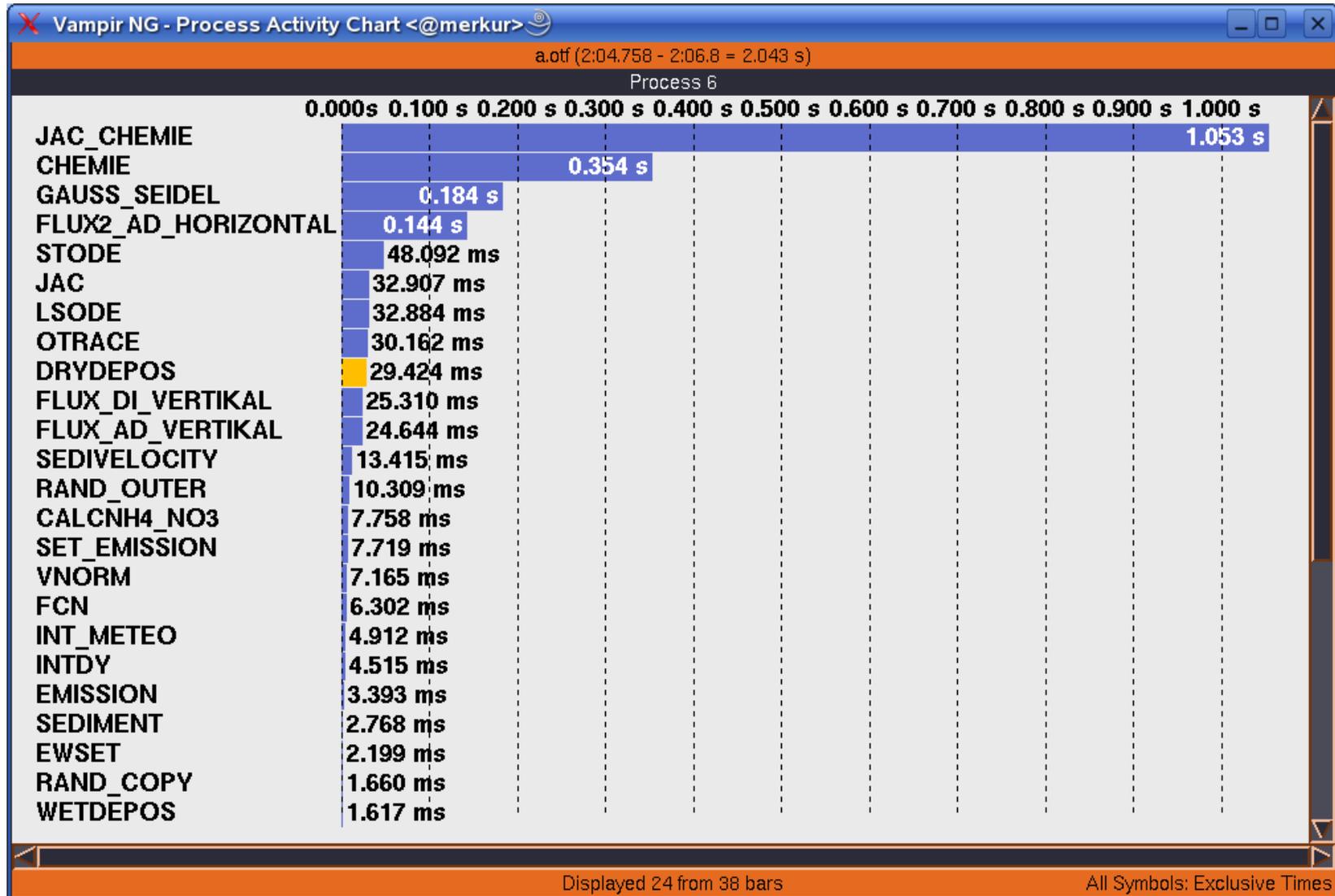
Diskussion

- Funktionen `chemie` und `jac_chemie` wenige FLOPS
- Jedoch sehr hohen Anteil an der Laufzeit
- Genaue Untersuchung der Ursachen notwendig!
 - Werden in diesen Funktionen viele floating point Operationen ausgeführt? → ja, d.h. Optimierung lohnenswert
 - Warum dann so wenige FLOPS? → Cache Misses einschalten, Quellcode studieren, ...
- Mit neuer Version 9.1 des Intel-Compilers...

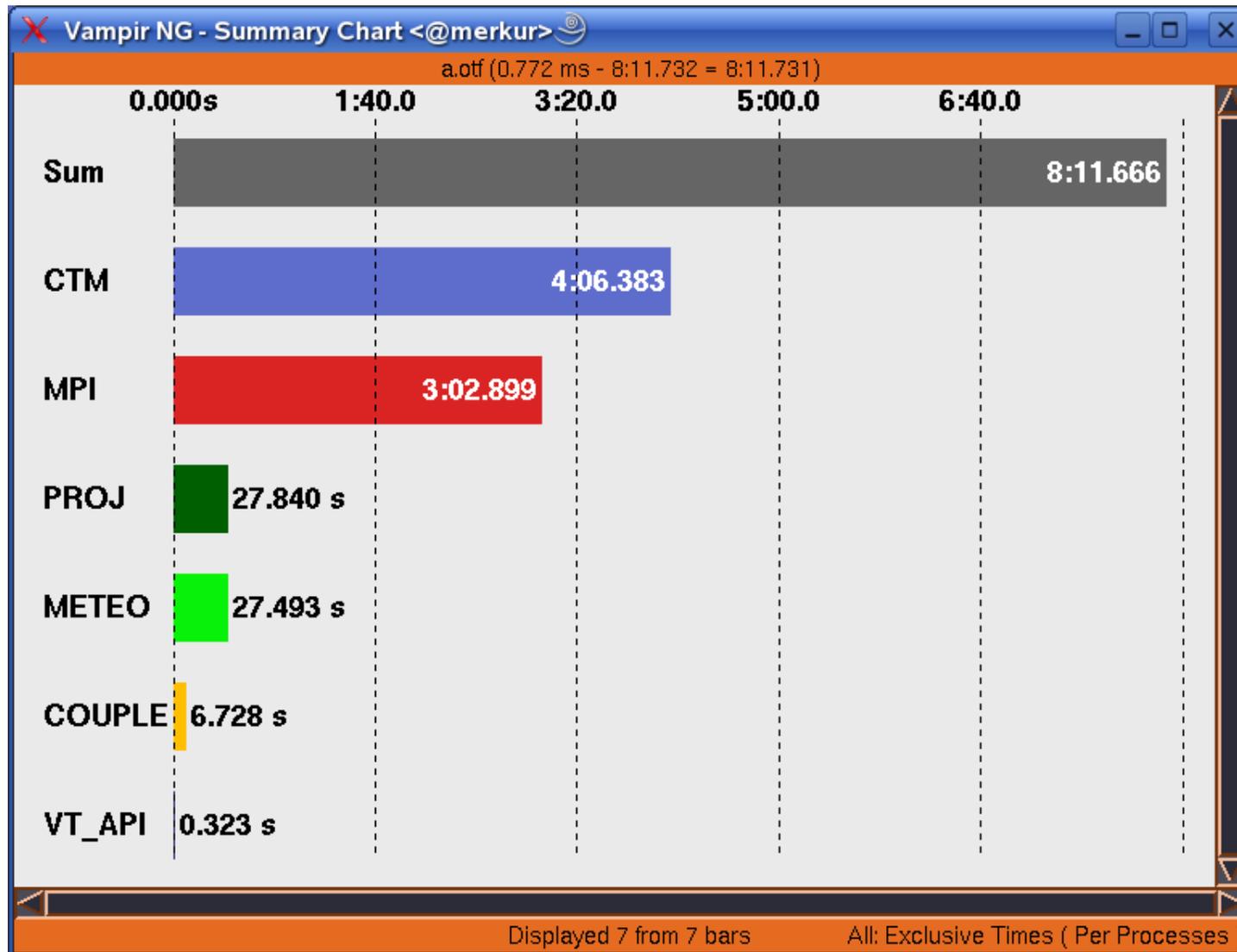
Process Timeline - Prozess 6 (CTM) - Zoom #3



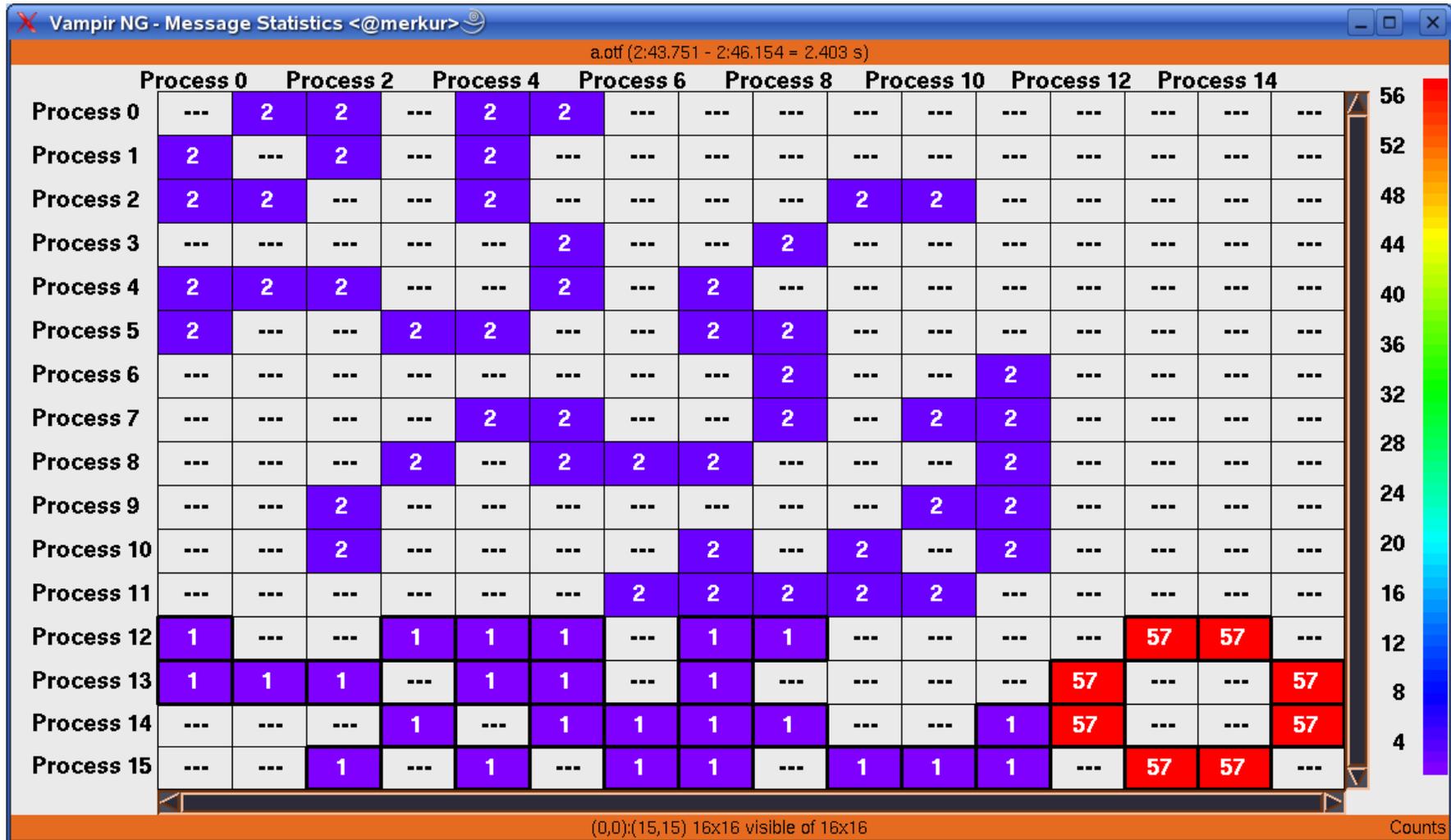
Process Activity Chart - Prozess 6 (CTM)



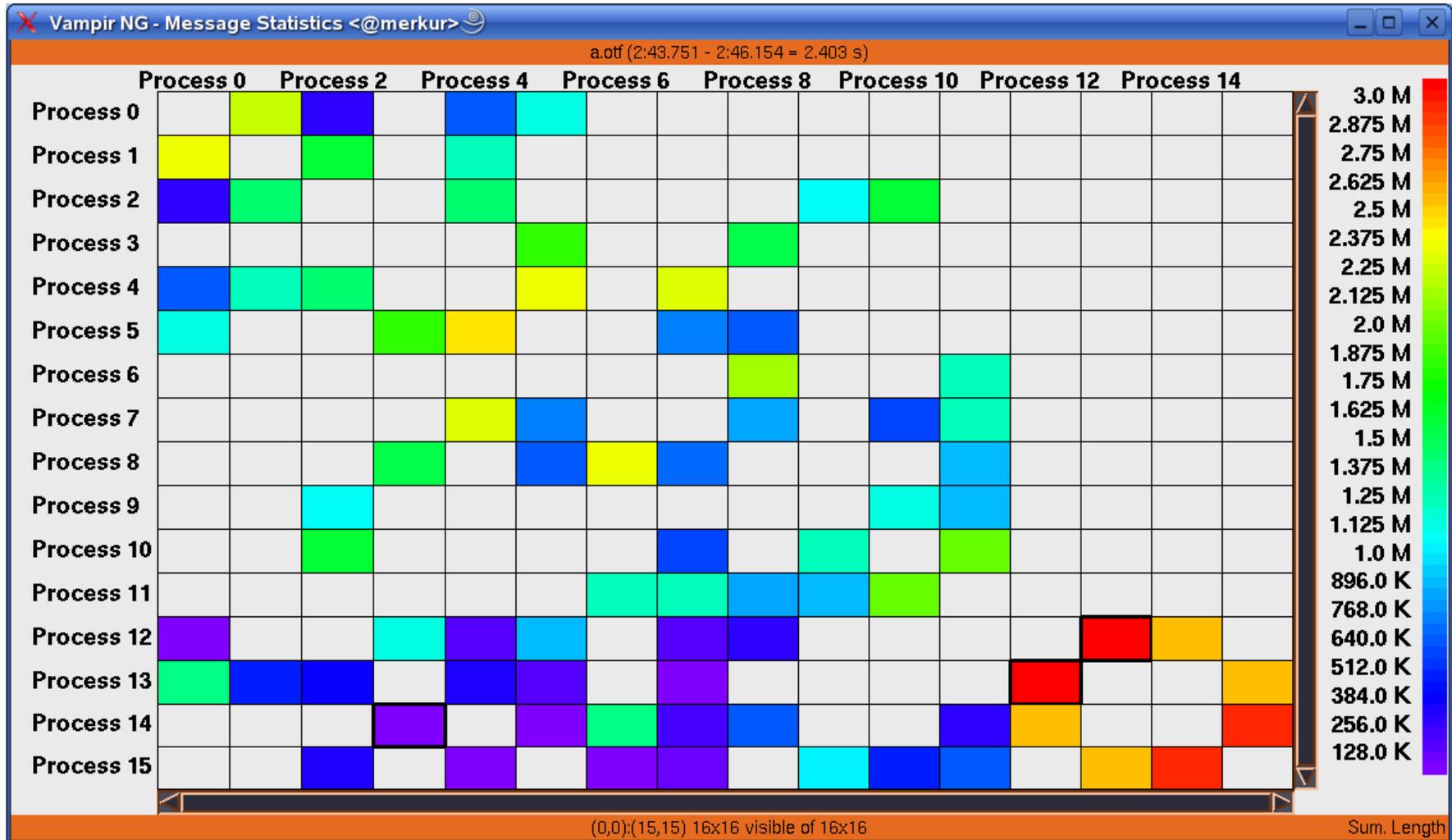
Summary Chart - Gesamtlaufzeit



Message Statistics - Anzahl der Nachrichten



Message Statistics - Summe der übertragenen Daten



VampirTrace Event Monitor

Software-Tuning on the Phobos PC-Farm

July 4th, 2006

Zellescher Weg 12

Willers-Bau A106

Tel. +49 351 - 463 - 31945

Matthias Jurenz (matthias.jurenz@tu-dresden.de)

Overview

- What is „VampirTrace“
- Instrumentation of user applications with VampirTrace
 - MPI function calls
 - Manually by Instrumentation API
 - Automatically by compiler
- Configuration
 - Region filter
 - Region grouping
 - Hardware counter
 - Environment variables

What is „VampirTrace“

The VampirTrace event tracing tool for serial und MPI parallelized applications produces tracefiles that can be analyzed with the Vampir performance analysis tool. It records all calls to the MPI library and all transmitted messages, and allows arbitrary user defined events to be recorded.

recordable events:

- MPI function calls
- transmitted messages
(point to point and collective operations)
- user function calls
- hardware counter

Instrumentation of user applications with VampirTrace

- MPI calls: Automatically by PMPI Wrapper Library
- User functions
 - Manually using VampirTrace Instrumentation API
 - Automatically by Compiler (PGI, Sunf90, IBM, GNU)

Instrumentation – MPI calls (1)

- The MPI profiling interface is a part of the MPI standard specification
- Each MPI-routine call has a corresponding PMPI-call with identical syntax and functionality

`MPI_Send(...)` \Leftrightarrow `PMPI_Send(...)`

- This interface can be used to intercept all MPI-calls with wrapper routines:

```
int MPI_Send(...) {  
    int rc;  
    log_start_of_send(...);  
    rc = PMPI_Send(...);  
    log_end_of_send(...);  
    return rc;  
}
```

Instrumentation – MPI calls (2)

- Relink MPI application replacing `-lmpi` by

C/C++:

```
-L$(VT_LIB) -lelg.otf.mpi -lotf -lmpi
```

Fortran:

```
-L$(VT_LIB) -lfmpi -lelg.otf.mpi -lotf -lmpi
```

- No recompilation necessary

Instrumentation – MPI calls (3)

- Advantages

- All MPI function calls and communication will be recorded
- No recompilation necessary

- Disadvantages

- **No user functions or source code regions will be recorded**

Instrumentation – Manually by Instrumentation API (1)

- Include header file for using VampirTrace instrumentation API

C/C++: `#include "elg_user.h"`

Fortran: `#include "elg_user.inc"`

- Place instrumentation macros/function calls at start and end of every function or region

C: `ELG_USER_START("function-name");`

`ELG_USER_END("function-name");`

C++: `ELG_TRACER("function-name");`

Fortran: `call ELG_USER_START('function-name')`

`call ELG_USER_END('function-name')`

Instrumentation – Manually by Instrumentation API (2)

- Compile manually instrumented application with
`-I$(VT_INC) -DEPILOG`
- Without `-DEPILOG` the `ELG_*` calls expand to nothing and will be ignored
- Fortran source files have to be **preprocessed**

Instrumentation – Manually by Instrumentation API (3)

- Link manually instrumented application as follows:

Sequential (no MPI):

```
-L$(VT_LIB) -lelg.otf -lotf
```

MPI:

```
-L$(VT_LIB) (-lfmpi) -lelg.otf.mpi -lotf -lmpi
```

Instrumentation – Manually by Instrumentation API (4)

- Advantages

- Independent of used compiler
- Also source code regions are instrumentable (e.g. loops)
- Only functions/regions will be recorded, which are interesting for user

- Disadvantages

- Error-prone
(e.g. corrupt call stack due to missing return instrumentation call)
- **Very hard instrumentation of large source codes**

Instrumentation - Automatically by Compiler (1)

- Uses the compiler profiling interface
(e.g. PGI -> compiler flag: `-Mprof=func`)
- In the commands to build the application (e.g. in a Makefile), the user should precede all compile and link commands with „`kinst`“

Example:

Instead of the command

```
% mpif90 myprog1.f90 myprog2.f90 -o myprog
```

the command

```
% kinst mpif90 myprog1.f90 myprog2.f90 -o myprog
```

has to be executed

Instrumentation - Automatically by Compiler (2)

- Advantages

- **Very simple way of instrumentation**

- Disadvantages

- Depends on used compiler
(supported compiler: PGI, Sunf90, IBM, GNU)
- Tracefiles can get very large because **all** user functions will be instrumented

Configuration – Region Filter (1)

- Tracefile written by VampirTrace can easily get **very large**
 - => The powerful filter management of VampirTrace is very useful to reducing size of tracefile, which will be generated
- The user can decide how often an instrumented region has to be recorded (call limit) to tracefile
- Summary of filtered regions will be written at the end of tracing

Configuration – Region Filter (2)

- Create a file which contains filter directives in following format:

```
<regions> -- <limit>
```

regions semicolon-separated list of regions (wildcards allowed)

call limit assigned call limit (per CPU)

0 = region(s) denied

-1 = unlimited

n = set call limit to *n*

- Set path of filter configuration file by environment variable
ELG_FILTER_SPEC

Configuration – Region Filter - Example

- Objectives:

- Limit recording of function „search()“ to 100
- Deny recording of all remaining functions except MPI calls

- Create a file „filter.spec“ with following content:

```
search -- 100
```

```
MPI_* -- -1
```

```
* -- 0
```

- Set environment variable ELG_FILTER_SPEC to “filter.spec”

```
export ELG_FILTER_SPEC=./filter.spec
```

Configuration – Region Grouping (1)

- The region grouping management allows the user to assign functions/regions to a group
- Very useful to summarize program parts
- Vampir allows different colors for groups (activities)
- MPI functions will be assigned to group „MPI“ and remaining functions to group „USR“ by default

Configuration – Region Grouping (2)

- Create a file, which contains grouping assignments in following format:

```
<group>=<regions>
```

group group name

regions semicolon-separated list of regions (wildcards allowed)

- Set path to group configuration file by environment variable
ELG_GROUPS_SPEC

Configuration – Region Grouping - Example

- Objectives:

- The functions „add()“, „sub()“, „mul()“ and „div()“ should be assigned to group „CALC“
- All remaining functions should be assigned to group „Application“

- Create a file „groups.spec“ with following content:

```
CALC=add;sub;mul;div
```

```
Application=*
```

- Set environment variable ELG_GROUPS_SPEC to “groups.spec”

```
export ELG_GROUPS_SPEC=./groups.spec
```

Configuration – Hardware Counter

- If VampirTrace has been built with enabled hardware-counter support, VampirTrace is capable of recording hardware counter information
- Uses the PAPI hardware-counter library
- Will be recorded by each enter/leave region event
- The user must set the environment variable `ELG_METRICS`, which should contain a colon-separated list of PAPI counter names

- Example:

```
export ELG_METRICS=PAPI_FP_OPS:PAPI_TOT_INS
```

Configuration – Environment Variables (1)

<u>Variable</u>	<u>Meaning</u>
ELG_PFORM_GDIR	Name of global, cluster-wide directory to store final trace file
ELG_PFORM_LDIR	Name of node-local directory, that can be used to store temporary trace files
ELG_FILE_PREFIX	Prefix used for trace file names
ELG_BUFFER_SIZE	Size of internal event trace buffer in bytes
ELG_VERBOSE	Print VampirTrace related control information during measurement

Configuration – Environment Variables (2)

Variable

Meaning

ELG_METRICS

Specify counter metrics to be recorded with trace events as a colon-separated list of PAPI counter names

ELG_FILTER_SPEC

Name of file specifying region filter directives

ELG_GROUPS_SPEC

Name of file specifying region grouping assignments