# Some Results Obtained

Hendryk Bockelmann, DKRZ
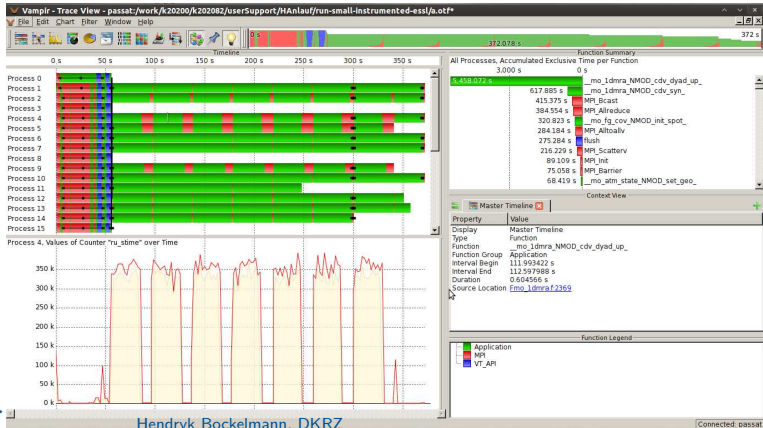
*"while testing some code that ran perfectly on our old POWER5, I noticed a very bad performance on blizzard. rusage/hpccount gave the following output:"*

```
Metrik              [Unit]  :     Average     Minimum     Maximum
maximum RSS         [Kbytes]: 1121529.67  1059896.0  2015480.0
time in user mode   [sec]   :      804.57      662.59      903.04
time in system mode [sec]   :     2703.09     2608.43     2860.22
inst per run cycle          :        0.42        0.42        0.43
peak performance    %   :        0.02        0.01        0.03
```

**DKRZ**
DEUTSCHES
KLIMARECHENZENTRUM

Hendryk Bockelmann, DKRZ

# Locate high sys time

- instrument code with vampirtrace and use `VT_RUSAGE=all`
- high sys time within module `mo_1dmra` localized in subroutine `cdv_dyad_up`
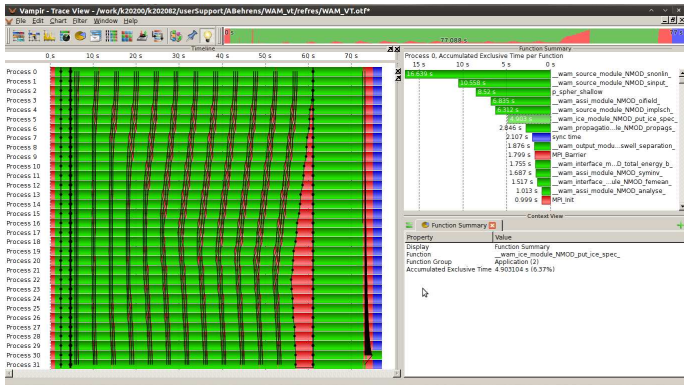


Hendryk Bockelmann, DKRZ

# Solution

- MATMUL uses multithreading by default:

  *"The default value for* `num_threads` *when using the* `MATMUL` *intrinsic equals the number of processors online. Changing the number of threads available to the* `MATMUL` *and* `RANDOM_NUMBER` *intrinsic procedures can influence performance."*
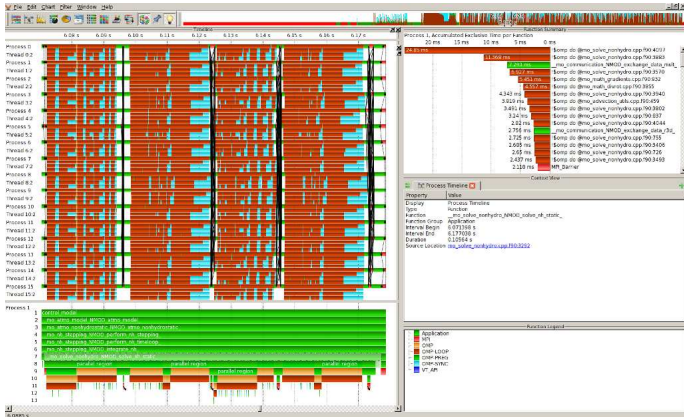
- set XLFRTEOPTS="intrinthds=1"

# MPI barrier slows down ...



solution: cpu-load due to additional ICE computations propagate with each nearest neighbour exchange

DKRZ
DEUTSCHES
KLIMARECHENZENTRUM

Hendryk Bockelmann, DKRZ

# OpenMP schedule optimized

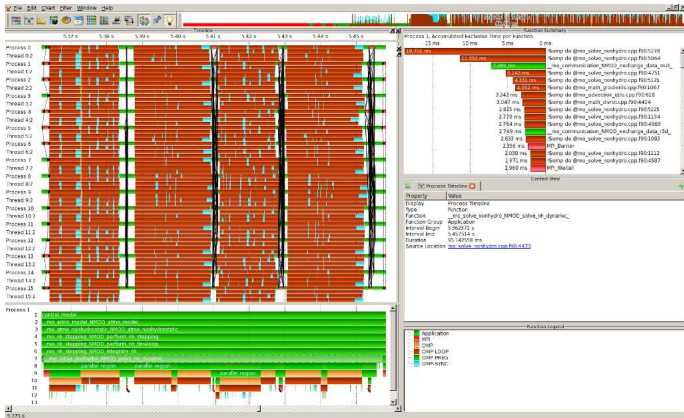what is the best setup of OpenMP threads in a hybrid code?



nproma=16, nchunksize=2, static schedule (default)

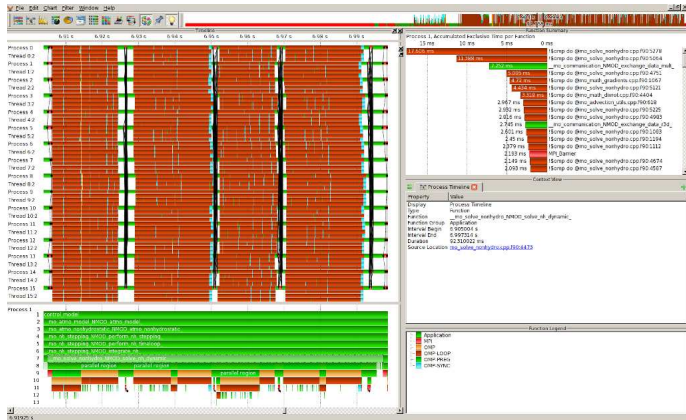Hendryk Bockelmann, DKRZ

# OpenMP schedule optimized



nproma=16, nchunksize=1, dynamic schedule

# OpenMP schedule optimized



nproma=4, nchunksize=1, dynamic schedule

ScalES

# **Program Analysis and Tuning Workshop - DKRZ 2012**

Jörg Behrens
Deutsches Klimarechenzentrum GmbH

Panagiotis Adamidis, Hendryk Bockelman,
Thomas Jahns (DKRZ)

DKRZ

# Example: ECHAM performance

In the ScalES project we looked at how model performance scales with the number of tasks.
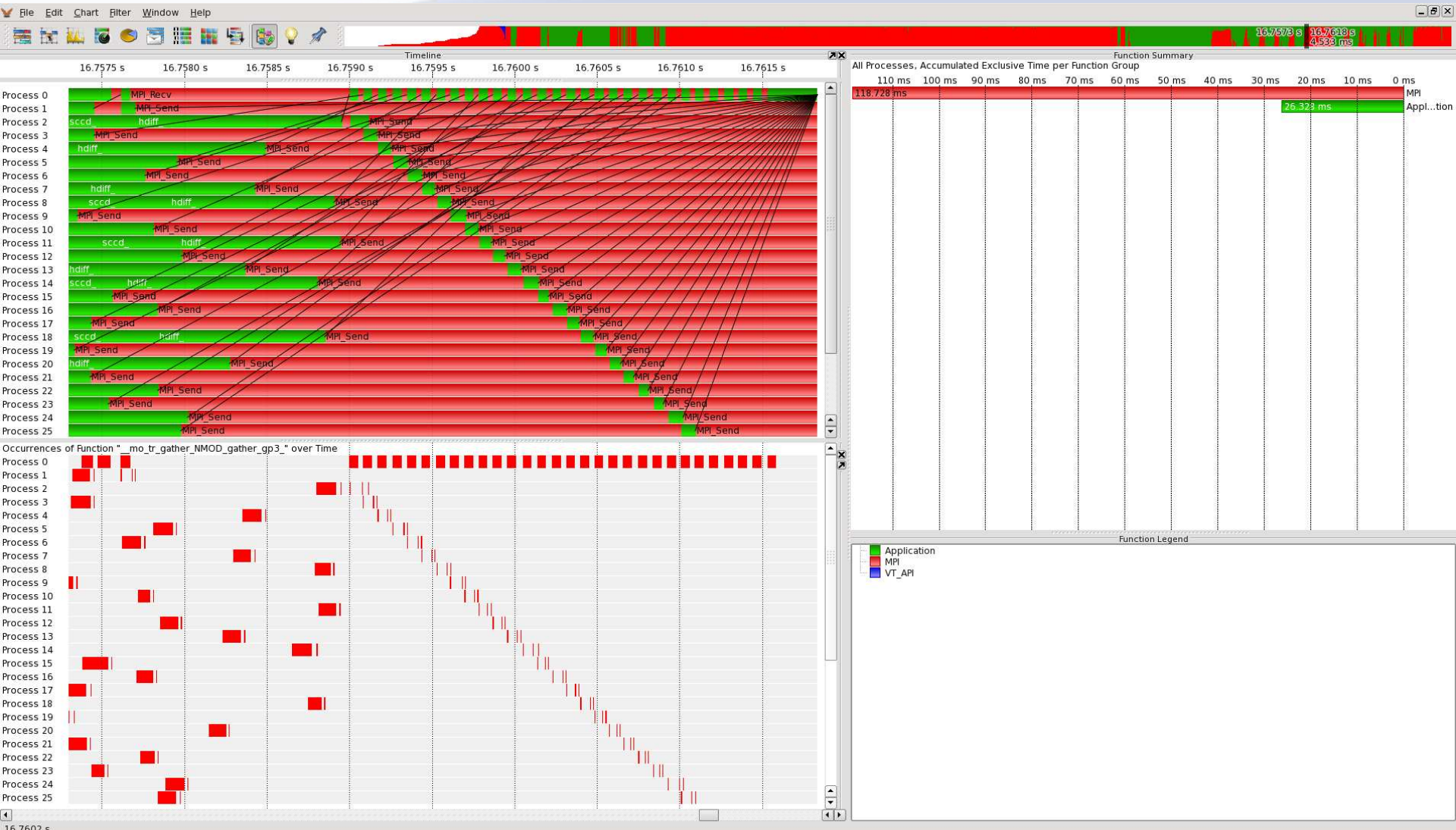
Of course, the worst part must be the serial part, e.g. serial IO.

But it still can be worse than that:

   identified bottleneck: hand-made p2p-gather

   becomes slower for #task > 512 (ECHAM T127L95)

p2p-gather required because the MPI_GATHERV is not general enough
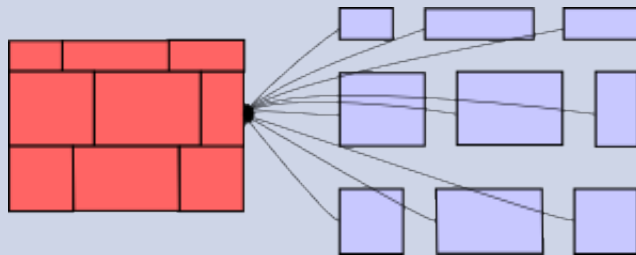
DKRZ, 2012

DKRZ

# Vampir view at p2p-gather

DKRZ, 2012

# Insertion of barriers
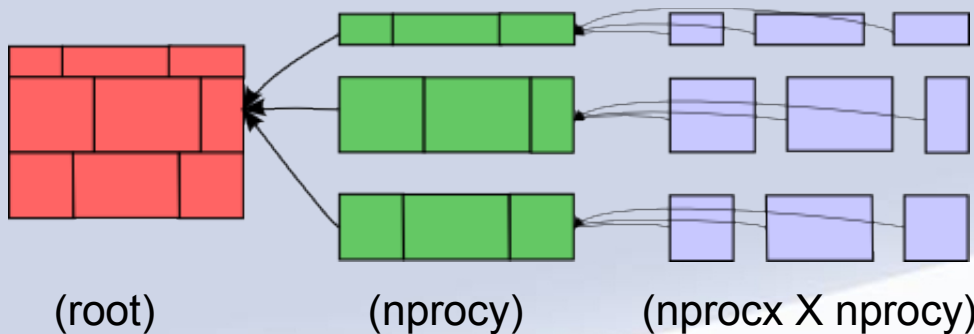# gives a compact context free view

ScalES

# New gather communication

Problem:
- non-uniform but latitude-aligned subarrays
- collective MPI_GATHERV unapplicable



Old:
- many global p2p communications
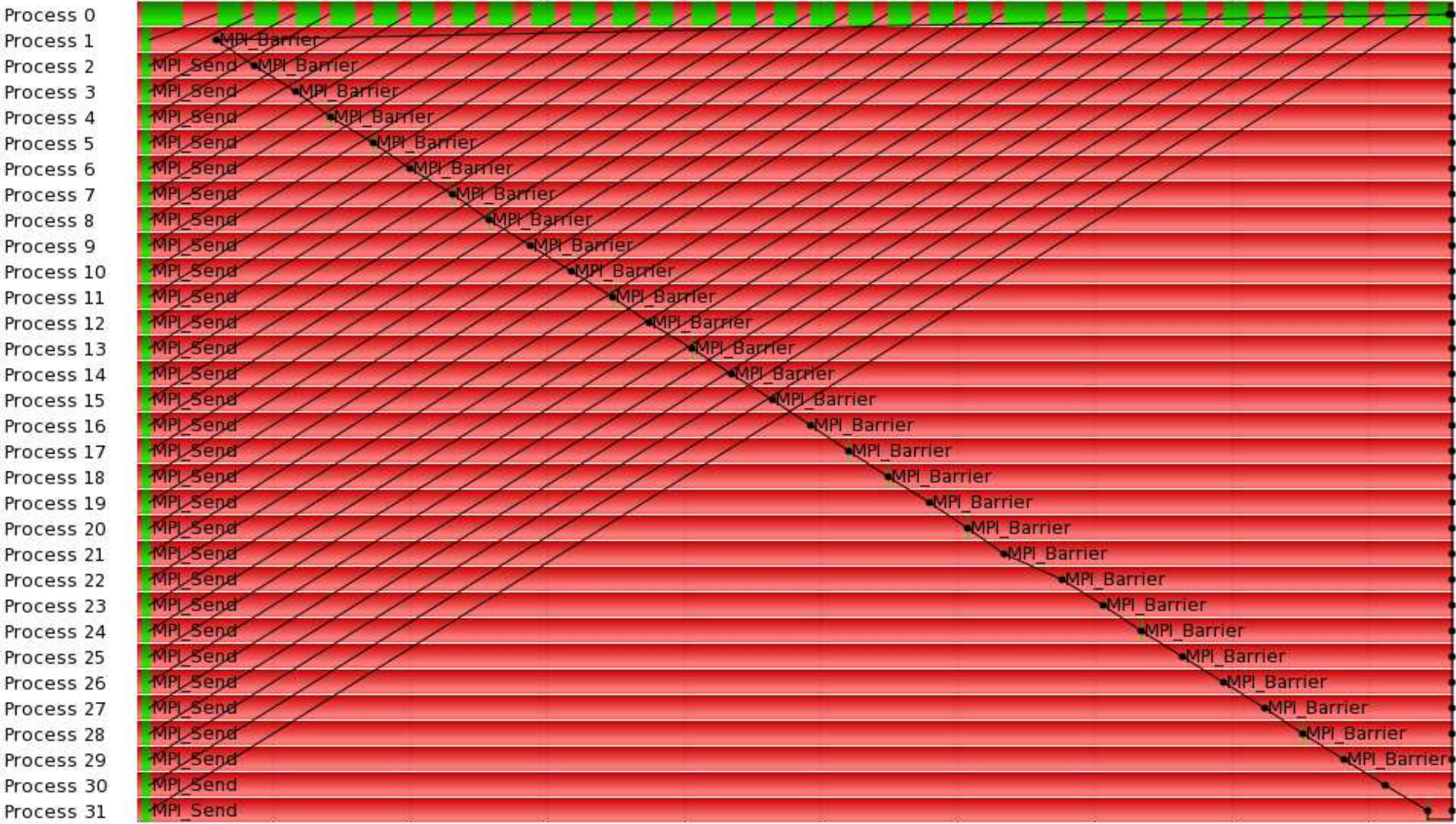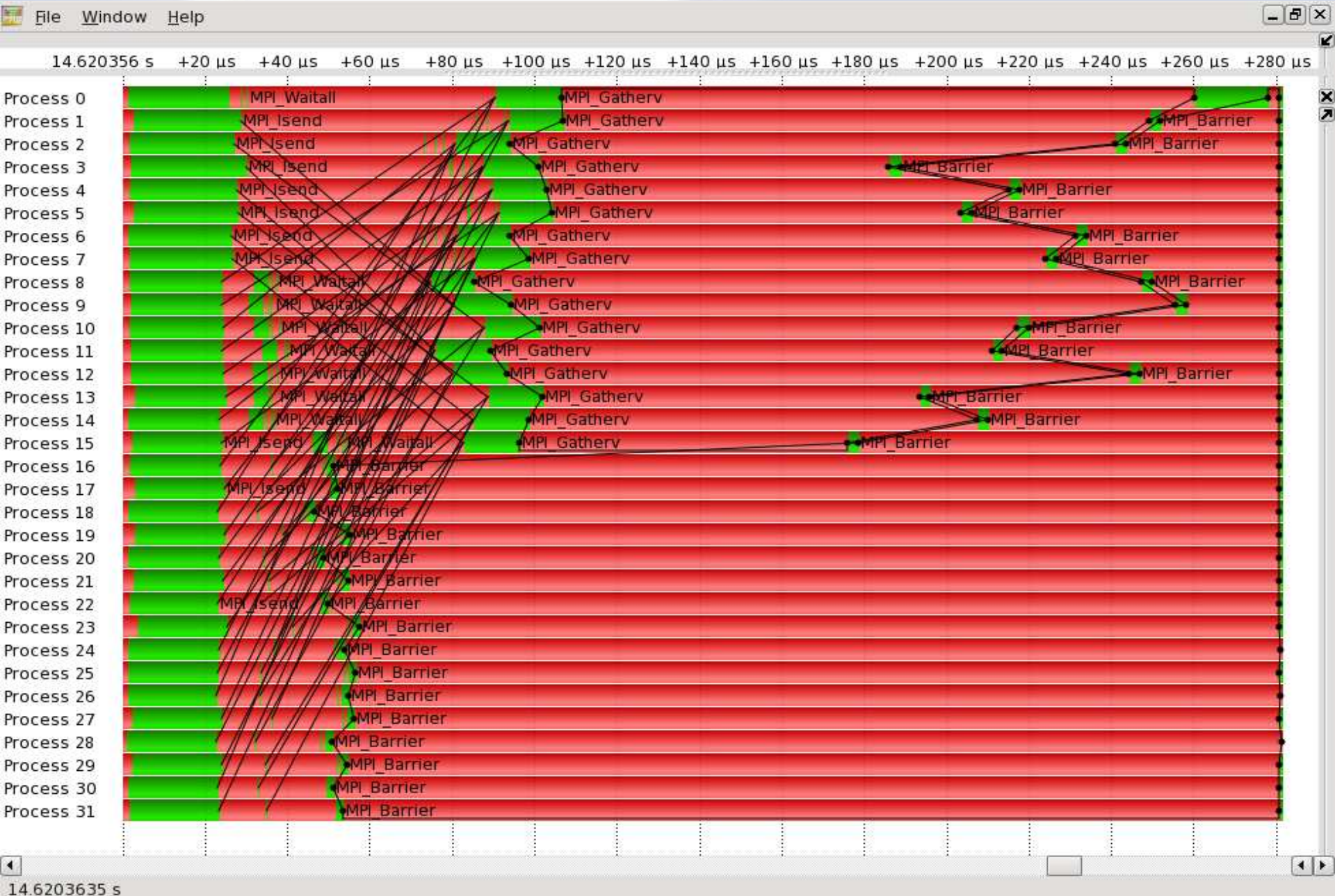- overhead concentrated on root
- increasing cost for high #tasks

New:
- build latitude-aligned subgroups
- distribute shape-overhead
- fast collective MPI_GATHERV for second phase applicable
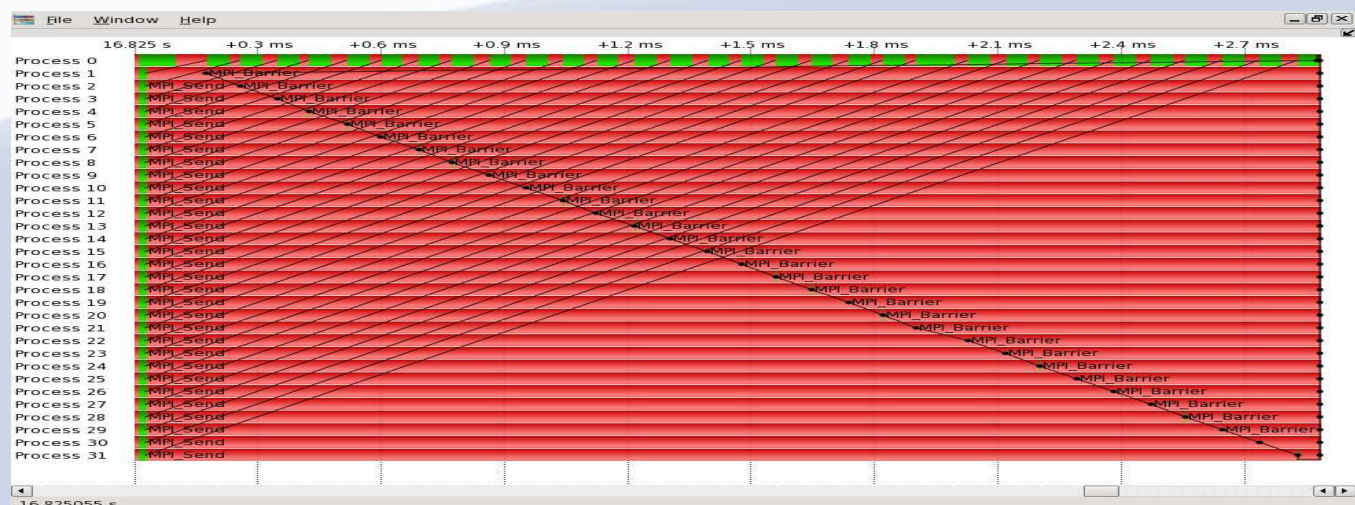- const. cost for high #tasks

(root)          (nprocy)          (nprocx X nprocy)
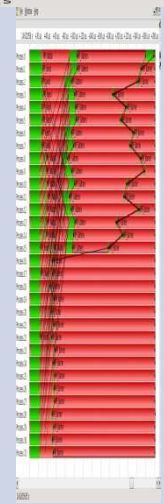
DKRZ

# Old gather: N → 1

# Gather: one-phase vs. two-phase
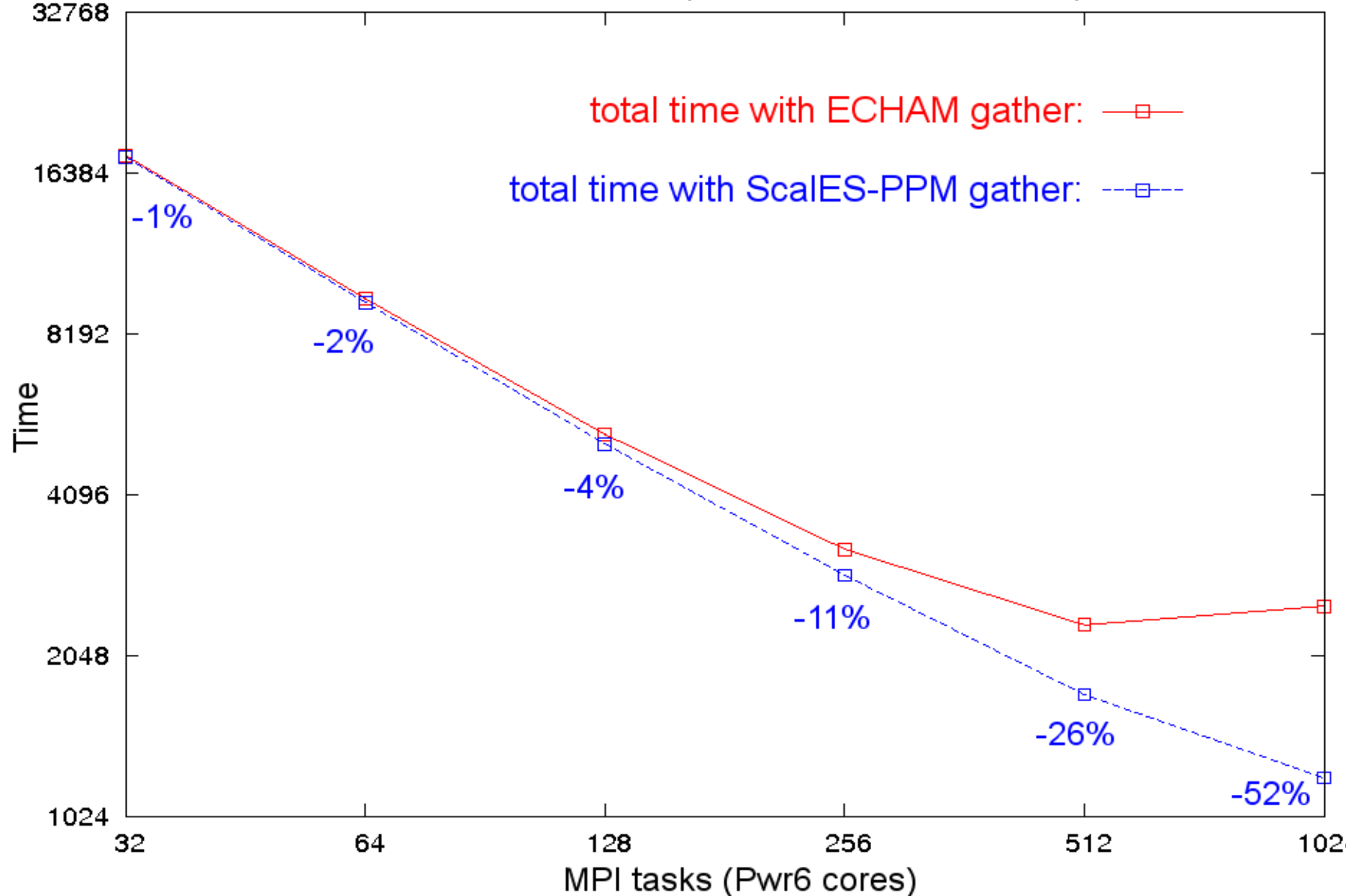
Old:
one-phase

New:
two-phase

0.28 ms

2.8 ms

DKRZ, 2012

# Total runtime measurement



Total runtime in ST mode, ECHAM6 T127L95, 1 month

# Example: MPIOM
# Aggregation of communication

Major communication in MPIOM: boundsexchange

Old implementation:
- Separated updates of x, y, special northborder
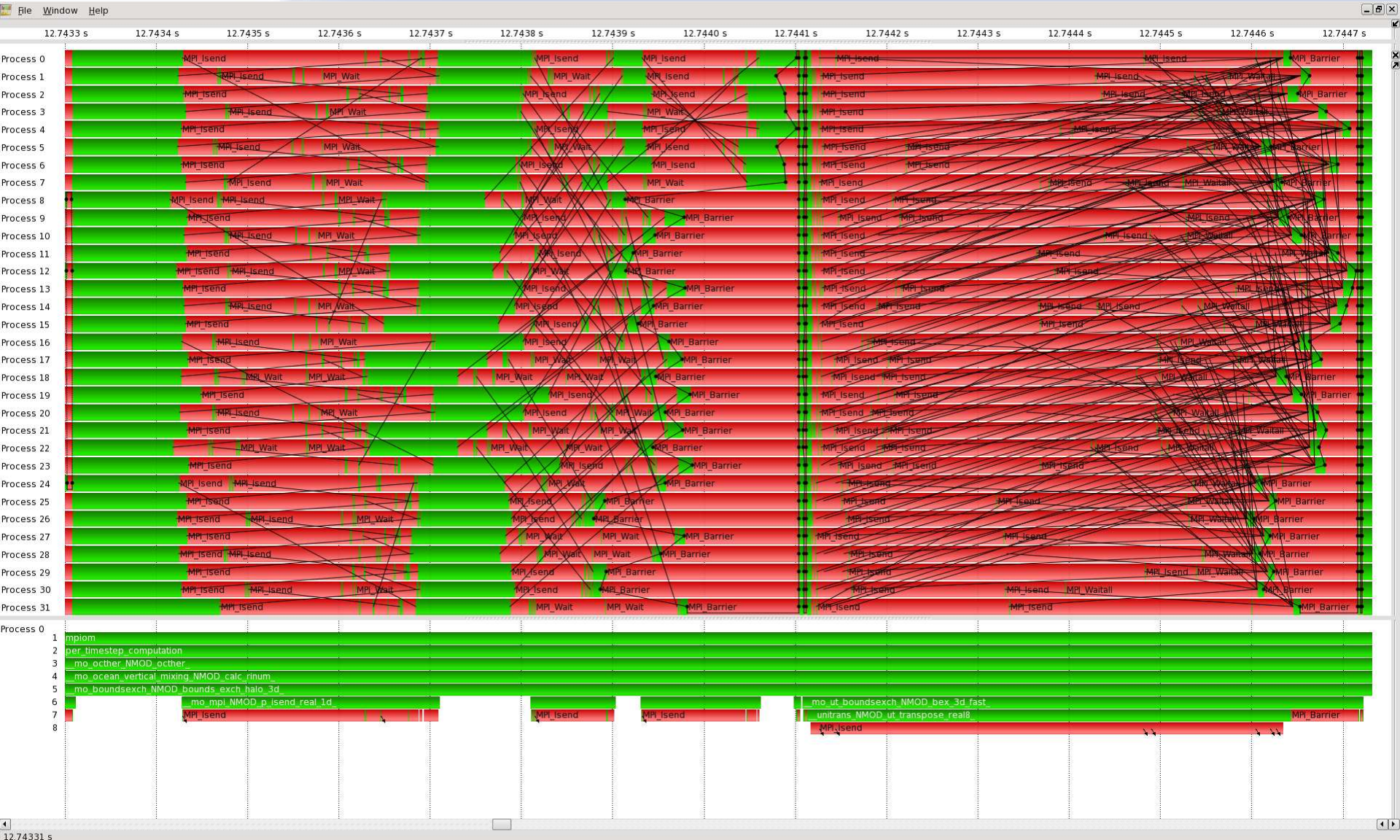- User-buffered messages within each phase

New implementation:
- reprogram communication using a more abstract formulation (using the *Unitrans* communication library)
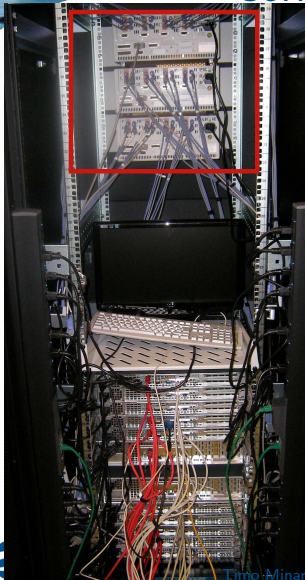- *Unitrans* uses MPI datatypes

Vampir is used to document changes

DKRZ, 2012

ScalES

DKRZ

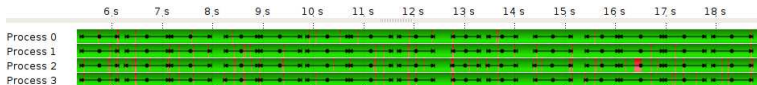# Comparison: old (left), new (right)

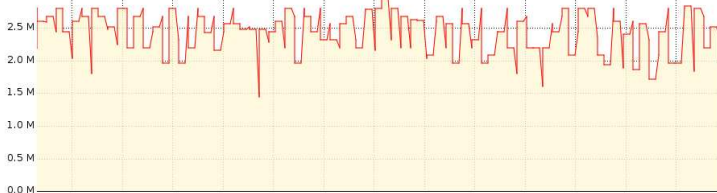# Correlating energy relevant metrics



## Idea

- ▶ Trace energy relevant metrics in database
    - ▶ Processor load, Performance Counter, ...
    - ▶ Processor frequency, ...
    - ▶ Power consumption
- ▶ Merge metrics post-mortem via VT Plugin Interface
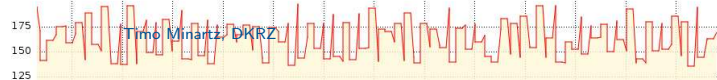- ▶ Switch processor frequency based on application phases

# GETM model power variations

# GETM model instrumented

# Experiences and open questions

## VT Plugin Interface

- ▶ Limitation to 256 counters
- ▶ Post-mortem integration takes a multiple of the application runtime...
- ▶ ... and crashes sometimes

```
[0]VampirTrace: FATAL: OTF_WStream_writeCounter failed:
ERROR in function OTF_WBuffer_setTimeAndProcess,
file: ../../../../extlib/otf/otflib/OTF_WBuffer.c, line: 296:
time not increasing. (t= 14835708593910, p= 1)
```