



CUBE: A success story -- But does it have a future?

Vi-HPS 10th Anniversary Workshop | Bernd Mohr

Context: 18 Years of Automatic Trace Analysis

- 1999 – 2004



- Working Group
 - EU ESPRIT + IST
- <http://webarchiv.fz-juelich.de/apart/>



- Sequential analyzer
 - EXPERT
- <http://webarchiv.fz-juelich.de/jsc/kojak.html>

- 2006 – now

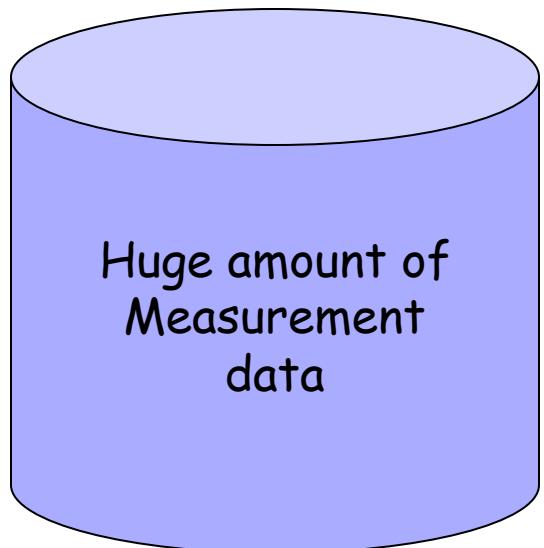


- Virtual Institute
 - Helmholtz
- <http://www.vi-hps.org/>

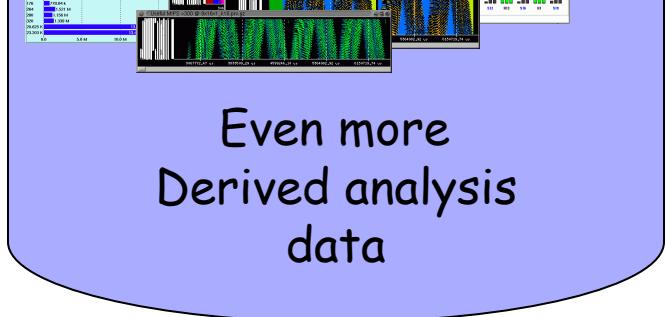
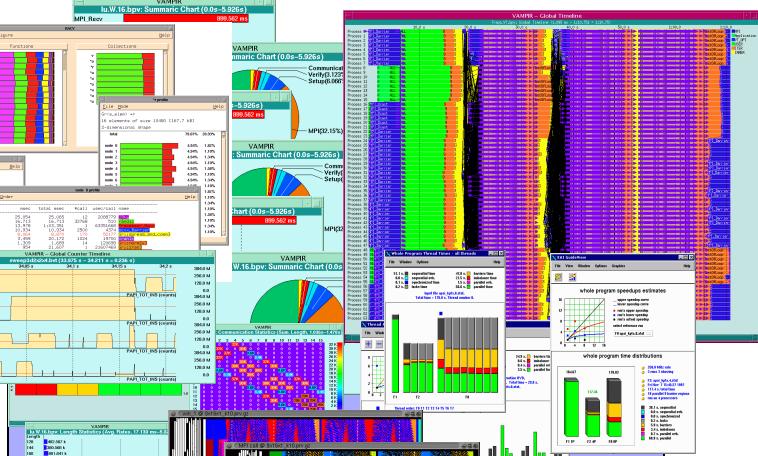


- Parallel analyzer
 - SCOUT
- <http://www.scalasca.org>

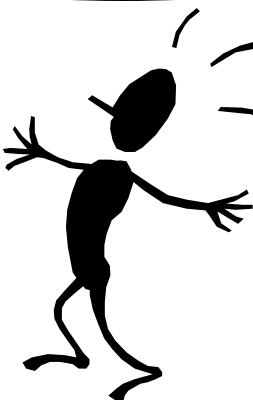
Traditional Performance Tools



Little Simple analysis

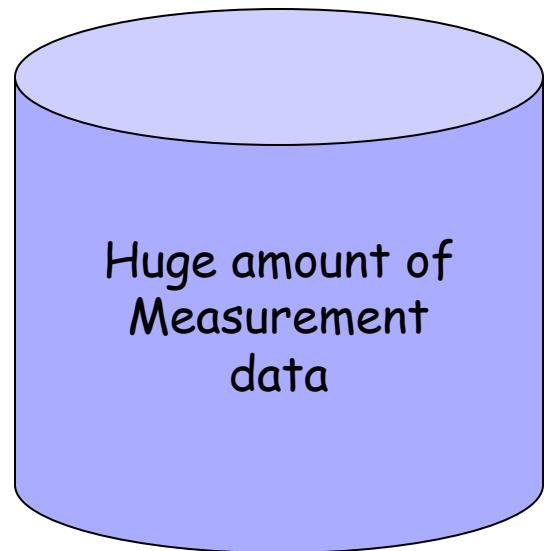


Poor User

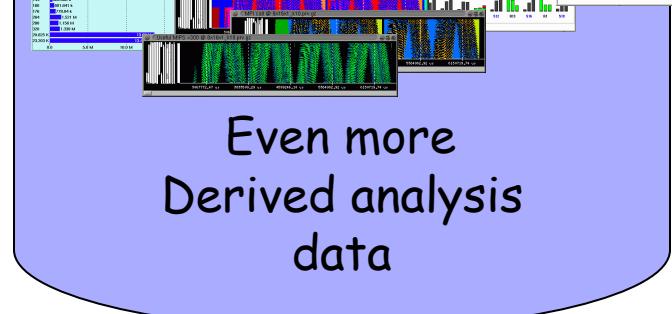
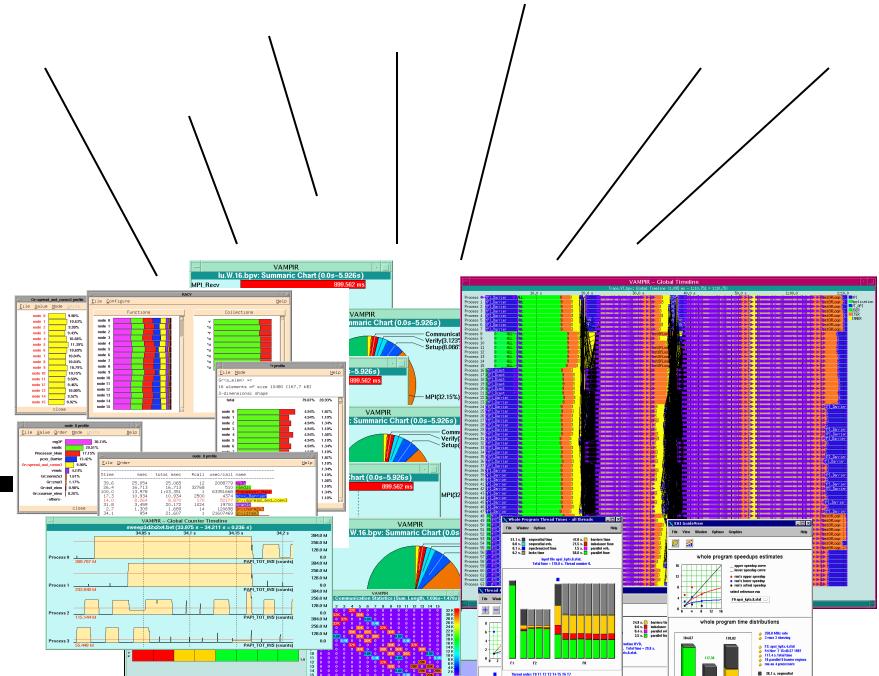


[2005]

Traditional Performance Tools

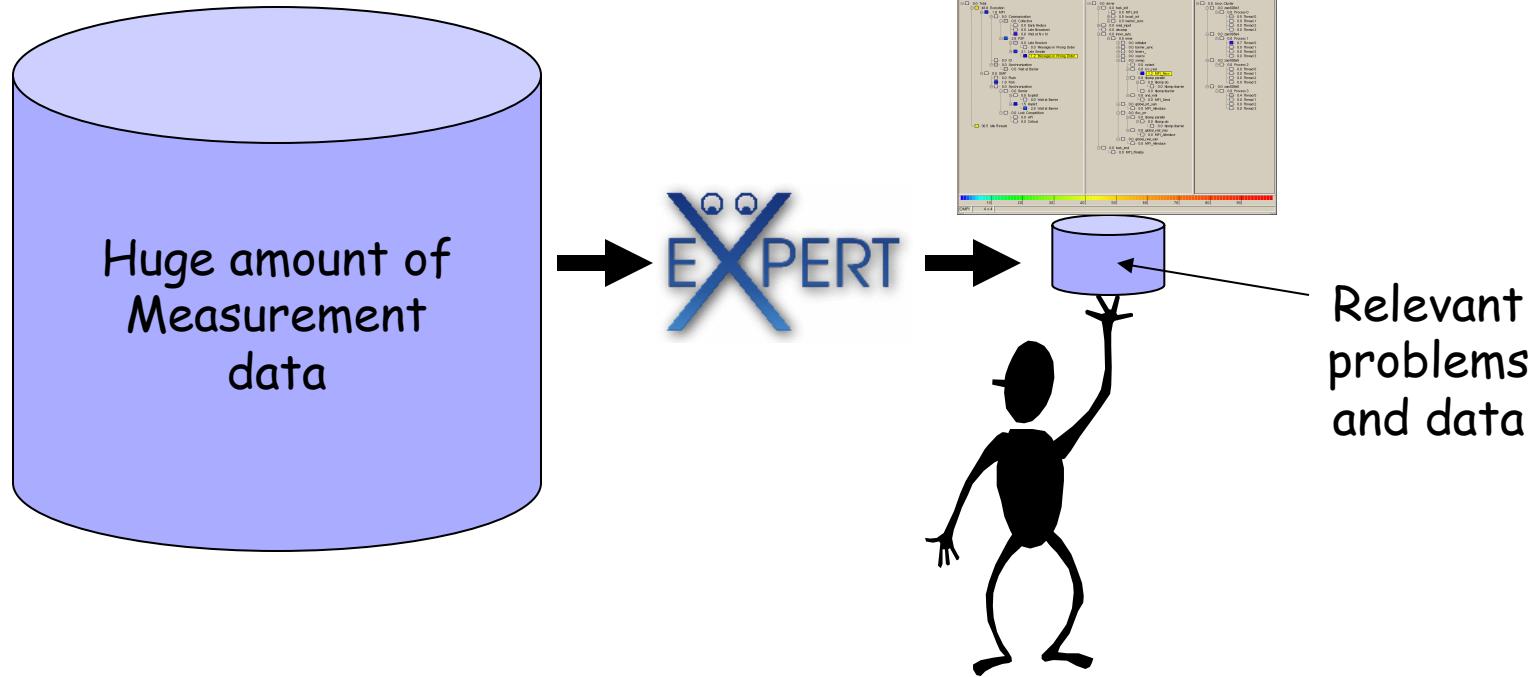


Little Simple analysis



[2005]

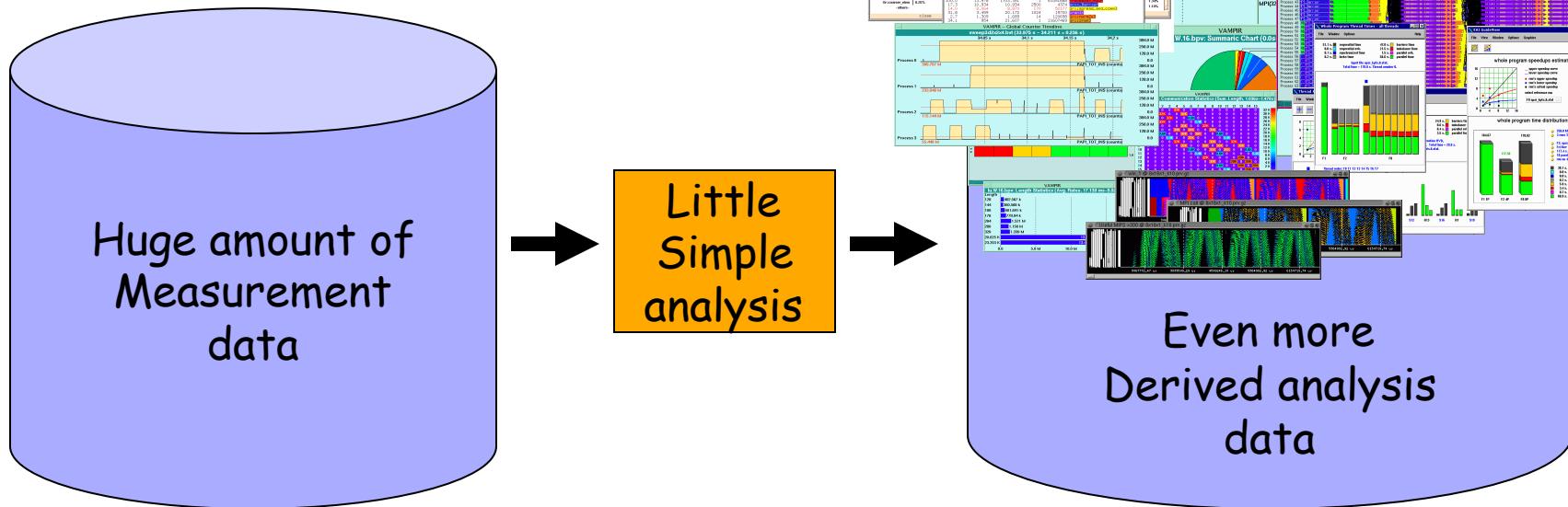
Solution Part 1: Automatic Tool



- For standard cases (90% ?!)
- For “normal” users
- Starting point for experts

[2005]

Solution Part 2: Expert Tools + Expert



- For non-standard / tricky cases (10%)
 - ⇒ More productivity for performance analysis process!

[2005]

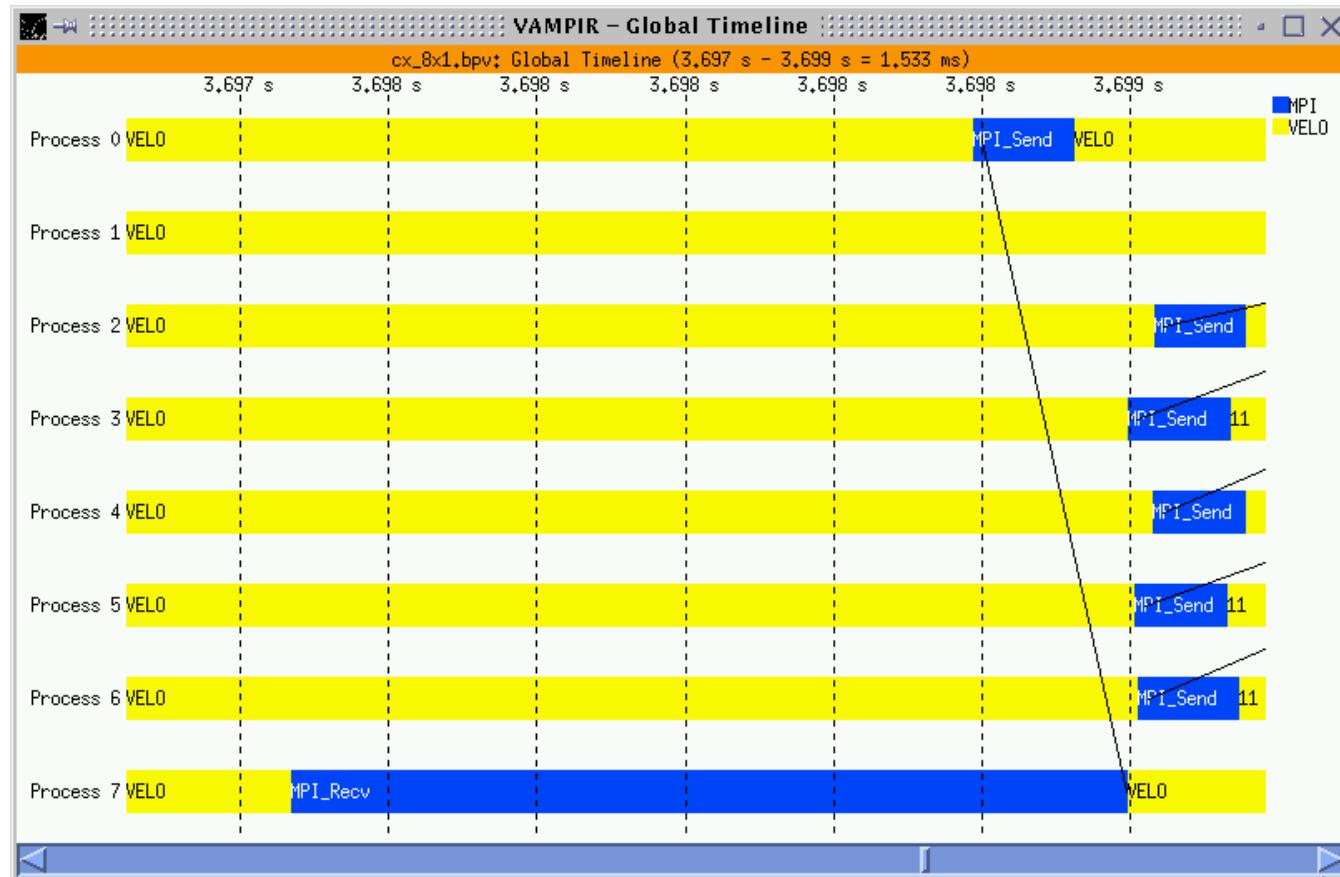
6



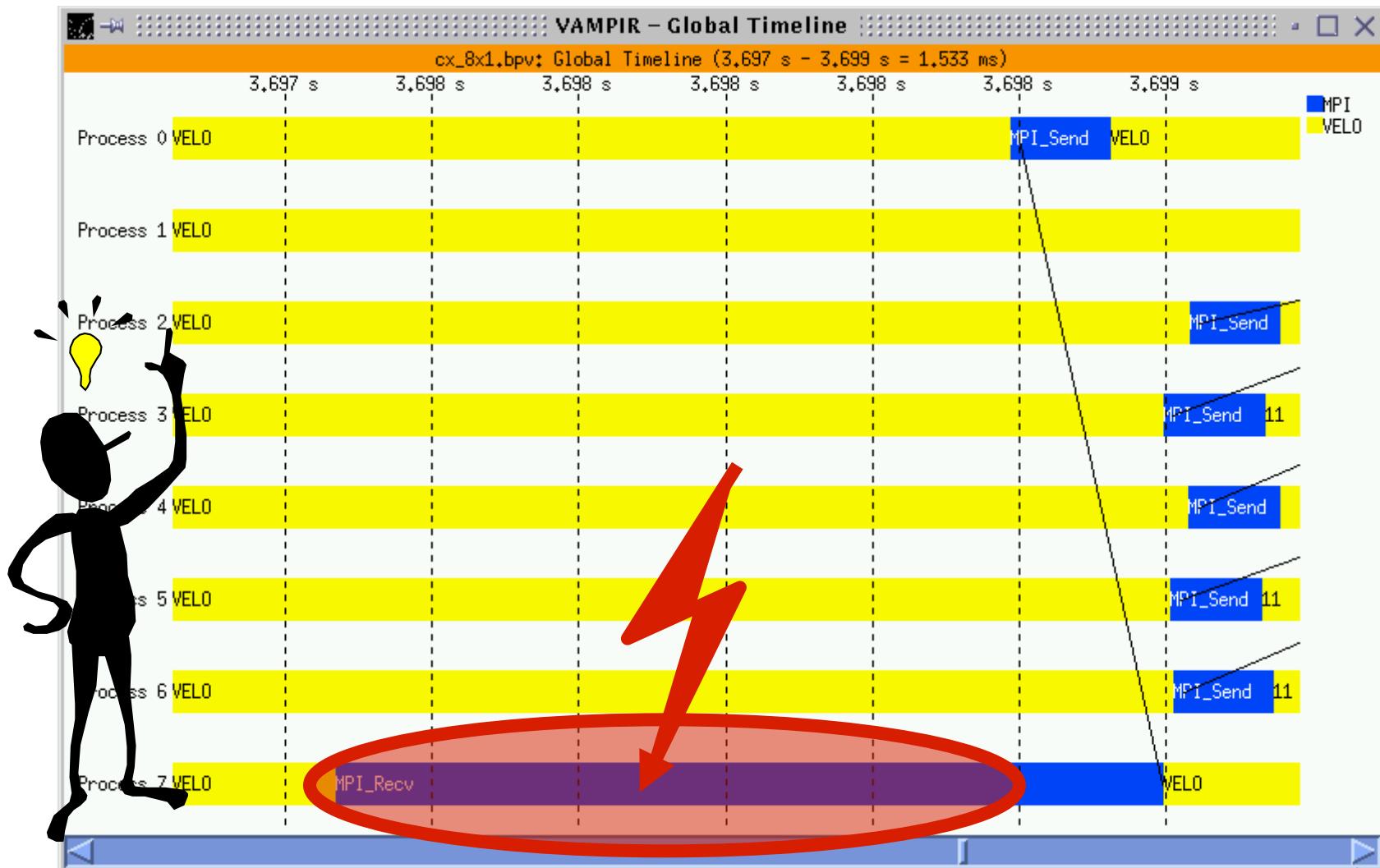
Example: Late Sender (blocked receiver)

[2001]

Czochralski crystal growth

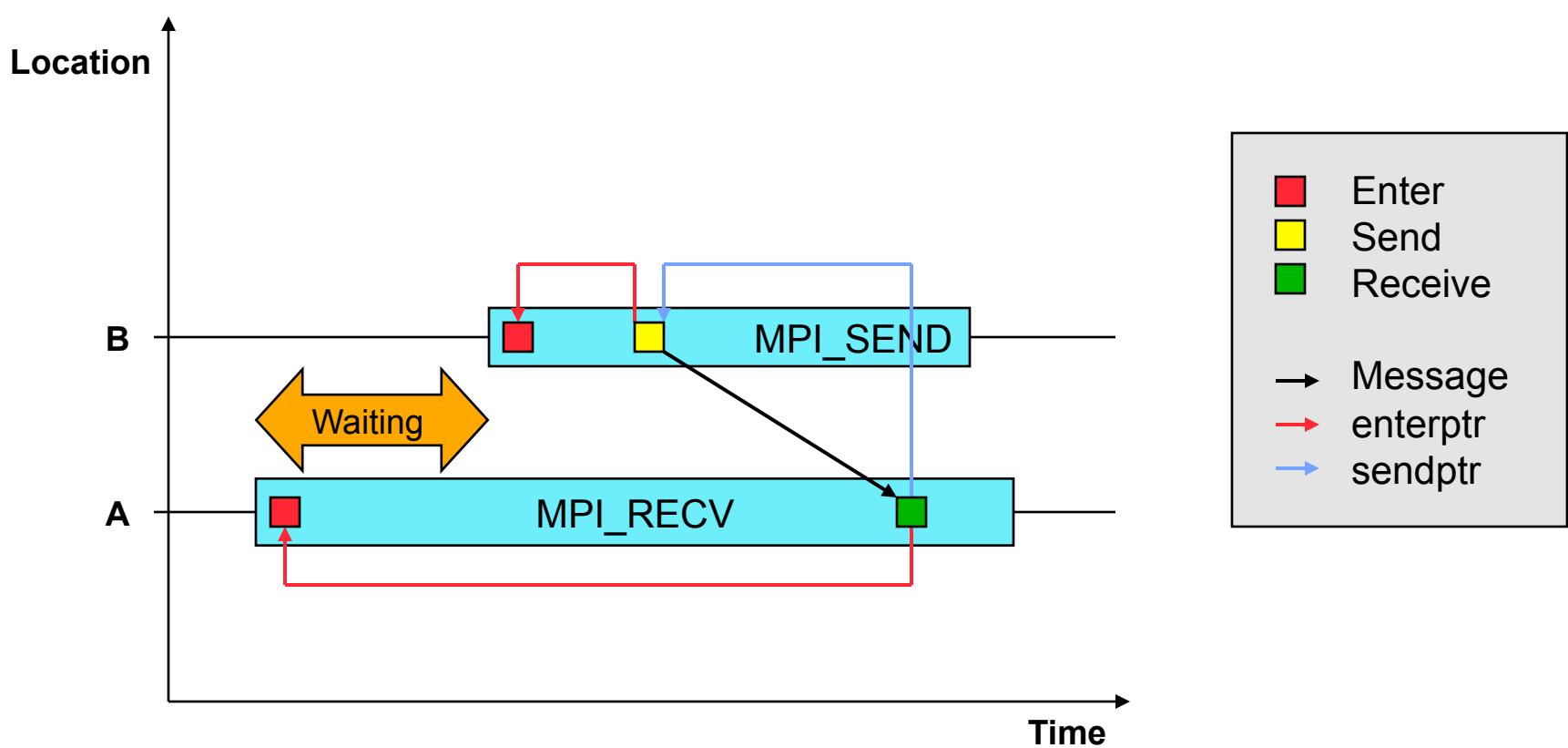


Example Automatic Analysis: Late Sender

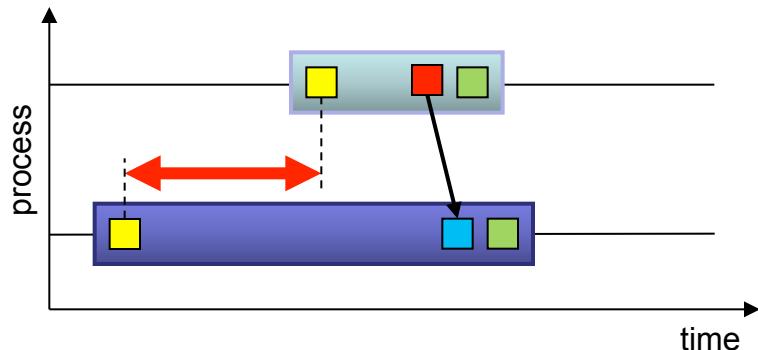


Example: Late Sender (blocked receiver)

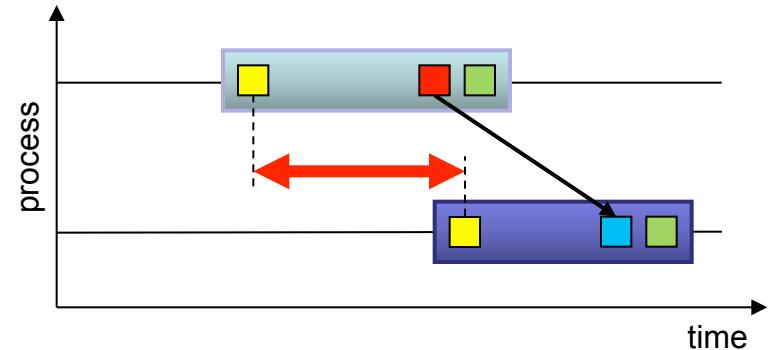
[2001]



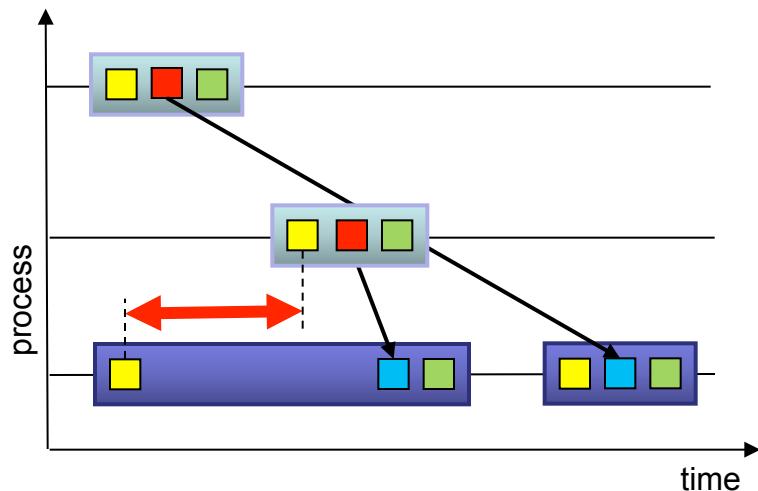
Scalasca: Example MPI Patterns



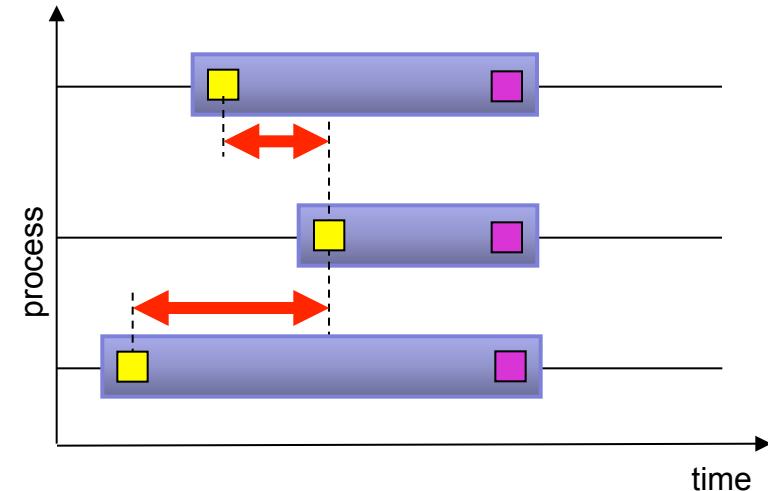
(a) Late Sender



(b) Late Receiver



(c) Late Sender / Wrong Order



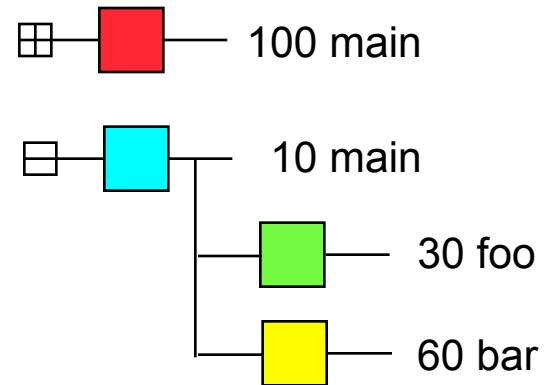
(d) Wait at $N \times N$

■ ENTER ■ EXIT ■ SEND ■ RECV ■ COLLEXIT

Presentation of Performance Behavior

[2001]

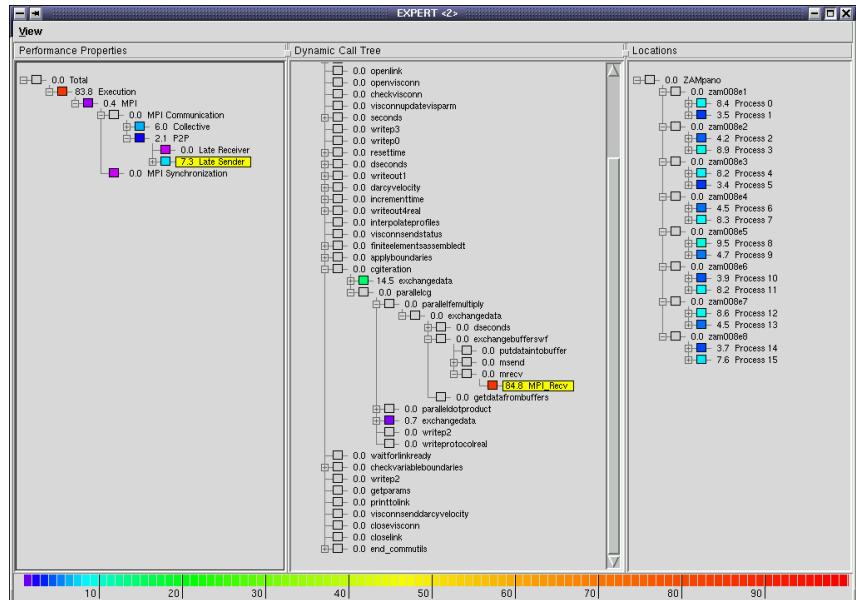
- **Performance behavior**
 - 3 dimensional matrix
 - Hierarchical dimensions
- **Weighted tree**
 - Tree browser
 - Each node has weight
 - * Percentage of CPU allocation time
 - * E.g. time spent in subtree of call tree
 - Displayed weight depends on state of node
 - * Collapsed (including weight of descendants)
 - * Expanded (without weight of descendants)
 - Displayed using
 - * Color
 - Allows to easily identify hot spots (bottlenecks)
 - * Numerical value
 - Detailed comparison

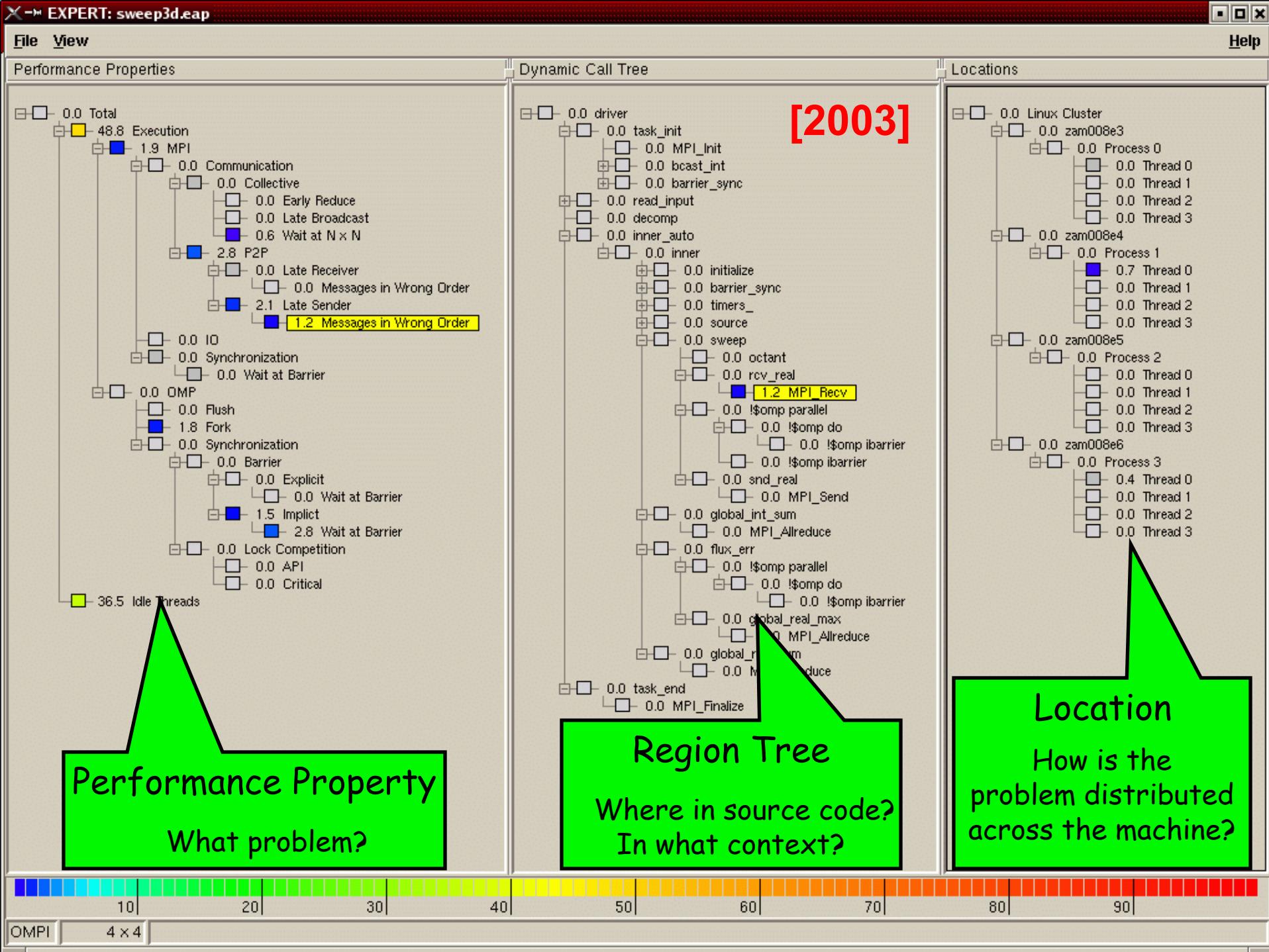


Presentation of Performance Behavior (2)

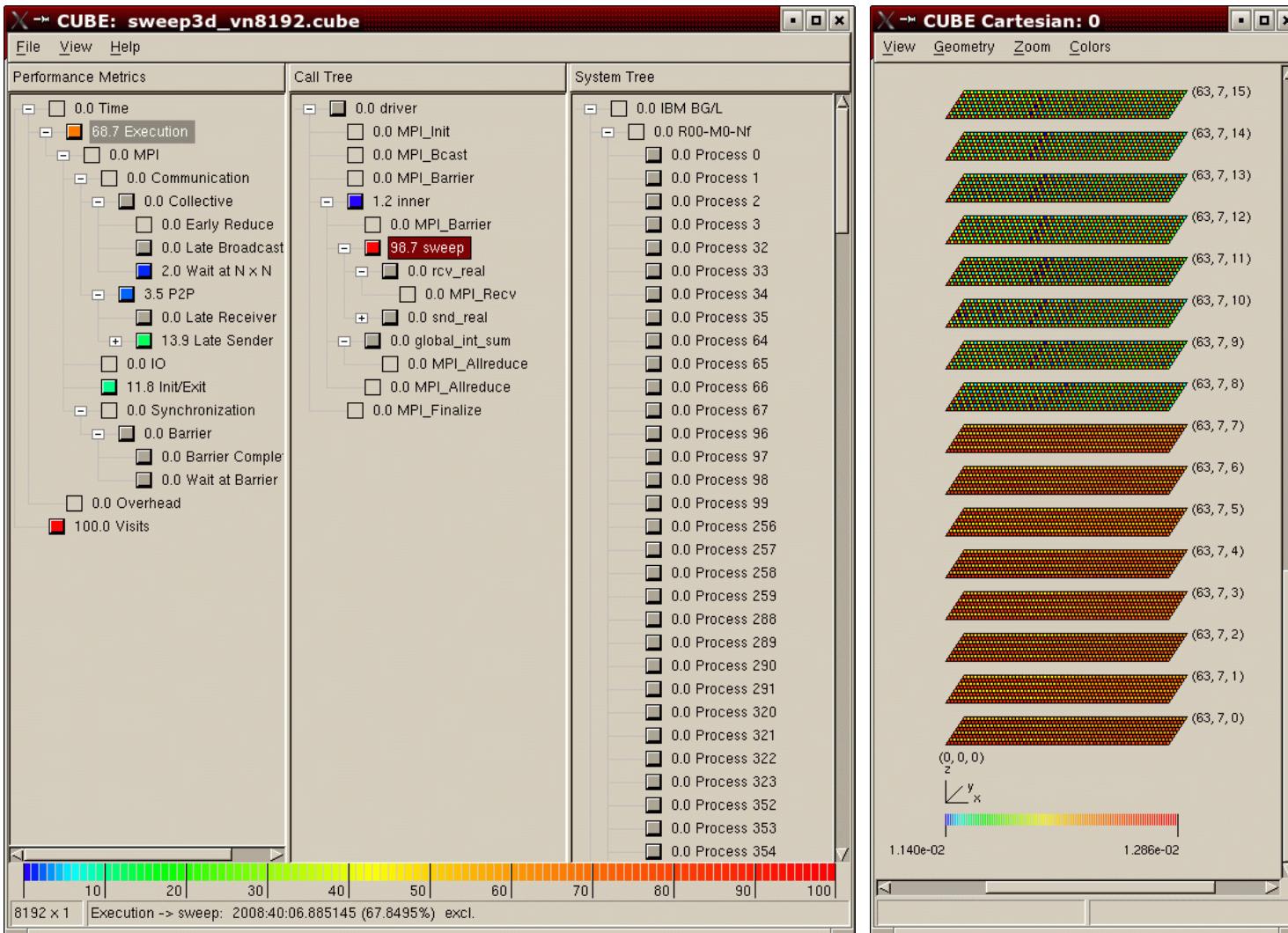
[2001]

- Three views
 - Performance property
 - Call tree
 - Locations
- Interconnected
 - View refers to selection in left neighbor
- Two modes
 - Absolute: percent of total CPU allocation time
 - Relative: percent of selection in left neighbor
- Collapsing/expanding of nodes
 - Analysis on all hierarchy levels



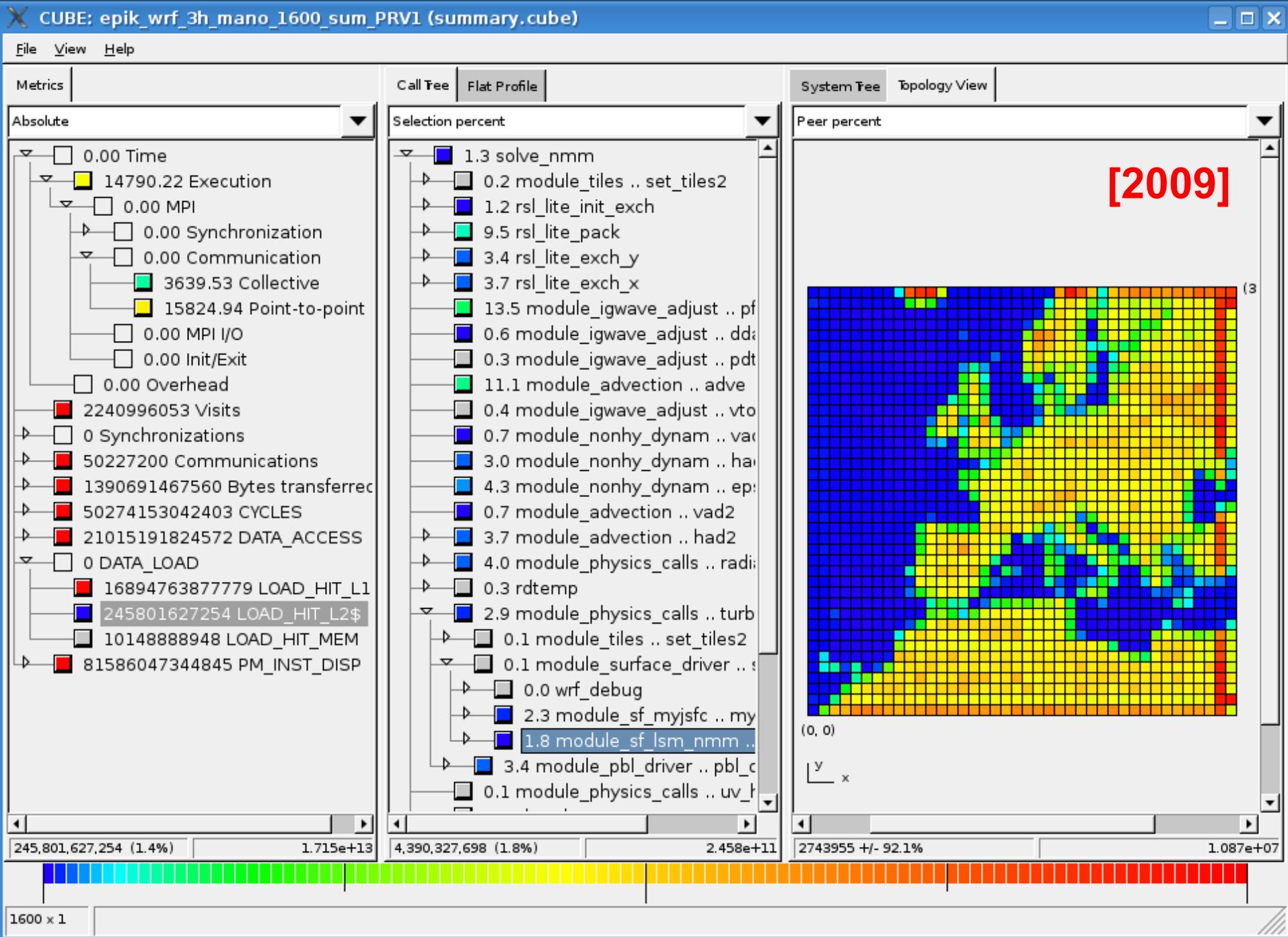


Example: sweep3D on 8192 BG/L PEs



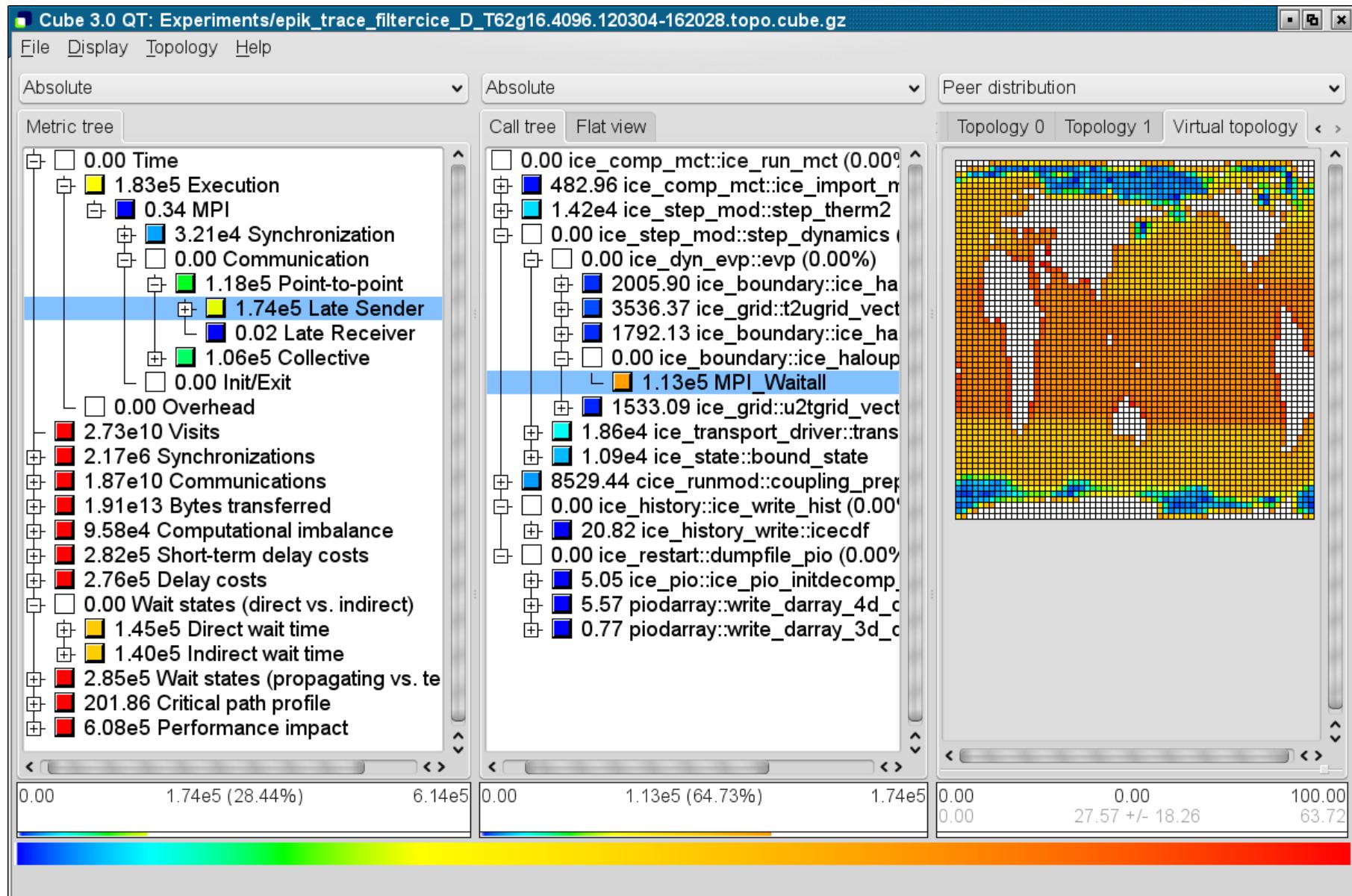
- New topology display
- Shows distribution of pattern over HW topology
- Scales to larger systems

[2007]



WRF-NMM weather prediction code on MareNostrum @ 1600 CPUS

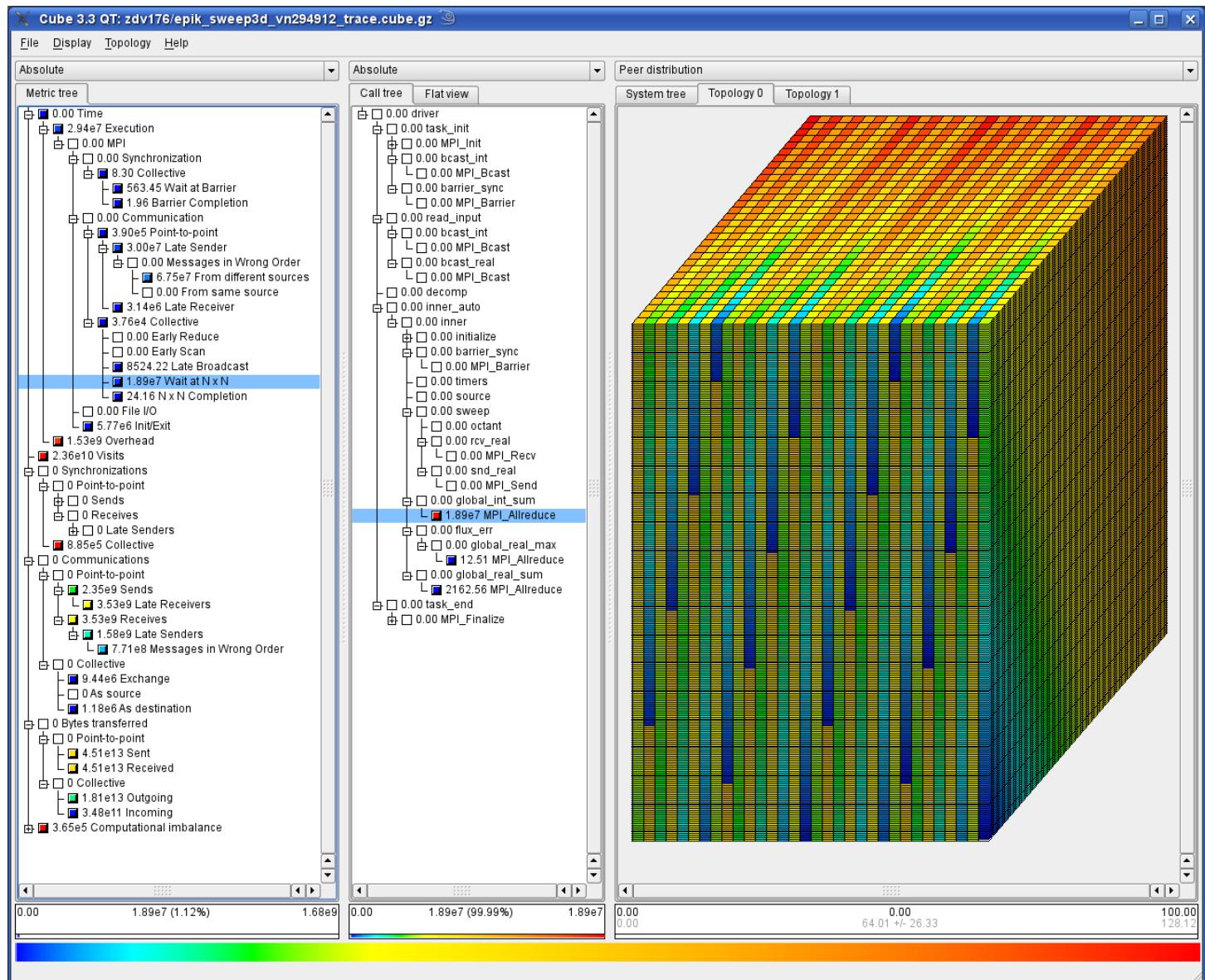
Scalasca Example: CESM Sea Ice Module Late Sender Analysis + Application Topology



Trace analysis sweep3D@294,912

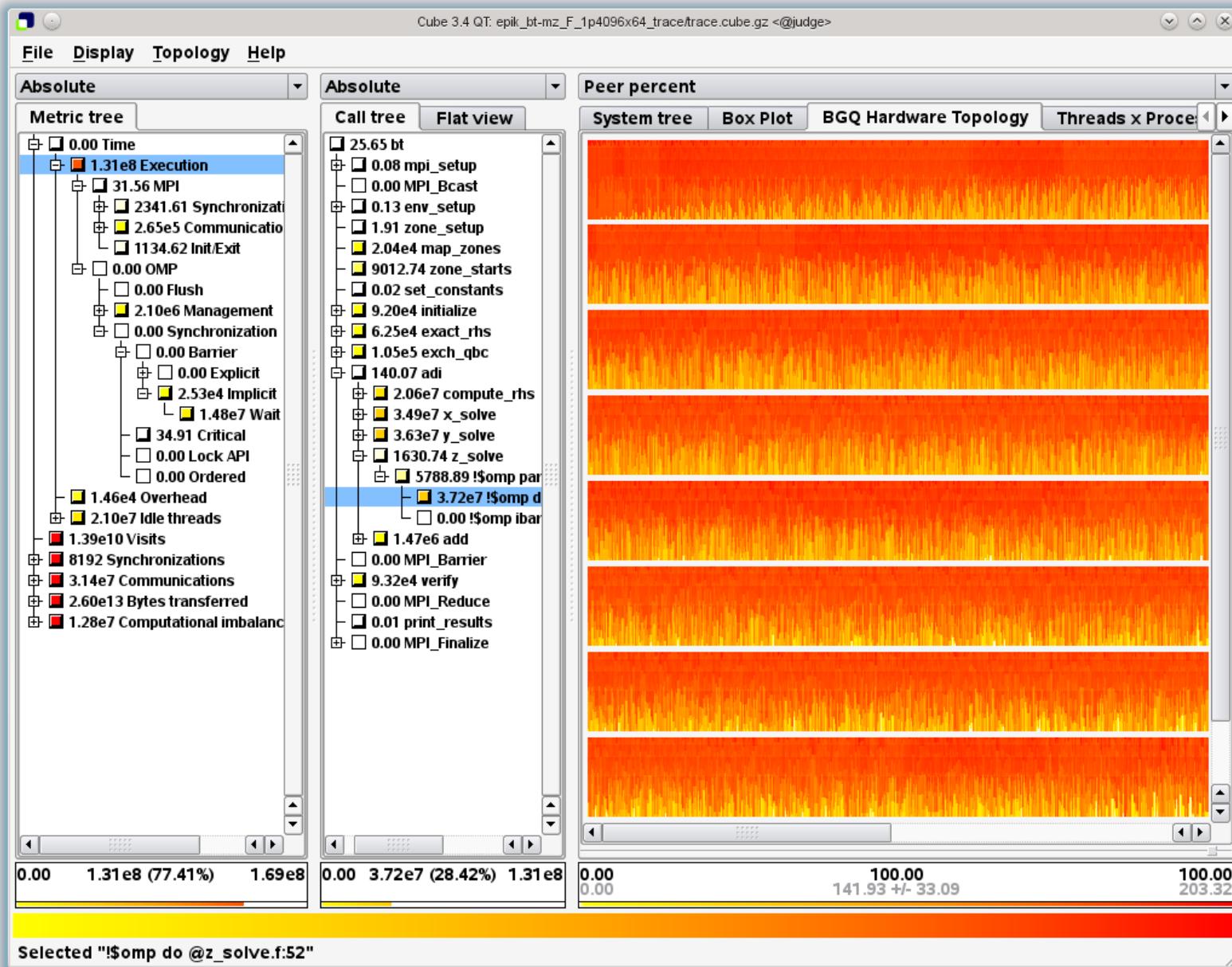
- 10 min sweep3D runtime
- 11 sec replay
- 4 min trace data write/read (576 files)
- 7.6 TB buffered trace data
- 510 billion events

[2010]

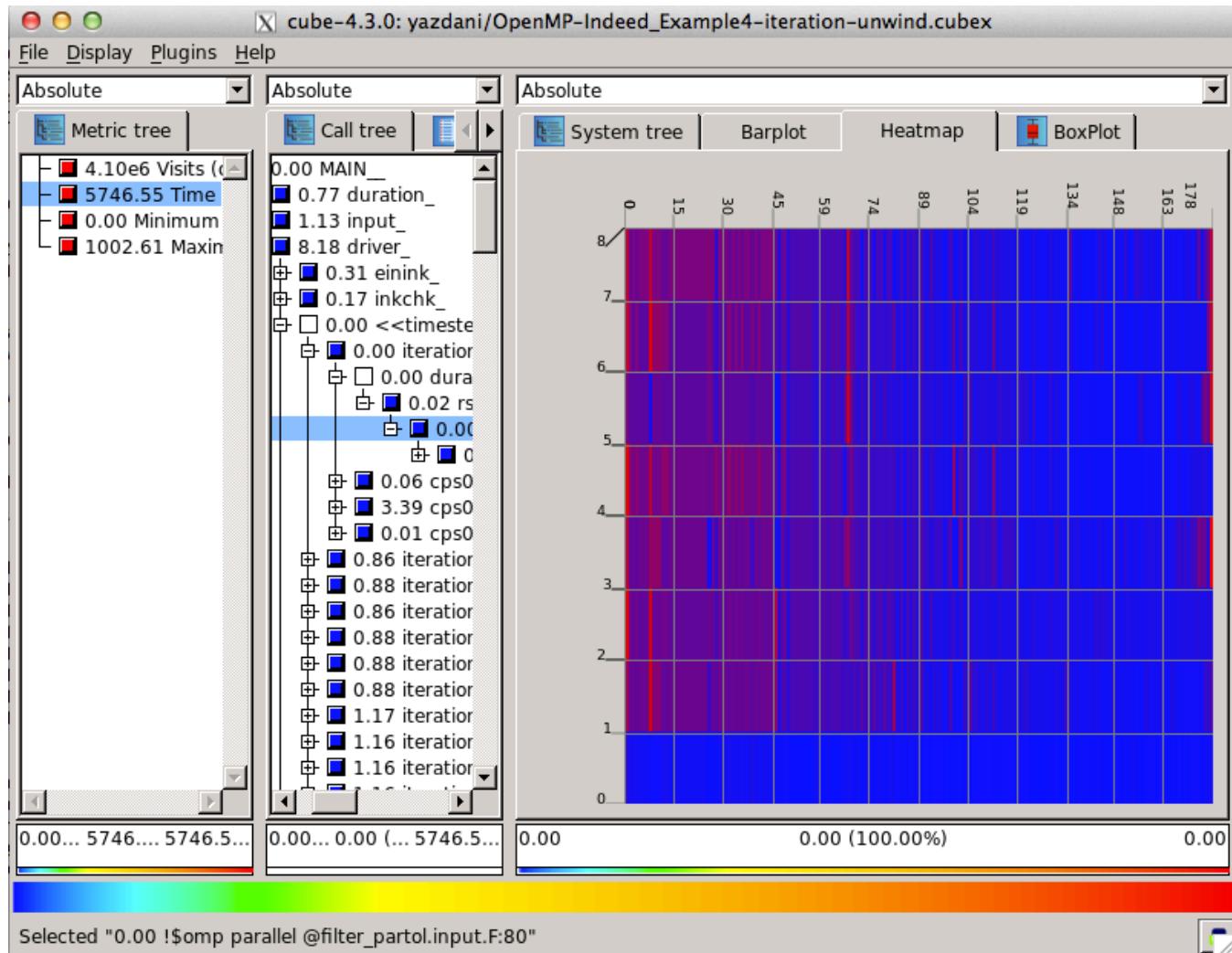


Scalasca trace analysis bt-mz@1,048,704 BGQ

[2013]



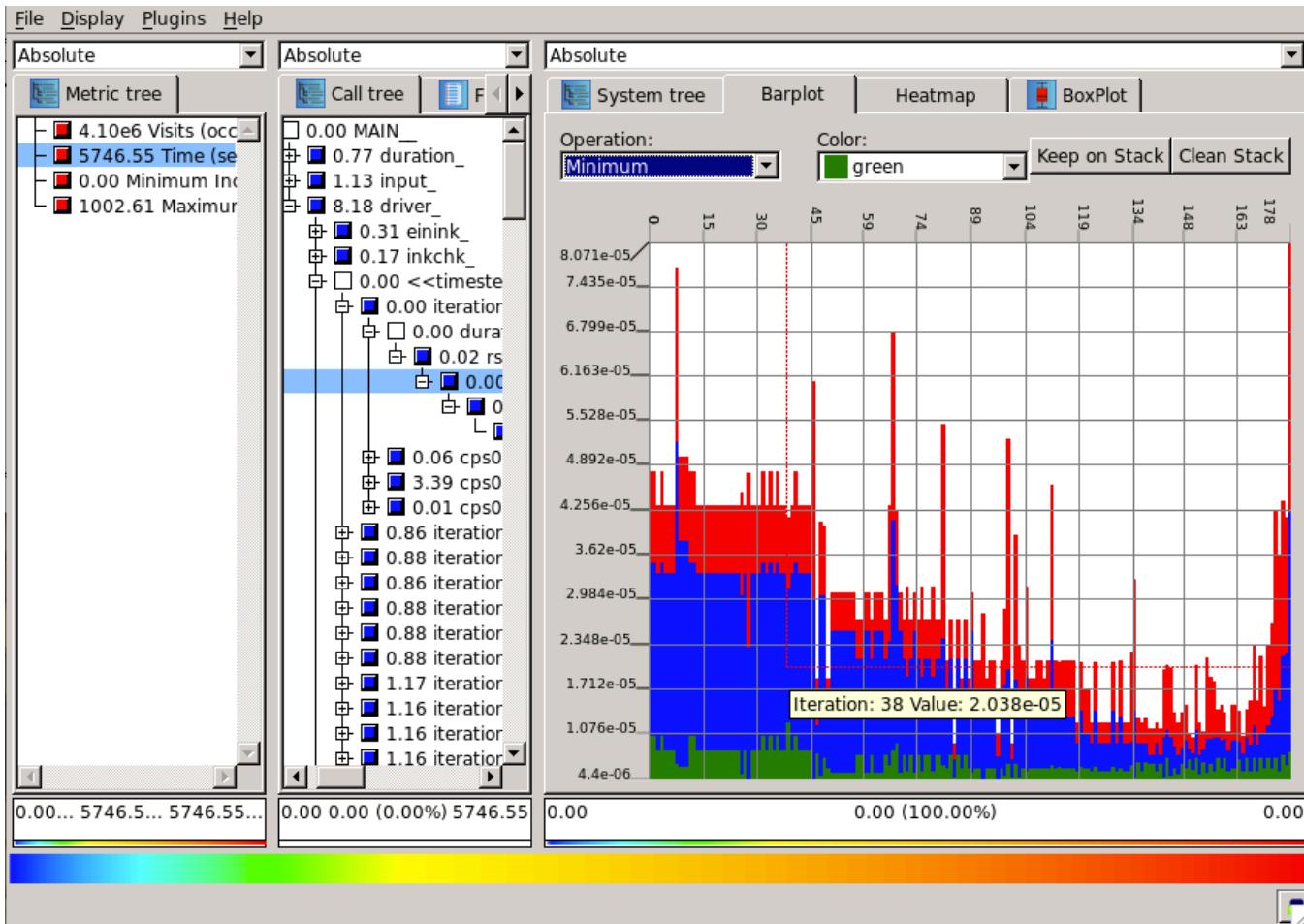
- **Phase profiling**
- Collects data for each instance of phases marked in program instead of aggregating it
- Shows data over “time” (phase instances) for each rank/thread



Cube Viz Plugins: Phase Barplot

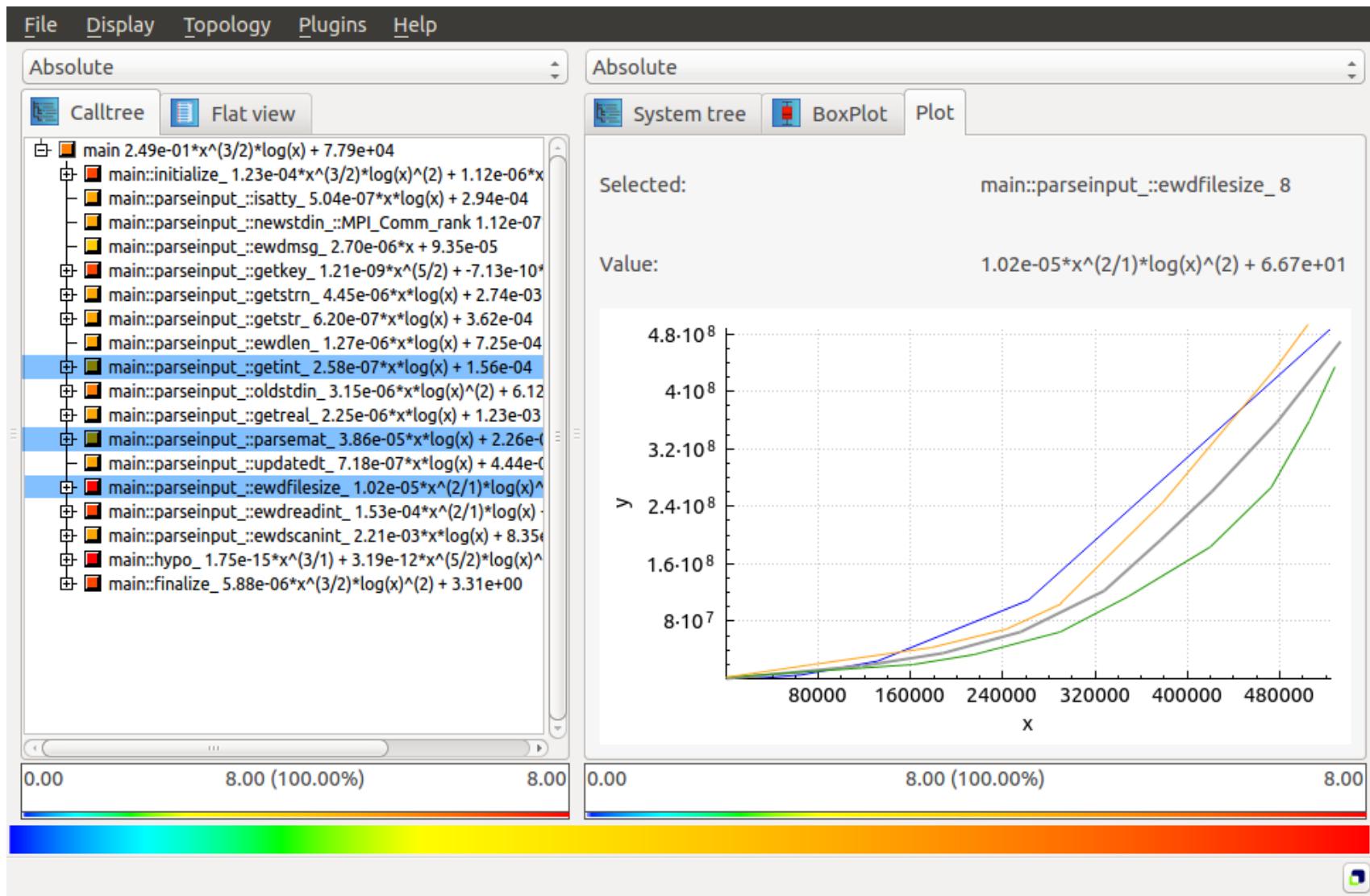
[2015]

- **Phase profiling**
- Collects data for each instance of phases marked in program instead of aggregating it
- Shows min/max/avg metric value over “time” (phase instances)



Catwalk: Result Visualization

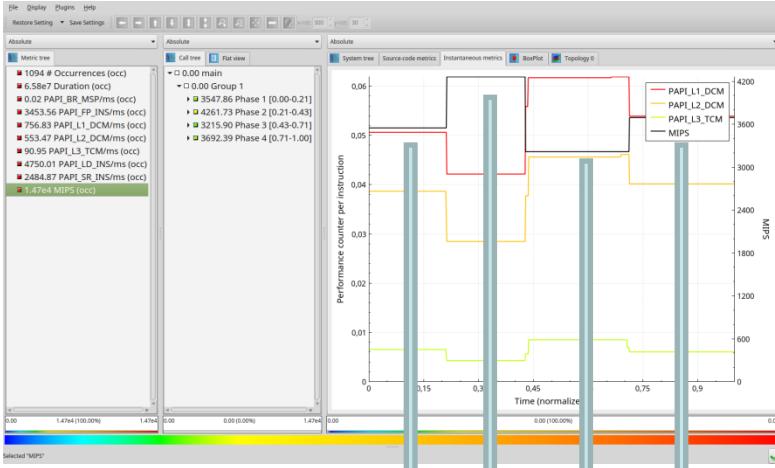
[2015]



T5.4 – Example: Stream Benchmark

[2015]

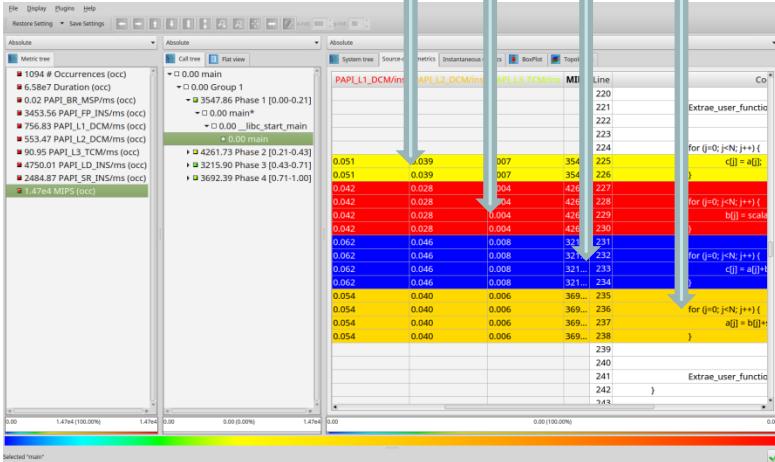
- Detailed performance progression (4 phases)



Stream code

```
for (i = 0; i < NITERS; i++)  
{  
    Extrae_event (BEGIN);  
    for (j = 0; j < N; j++)  
        c[j] = a[j];  
    for (j = 0; j < N; j++)  
        b[j] = s*c[j];  
    for (j = 0; j < N; j++)  
        c[j] = a[j]+b[j];  
    for (j = 0; j < N; j++)  
        a[j] = b[j]+s*c[j];  
    Extrae_event (END);  
}
```

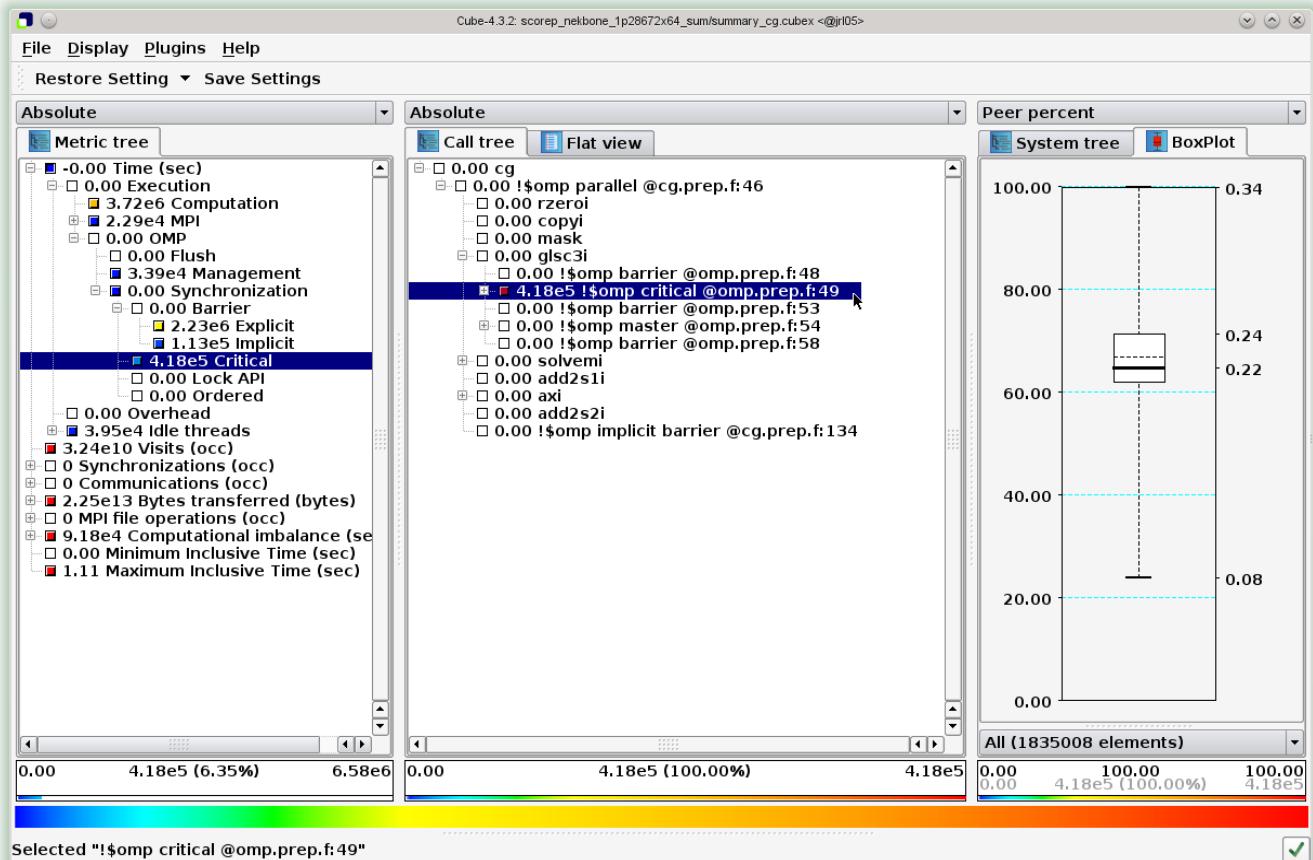
- Source code reference (4 kernels within the instrumented region)



Scalasca: 1,835,008 Threads Test Case

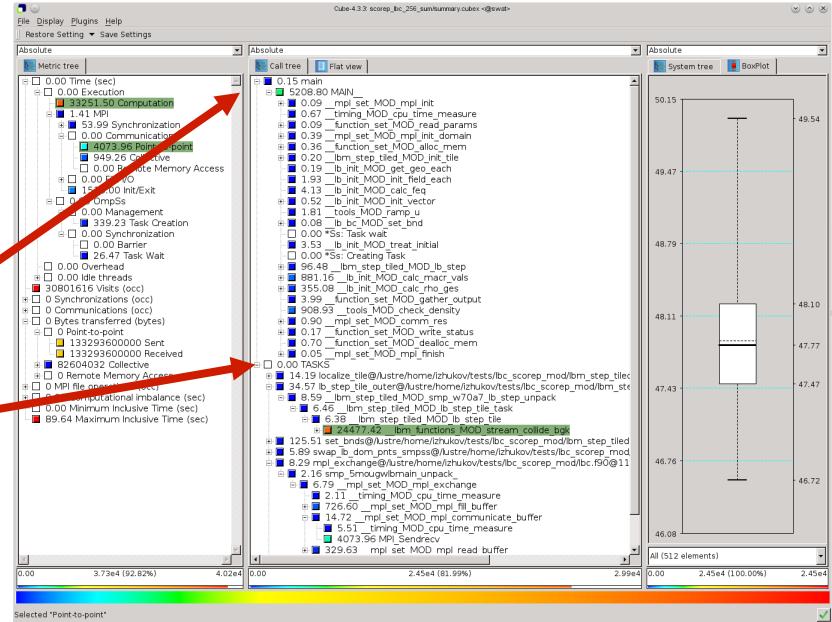
[2016]

- Nekbone
- CORAL benchmark
- JuQueen experiment
- $28,672 \times 64 = 1,835,008$ threads
- Load imbalance at OpenMP critical section



Issue in the Future?

- Trend to asynchronous programming models
 - Kernels (CUDA, OpenCL)
 - Tasks (OpenMP, OmpSs)
 - Threads (Pthreads, ...)
- CUBE representation of these items
 - Separate list of trees besides main
 - Non-scalable
 - Does not show dependencies
 - To creation / launch / spawn location
 - Between items (e.g. tasks)
- **Need new concept?**



**YOU KNOW YOU MADE IT ...
... IF LARGE COMPANIES
“COPY” YOUR STUFF**

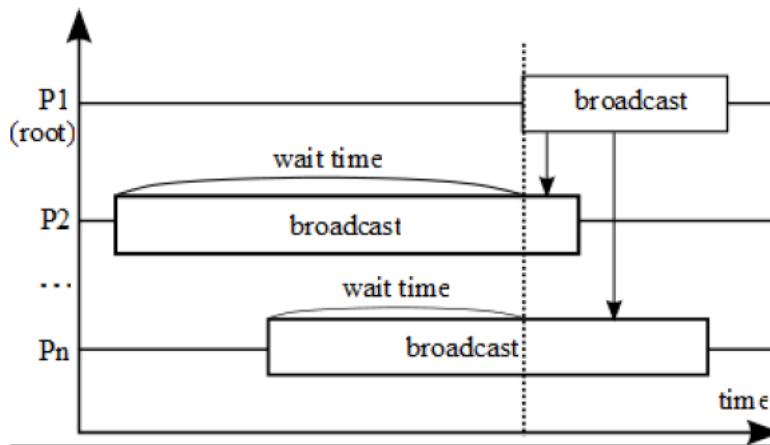
Introducing the Intel® Trace Analyzer and Collector Performance Assistant

Motivation: Improve method of performance analysis via the GUI

Solution:

- Define common/known performance problems
- Automate detection via the Intel® Trace Analyzer

Example: A “Late Broadcast” is not easy to identify with existing views



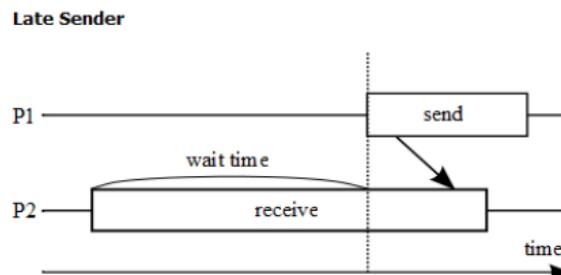
[2014]

Source: <https://software.intel.com/en-us/videos/quickly-discover-performance-issues-with-the-intel-trace-analyzer-and-collector-90-beta>

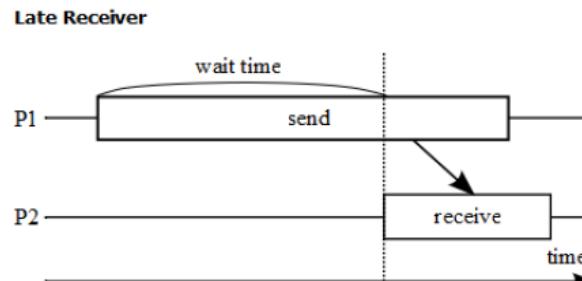
Which Performance Issues are automatically identified?

Point-to-point exchange problems:

- Late Sender

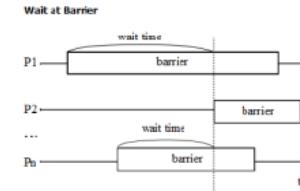


- Late Receiver

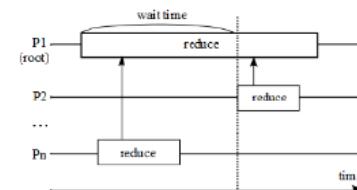


Problems with global collective operation performance:

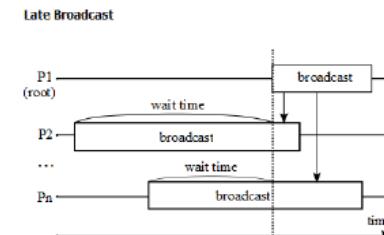
- Wait at Barrier



- Early Reduce

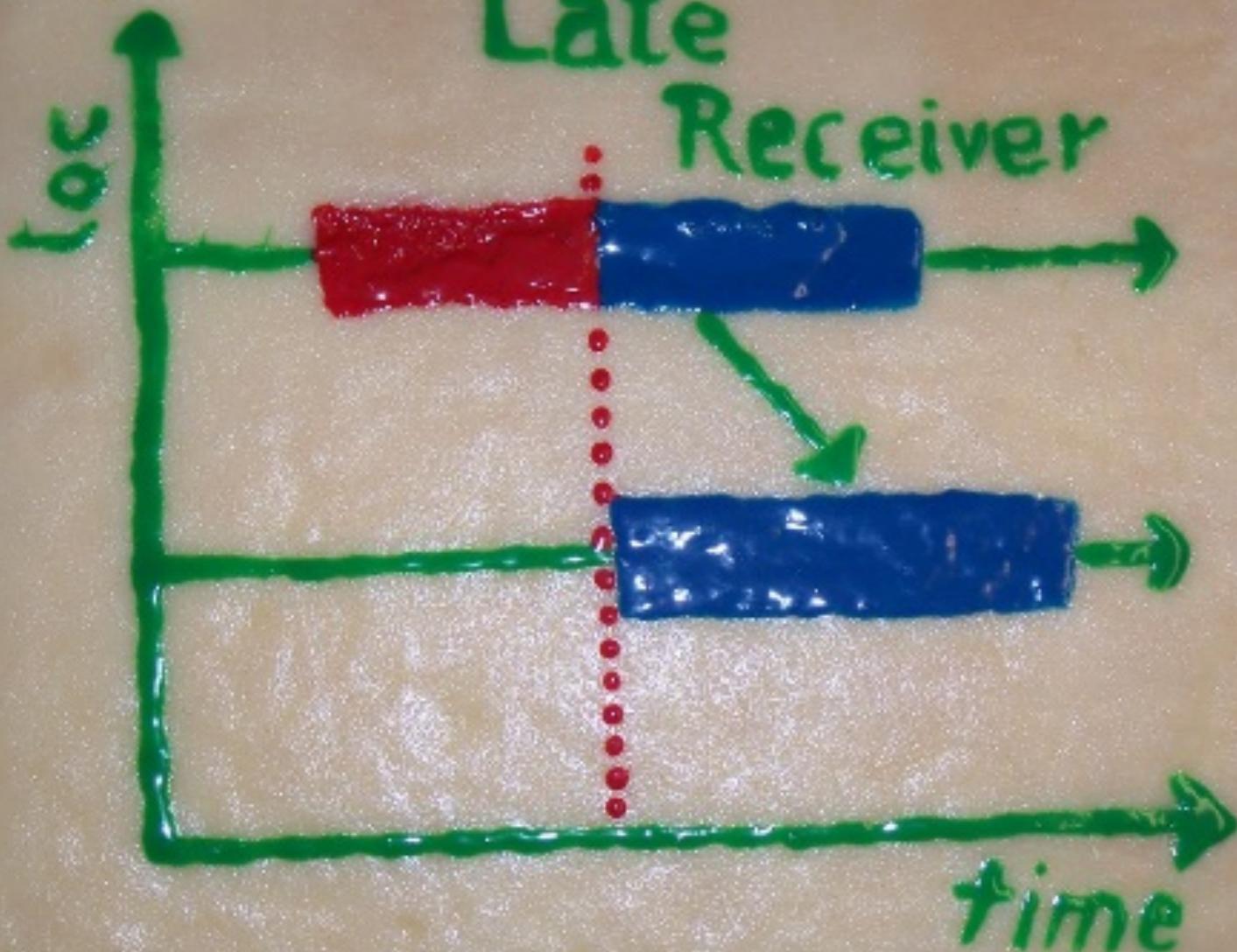


- Late Broadcast



[2014]

Late
Receiver



Questions?



scalasca 

- <http://www.scalasca.org>
- scalasca@fz-juelich.de



**Score-P**
Scalable performance measurement
infrastructure for parallel codes

- <http://www.score-p.org>
- support@score-p.org

