



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



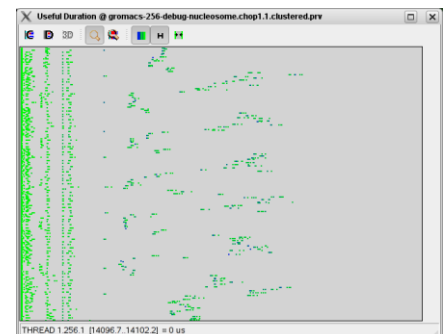
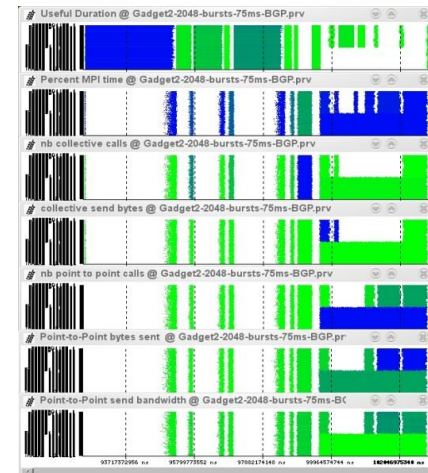
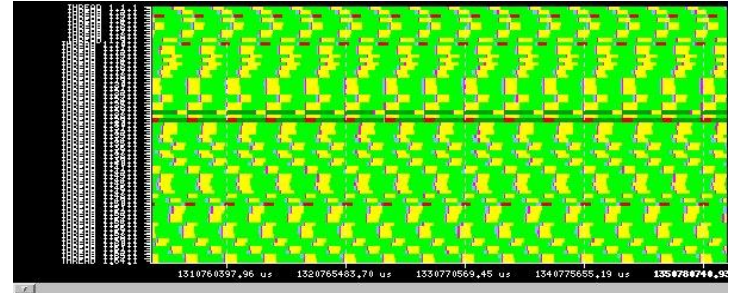
Understanding applications with Paraver

tools@bsc.es

2018

Our Tools

- Since 1991
- Based on traces
- Open Source
- <http://tools.bsc.es>
- Core tools:
 - Paraver (paramedir) – offline trace analysis
 - Dimemas – message passing simulator
 - Extrae – instrumentation
- Focus
 - Detail, variability, flexibility
 - Behavioral structure vs. syntactic structure
 - Intelligence: Performance Analytics



Paraver



**Barcelona
Supercomputing
Center**

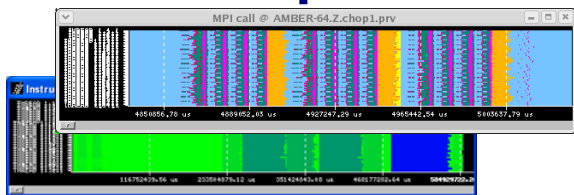
Centro Nacional de Supercomputación

Paraver – Performance data browser



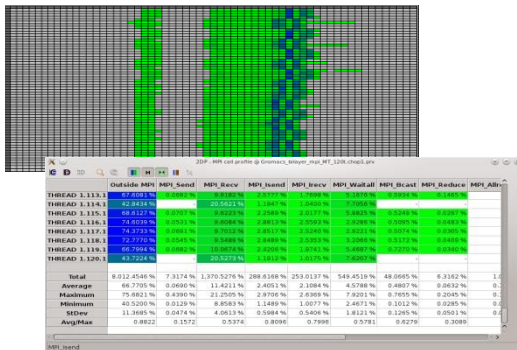
Trace visualization/analysis

+ trace manipulation



Timelines

Goal = Flexibility
No semantics
Programmable

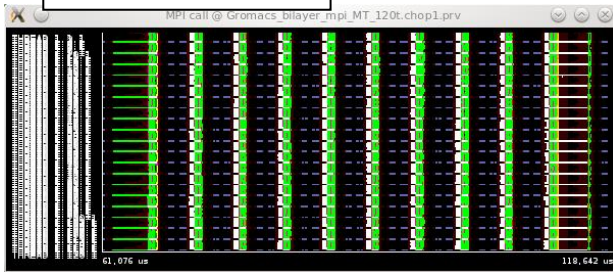


2/3D tables
(Statistics)

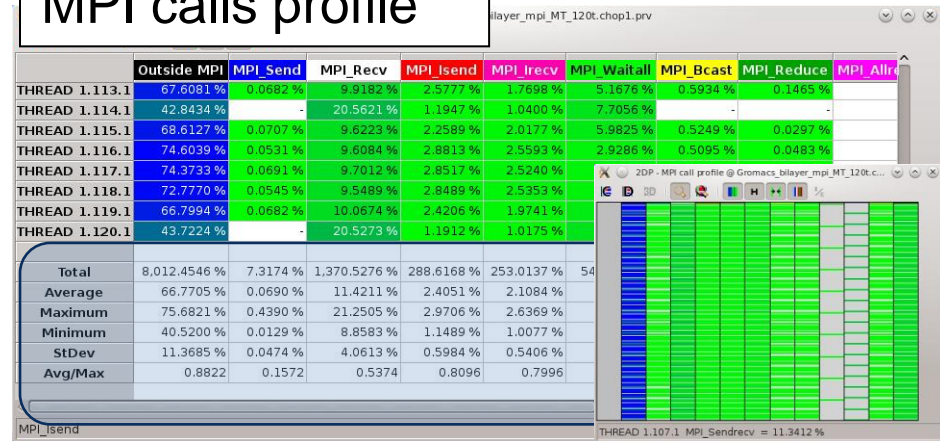
Comparative analyses
Multiple traces
Synchronize scales

From timelines to tables

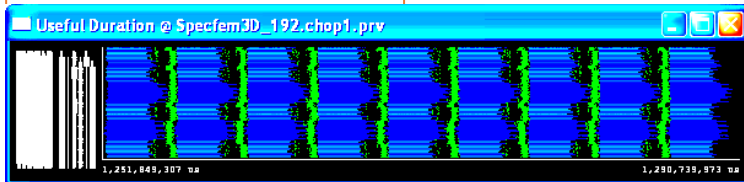
MPI calls



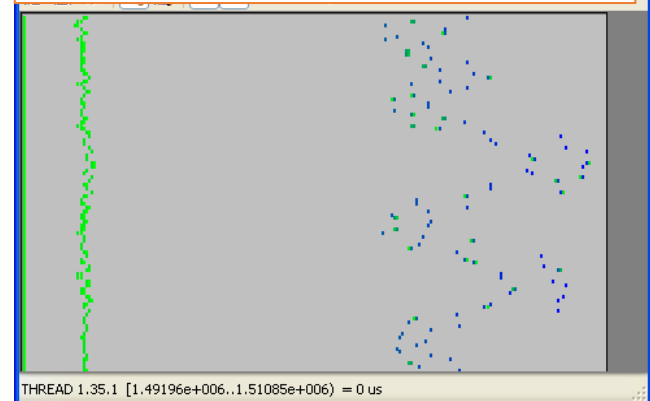
MPI calls profile



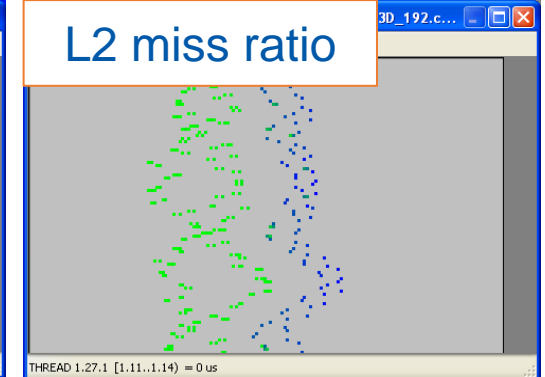
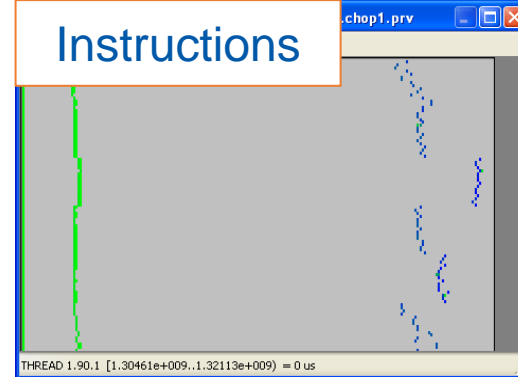
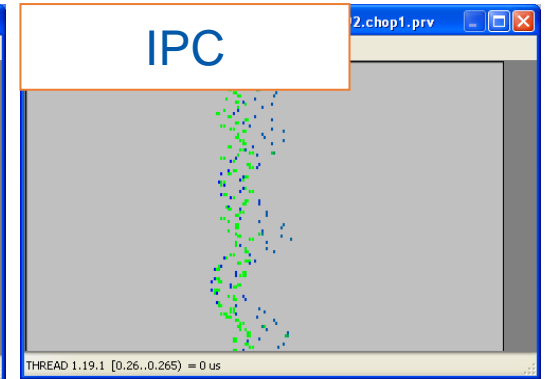
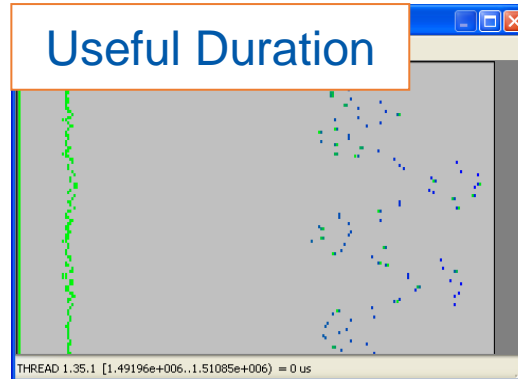
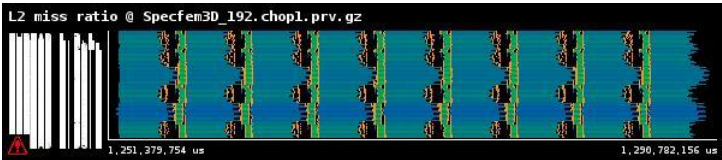
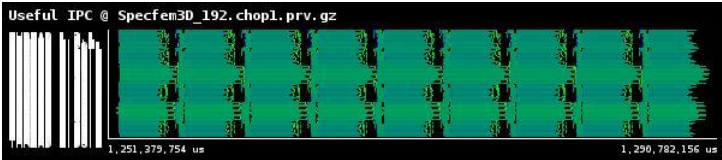
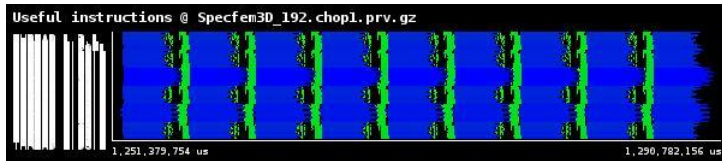
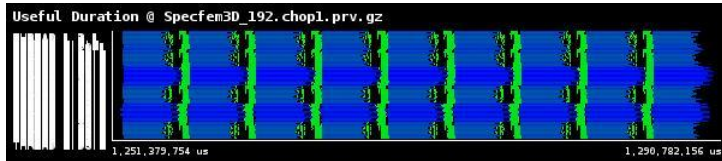
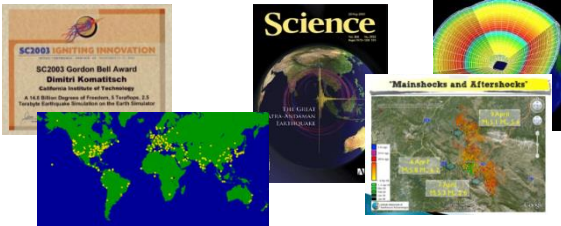
Useful Duration



Histogram Useful Duration

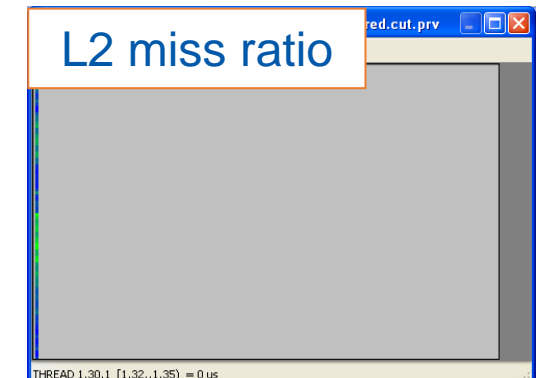
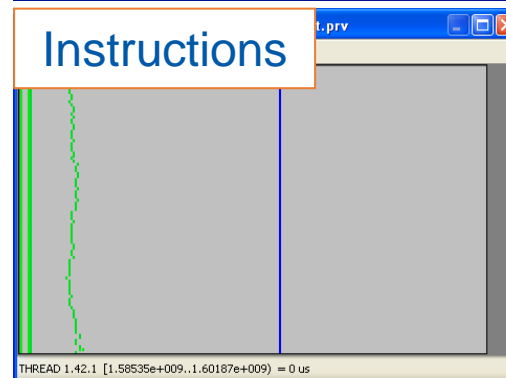
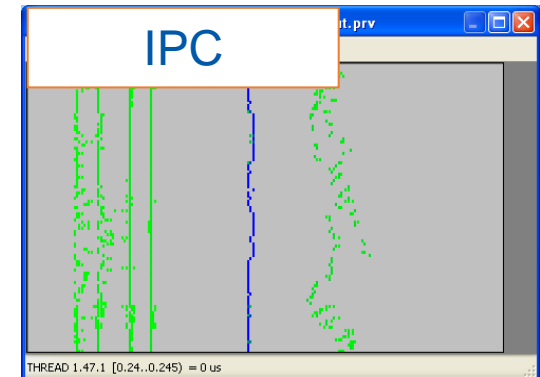
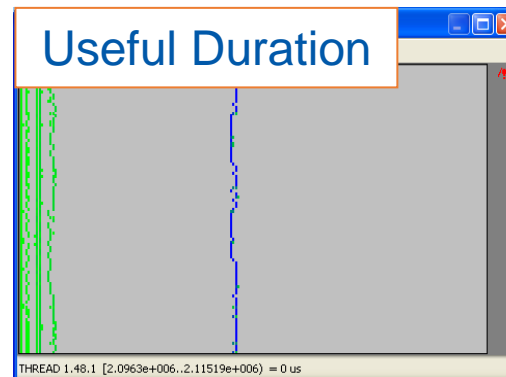


Analyzing variability



Analyzing variability

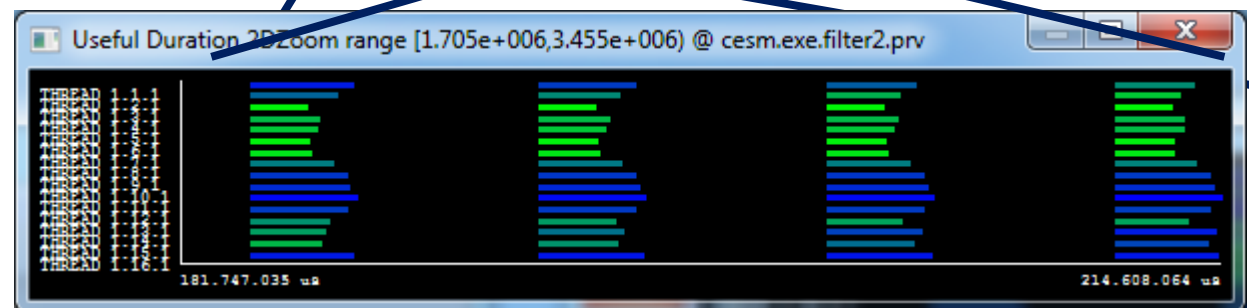
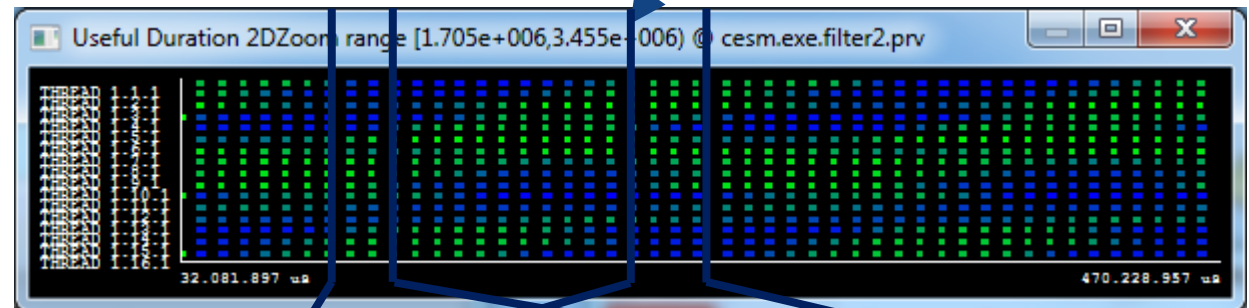
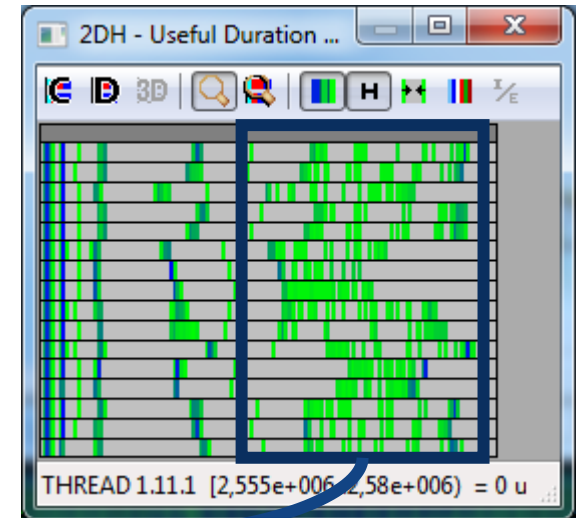
- By the way: six months later



From tables to timelines

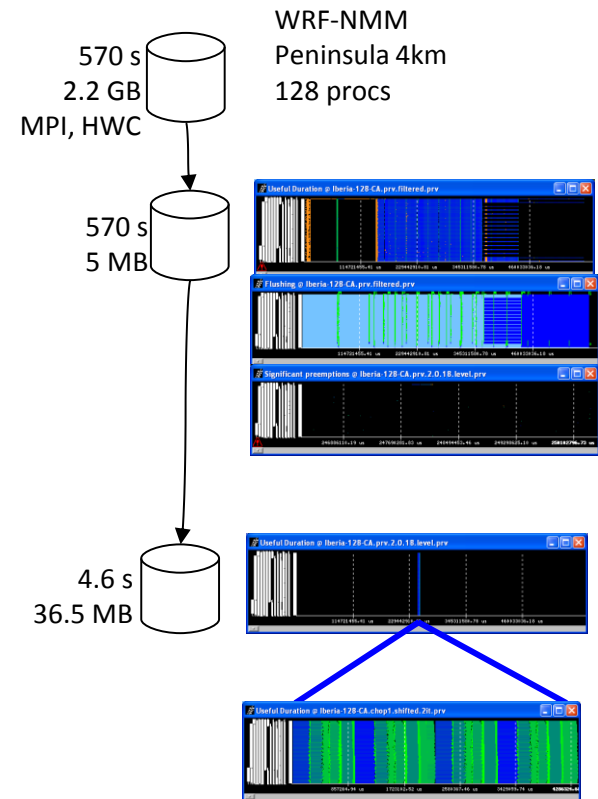
CESM: 16 processes, 2 simulated days

- Histogram useful computation duration shows high variability
- How is it distributed?
- Dynamic imbalance
 - In space and time
 - Day and night.
 - Season ? ☺



Trace manipulation

- Data handling/summarization capability
 - Filtering
 - Subset of records in original trace
 - By duration, type, value,...
 - Filtered trace IS a paraver trace and can be analysed with the same cfgs (as long as needed data kept)
 - Cutting
 - All records in a given time interval
 - Only some processes
 - Software counters
 - Summarized values computed from those in the original trace emitted as new even types
 - #MPI calls, total hardware count,...



Extrae



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Extræ features

- Platforms
 - Intel, Cray, BlueGene, MIC, ARM, Android, Fujitsu Sparc...
- Parallel programming models
 - MPI, OpenMP, pthreads, OmpSs, CUDA, OpenCL, Java, Python...
- Performance Counters
 - Using PAPI interface
- Link to source code
 - Callstack at MPI routines
 - OpenMP outlined routines
 - Selected user functions (Dyninst)
- Periodic sampling
- User events (Extræ API)

**No need
to
recompile
/ relink!**

Extrae overheads

	Average values
Event	150 – 200 ns
Event + PAPI	750 ns – 1.5 us
Event + callstack (1 level)	1 us
Event + callstack (6 levels)	2 us

How does Extrae work?

- Symbol substitution through LD_PRELOAD
 - Specific libraries for each combination of runtimes
 - MPI
 - OpenMP
 - OpenMP+MPI
 - ...
- Dynamic instrumentation
 - Based on Dyninst (developed by U.Wisconsin / U.Maryland)
 - Instrumentation in memory
 - Binary rewriting
- Alternatives
 - Static link (i.e., PMPI, Extrae API)

Recommended

Extrae XML configuration

```
<mpi enabled="yes">
  <counters enabled="yes" />
</mpi>

<openmp enabled="yes">
  <locks enabled="no" />
  <counters enabled="yes" />
</openmp>

<pthread enabled="no">
  <locks enabled="no" />
  <counters enabled="yes" />
</pthread>

<callers enabled="yes">
  <mpi enabled="yes">1-3</mpi>
  <sampling enabled="no">1-5</sampling>
</callers>
```

Trace the MPI calls
(What's the program doing?)

Trace the call-stack
(Where in my code?)

Extrae XML configuration (II)

```
<counters enabled="yes">
  <cpu enabled="yes" starting-set-distribution="1">
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_L1_DCM, PAPI_L2_DCM, PAPI_L3_TCM
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_BR_MSP, PAPI_BR_UCN,
      PAPI_BR_CN, RESOURCE_STALLS
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_VEC_DP, PAPI_VEC_SP,
      PAPI_FP_INS
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_LD_INS, PAPI_SR_INS
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, RESOURCE_STALLS:LOAD,
      RESOURCE_STALLS:STORE, RESOURCE_STALLS:ROB_FULL, RESOURCE_STALLS:RS_FULL
    </set>
  </cpu>
  <network enabled="no" />
  <resource-usage enabled="no" />
  <memory-usage enabled="no" />
</counters>
```

Select which
HW counters
are measured
(How's the machine doing?)

Extrae XML configuration (III)

```
<buffer enabled="yes">  
  <size enabled="yes">500000</size>  
  <circular enabled="no" />  
</buffer>
```

Trace buffer size
(Flush/memory trade-off)

```
<sampling enabled="no" type="default" period="50m" variability="10m" />
```

Enable sampling
(Want more details?)

```
<merge enabled="yes"  
  synchronization="default"  
  tree-fan-out="16"  
  max-memory="512"  
  joint-states="yes"  
  keep-mpits="yes"  
  sort-addresses="yes"  
  overwrite="yes"
```

**Automatic
post-processing
to generate the
Paraver trace**

```
>  
  $TRACE_NAME$  
</merge>
```


Dimemas

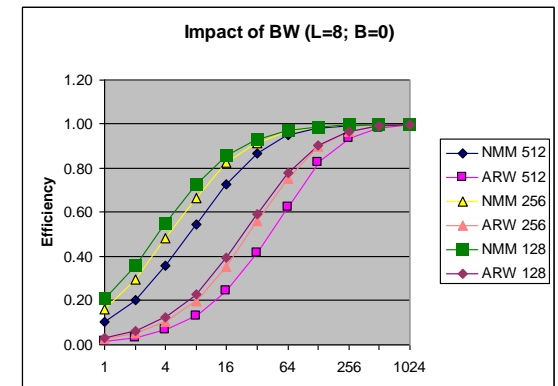
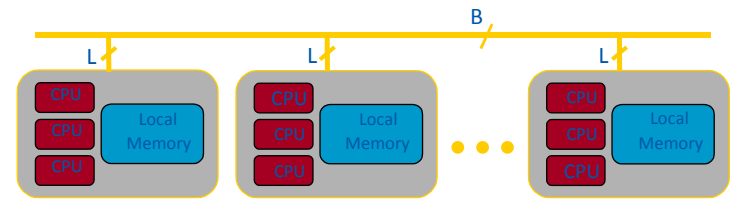


**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

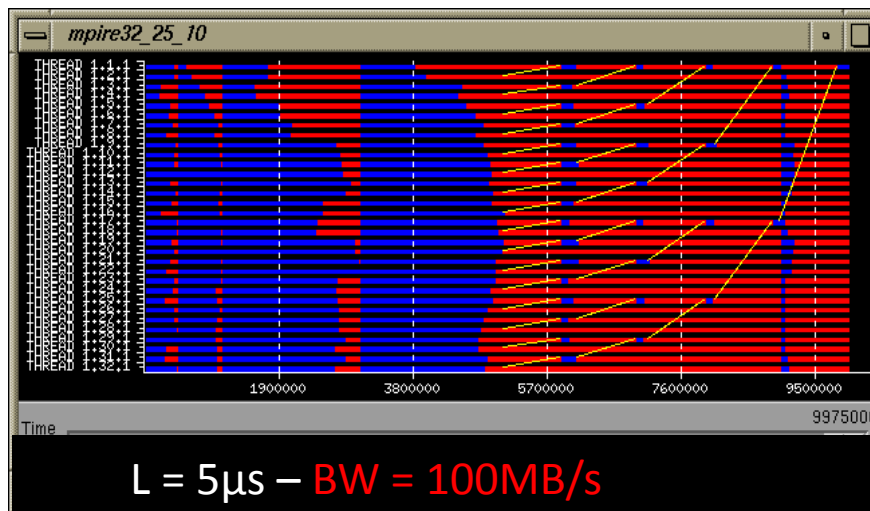
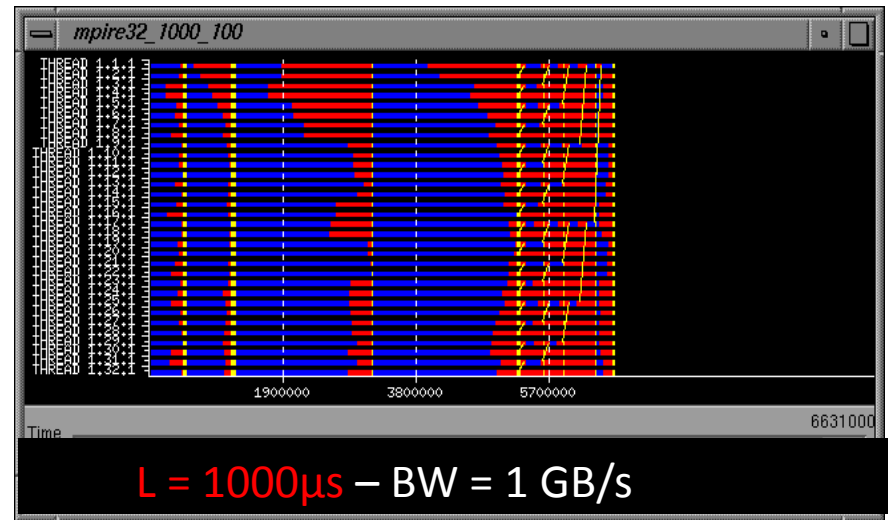
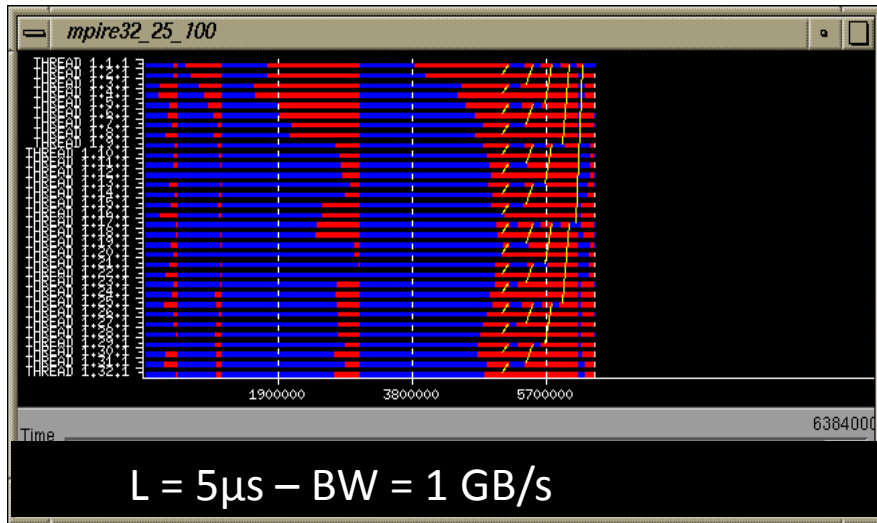
Dimemas – Coarse grain, Trace driven simulation

- Simulation: Highly non linear model
 - MPI protocols, resource contention...
- Parametric sweeps
 - On abstract architectures
 - On application computational regions
- What if analysis
 - Ideal machine (instantaneous network)
 - Estimating impact of ports to MPI+OpenMP/CUDA/...
 - Should I use asynchronous communications?
 - Are all parts equally sensitive to network?
- MPI sanity check
 - Modeling nominal
- Paraver – Dimemas tandem
 - Analysis and prediction
 - What-if from selected time window



Network sensitivity

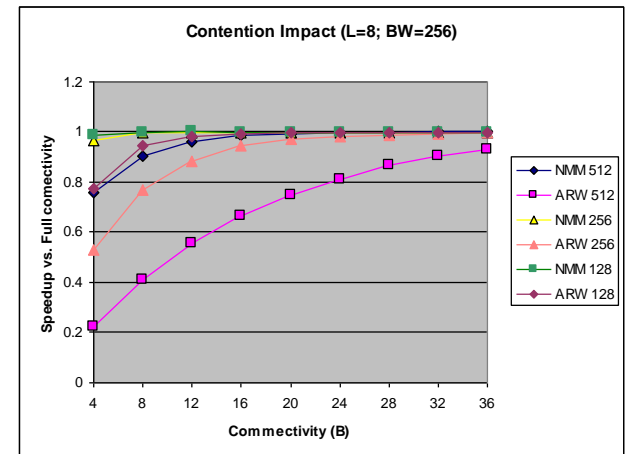
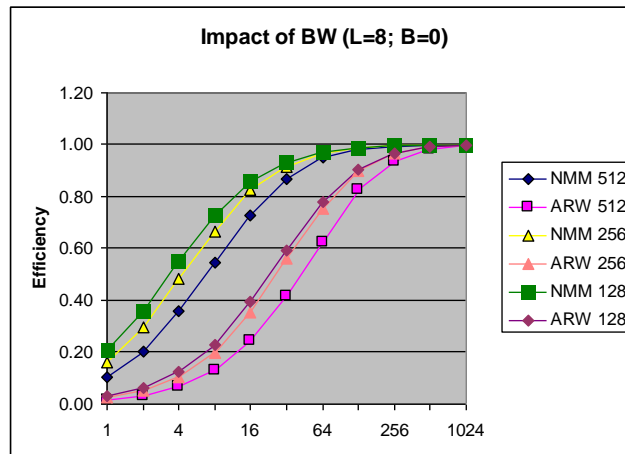
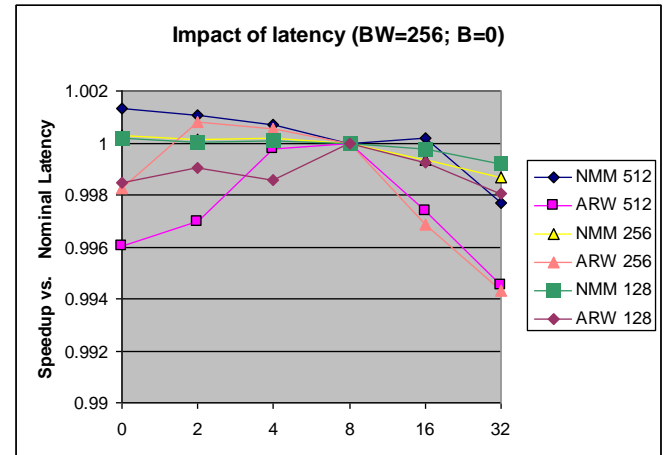
- MPIRE 32 tasks, no network contention



All windows same scale

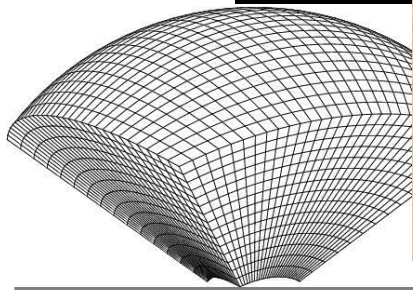
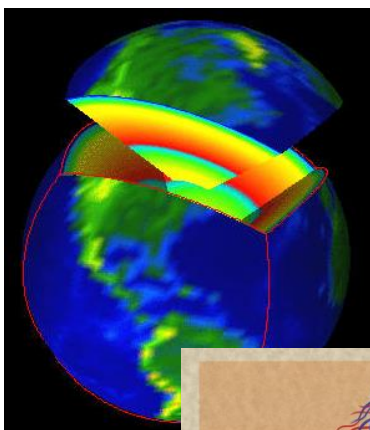
Network sensitivity

- WRF, Iberia 4Km, 4 procs/node
 - Not sensitive to latency
 - NMM
 - BW – 256MB/s
 - 512 – sensitive to contention
 - ARW
 - BW - 1GB/s
 - Sensitive to contention

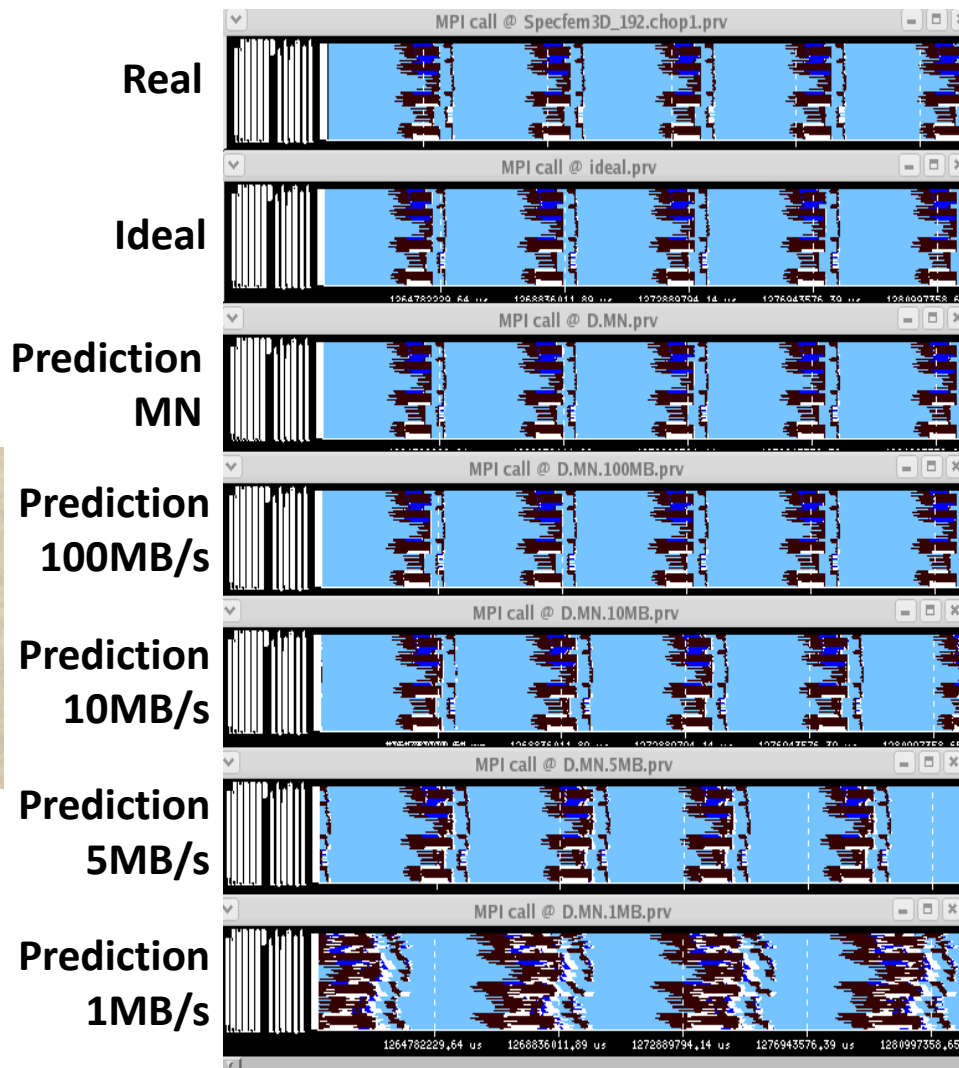


Would I will benefit from asynchronous communications?

SPECFEM3D



Courtesy Dimitri Komatitsch



Ideal machine

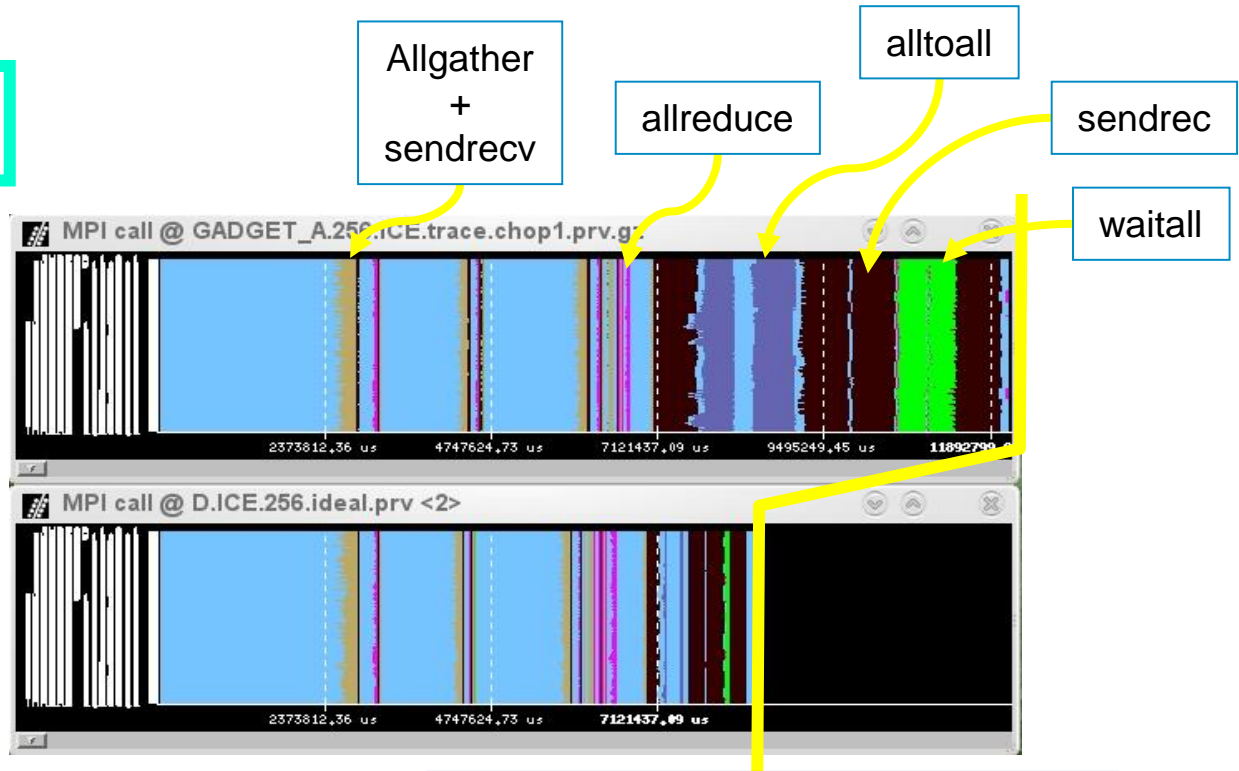
The impossible machine: $BW = \infty$, $L = 0$

- Actually describes/characterizes Intrinsic application behavior
 - Load balance problems?
 - Dependence problems?

GADGET @ Nehalem cluster
256 processes

Real
run

Ideal
network

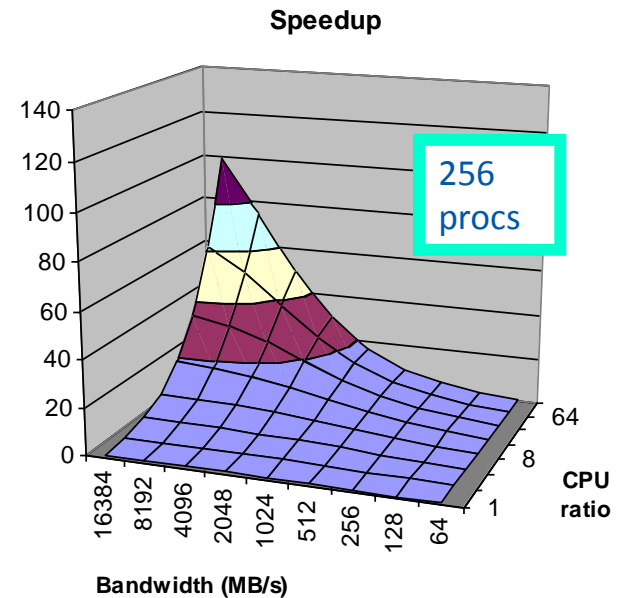
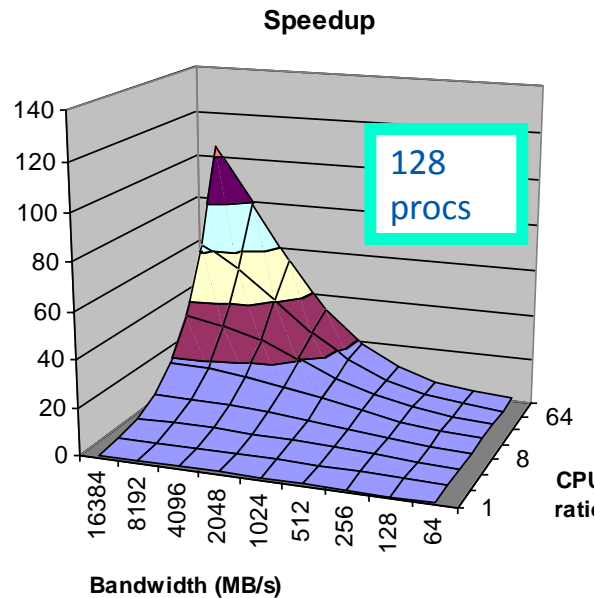
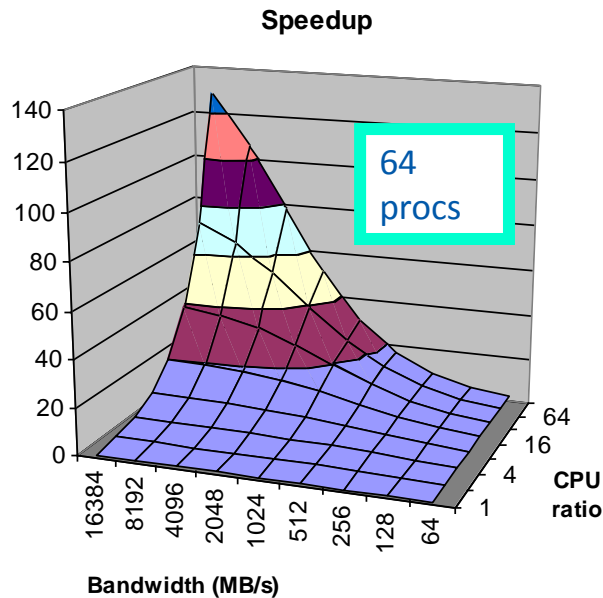


Impact on practical machines?

Impact of architectural parameters

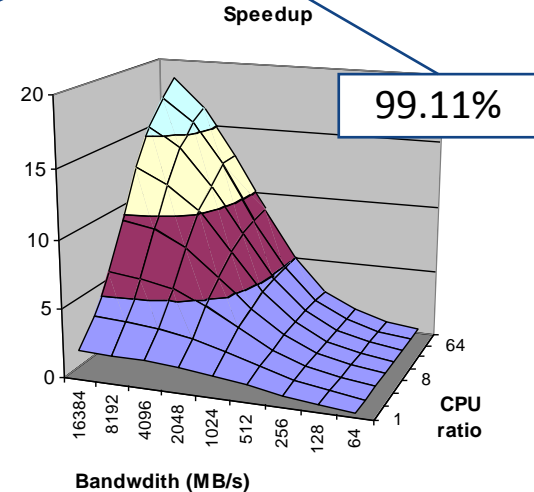
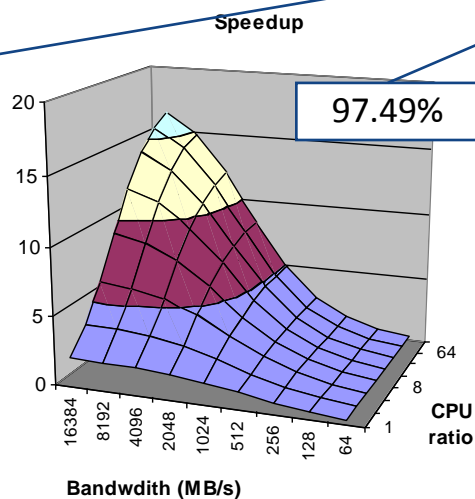
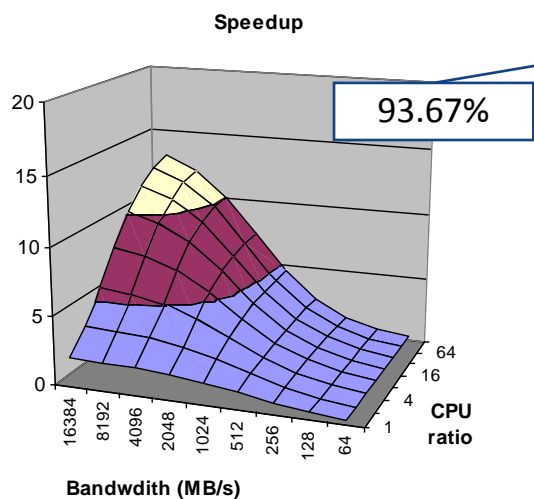
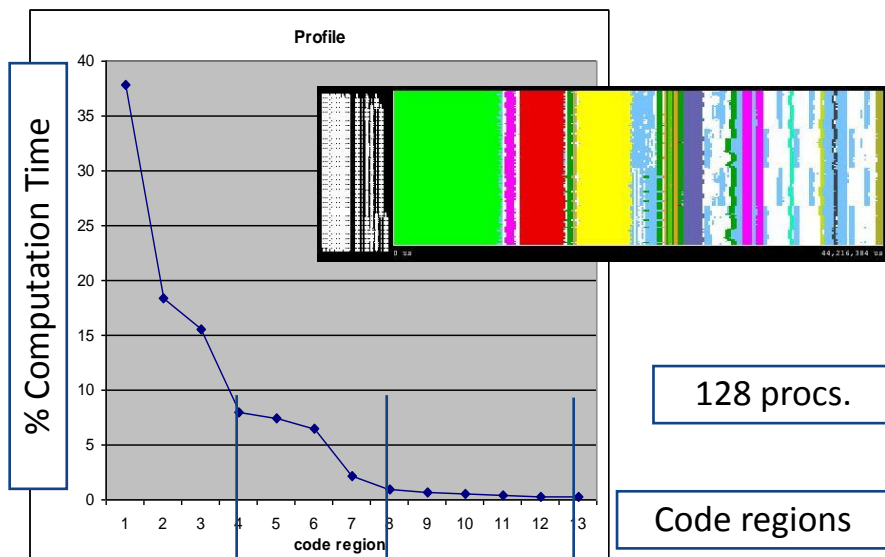
- **Ideal speeding up ALL** the computation bursts by the CPU ratio factor
 - The more processes the less speedup (higher impact of bandwidth limitations) !!

GADGET



Hybrid parallelization

- Hybrid/accelerator parallelization
 - Speed-up SELECTED regions by the CPUratio factor



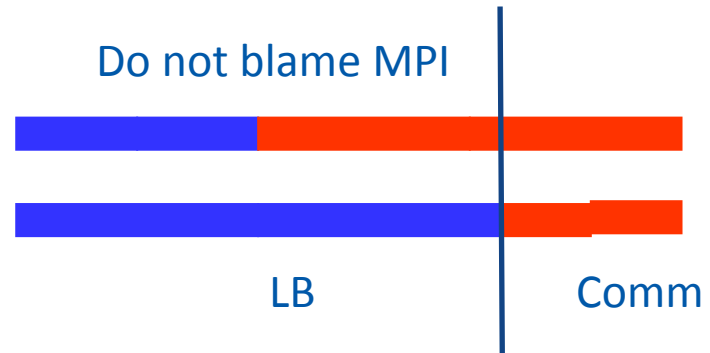
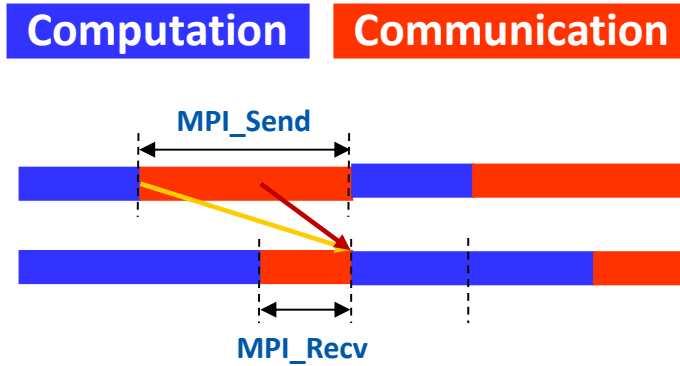
(Previous slide: speedups up to 100x)

Efficiency Models



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación

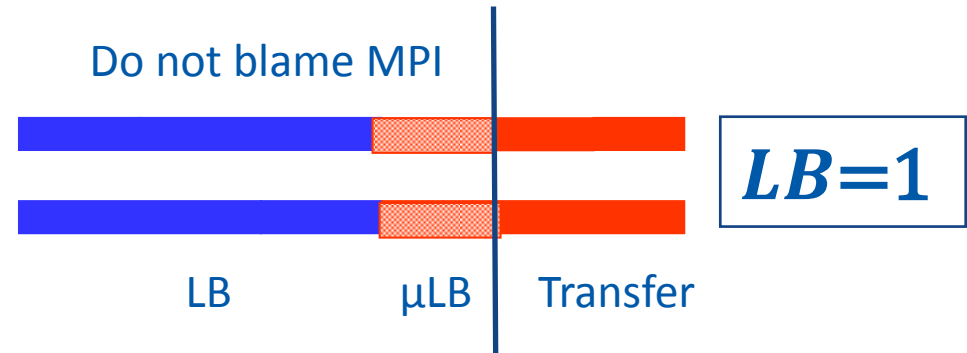
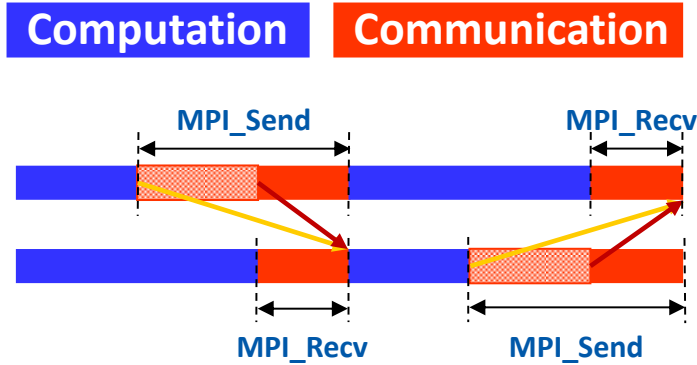
Parallel efficiency model



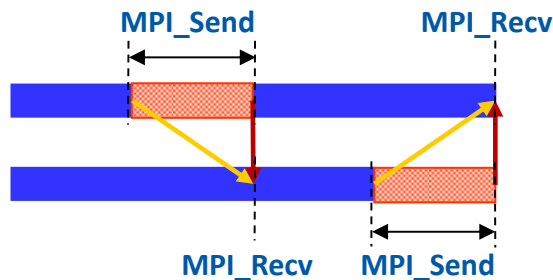
- Parallel efficiency = LB eff * Comm eff

Parallel efficiency refinement:

$$LB * \mu LB * Tr$$



- Serializations / dependences (μLB)
- Dimemas ideal network \rightarrow Transfer (efficiency) = 1

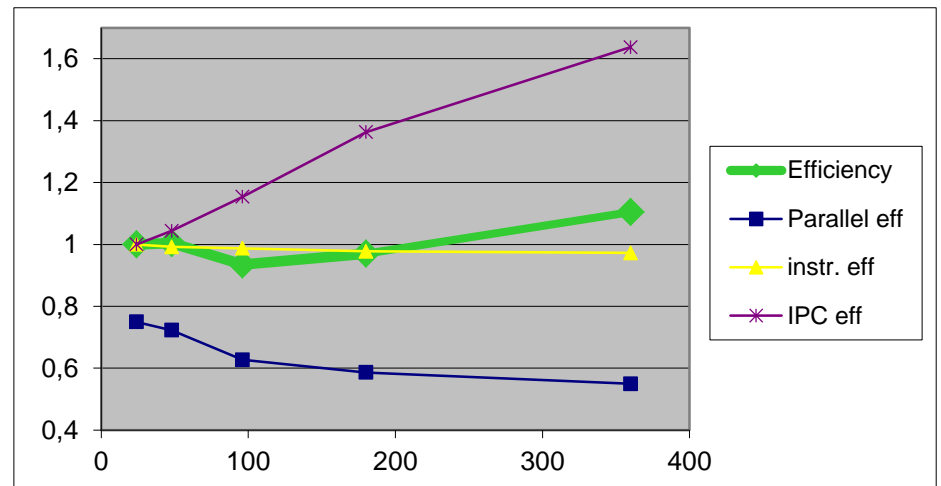
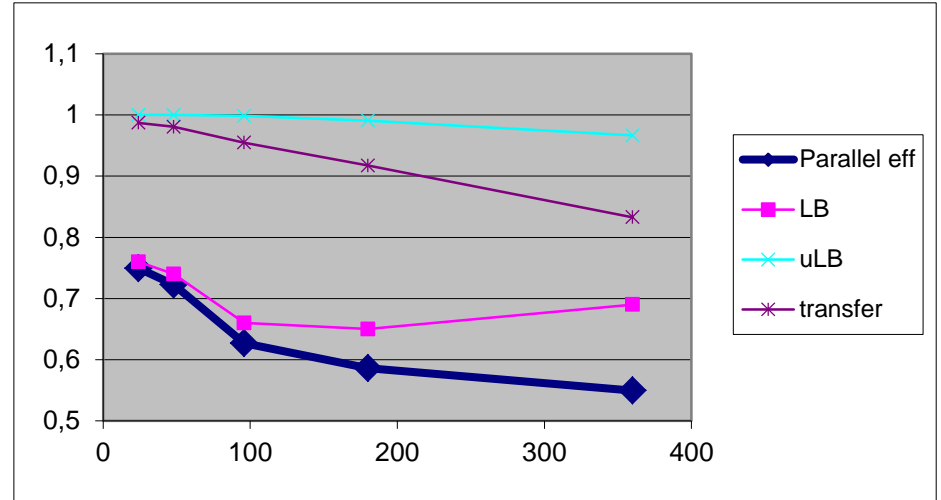
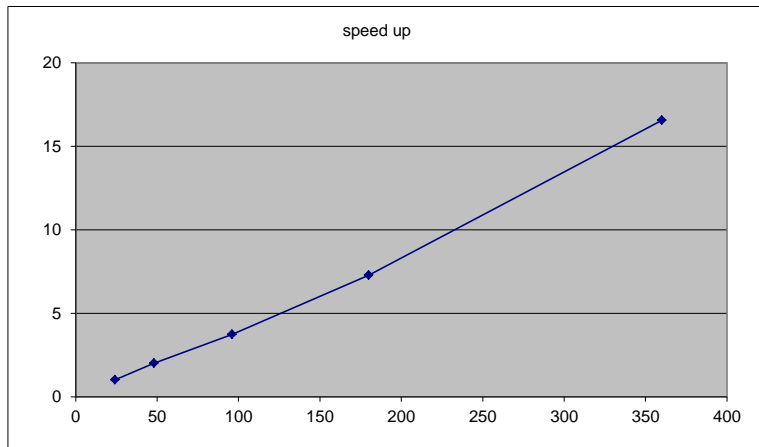


Why scaling?

$$\eta_{\parallel} = LB * Ser * Trf$$

CG-POP mpi2s1D - 180x120

Good scalability !!
Should we be happy?



Some examples of efficiencies

Code	Parallel efficiency	Communication efficiency	Load Balance efficiency
Gromacs@mt	66.77	75.68	88.22
BigDFT@altamira	59.64	78.97	75.52
CG-POP@mt	80.98	98.92	81.86
ntchem_mini@pi	92.56	94.94	97.49
nicam@pi	87.10	75.97	89.22
cp2k@jureca	75.34	81.07	92.93
icon@mistral	79.86	84.02	95.05
k-Wave@salomon	89.08	92.84	95.96
fleur@claix	76.22	90.66	84.07

Same code, different behaviour

Code	Parallel efficiency	Communication efficiency	Load Balance efficiency
lulesh@mn3	90.55	99.22	91.26
lulesh@leftraru	69.15	99.12	69.76
lulesh@uv2 (mpt)	70.55	96.56	73.06
lulesh@uv2 (impi)	85.65	95.09	90.07
lulesh@mt	83.68	95.48	87.64
lulesh@cori	90.92	98.59	92.20
lulesh@thunderX	73.96	97.56	75.81
lulesh@jetson	75.48	88.84	84.06
lulesh@claix	77.28	92.33	83.70
lulesh@jureca	88.20	98.45	89.57
lulesh@mn4	86.59	98.77	87.67
lulesh@inti	88.16	98.65	89.36

Warning::: Higher parallel efficiency does not mean faster!

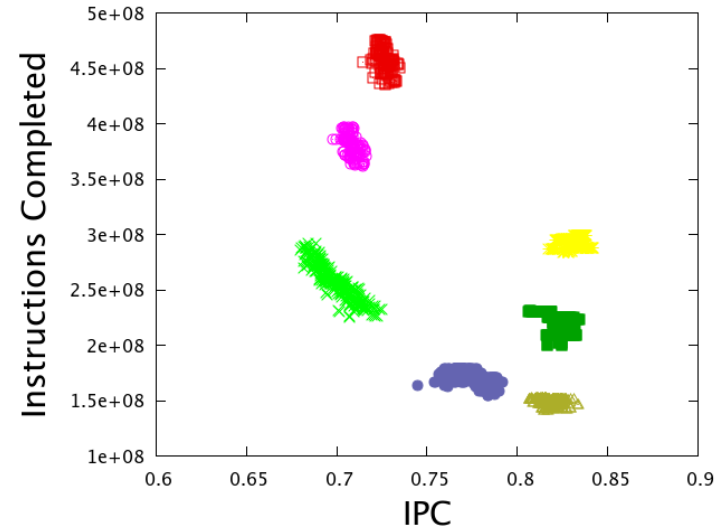
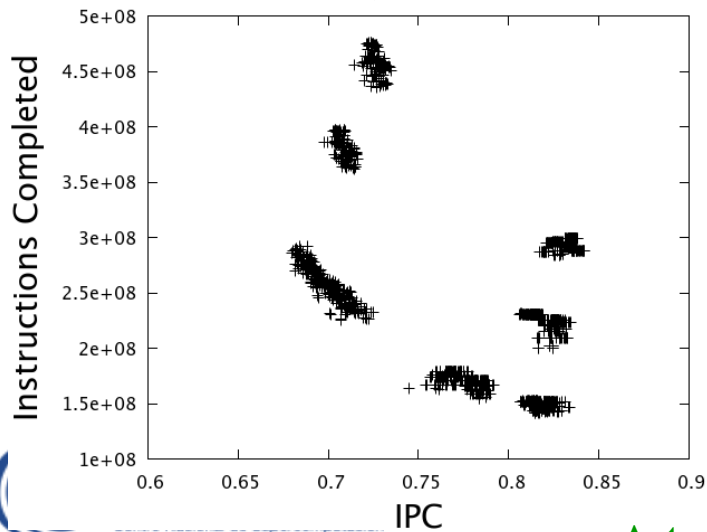
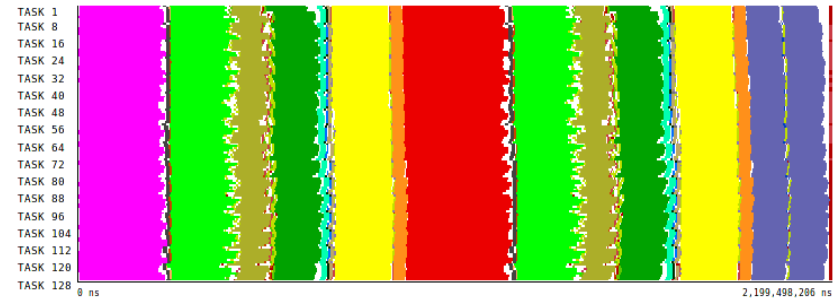
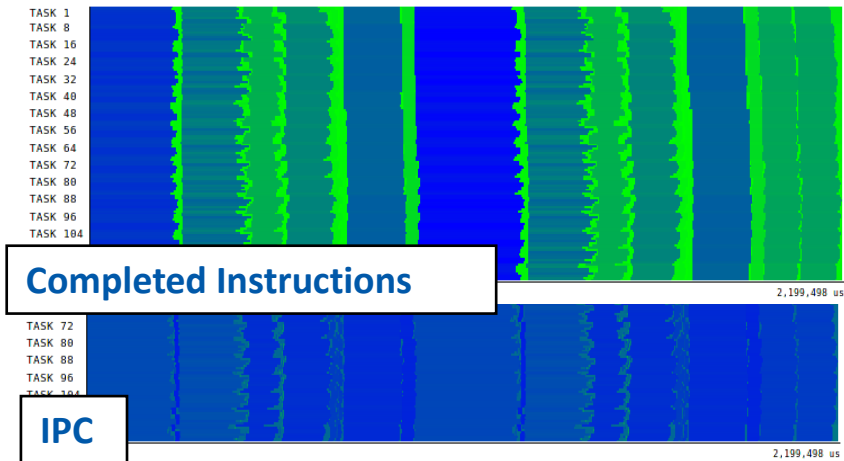
Analytics



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Using Clustering to identify structure



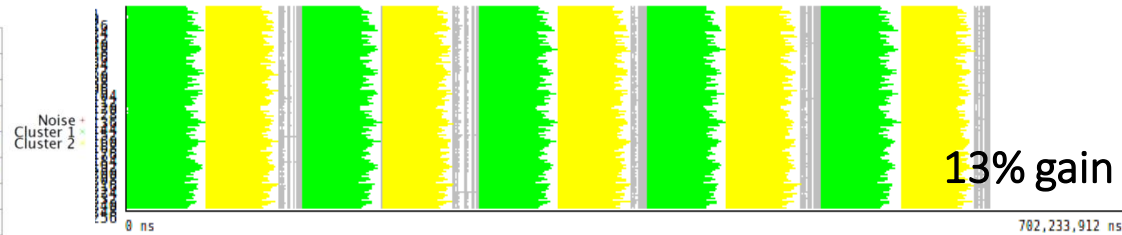
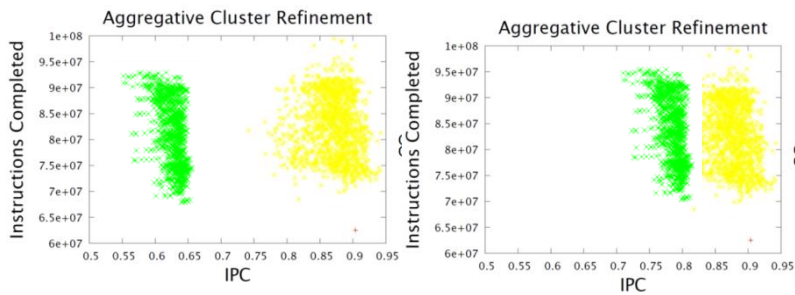
What should I improve?

What if ...

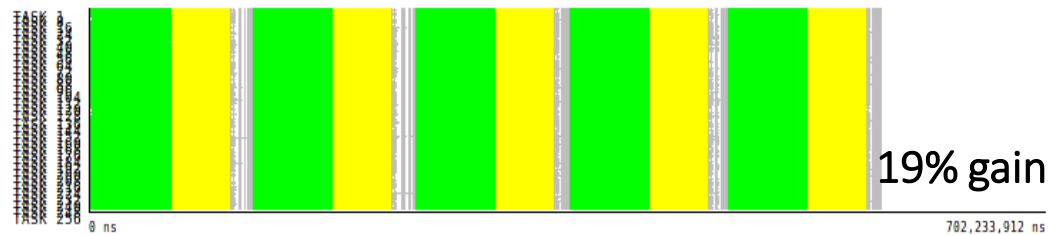
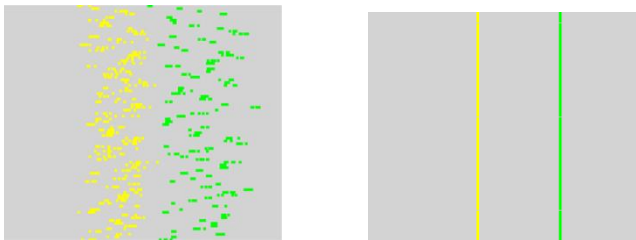
PEPC



... we increase the IPC of Cluster1?

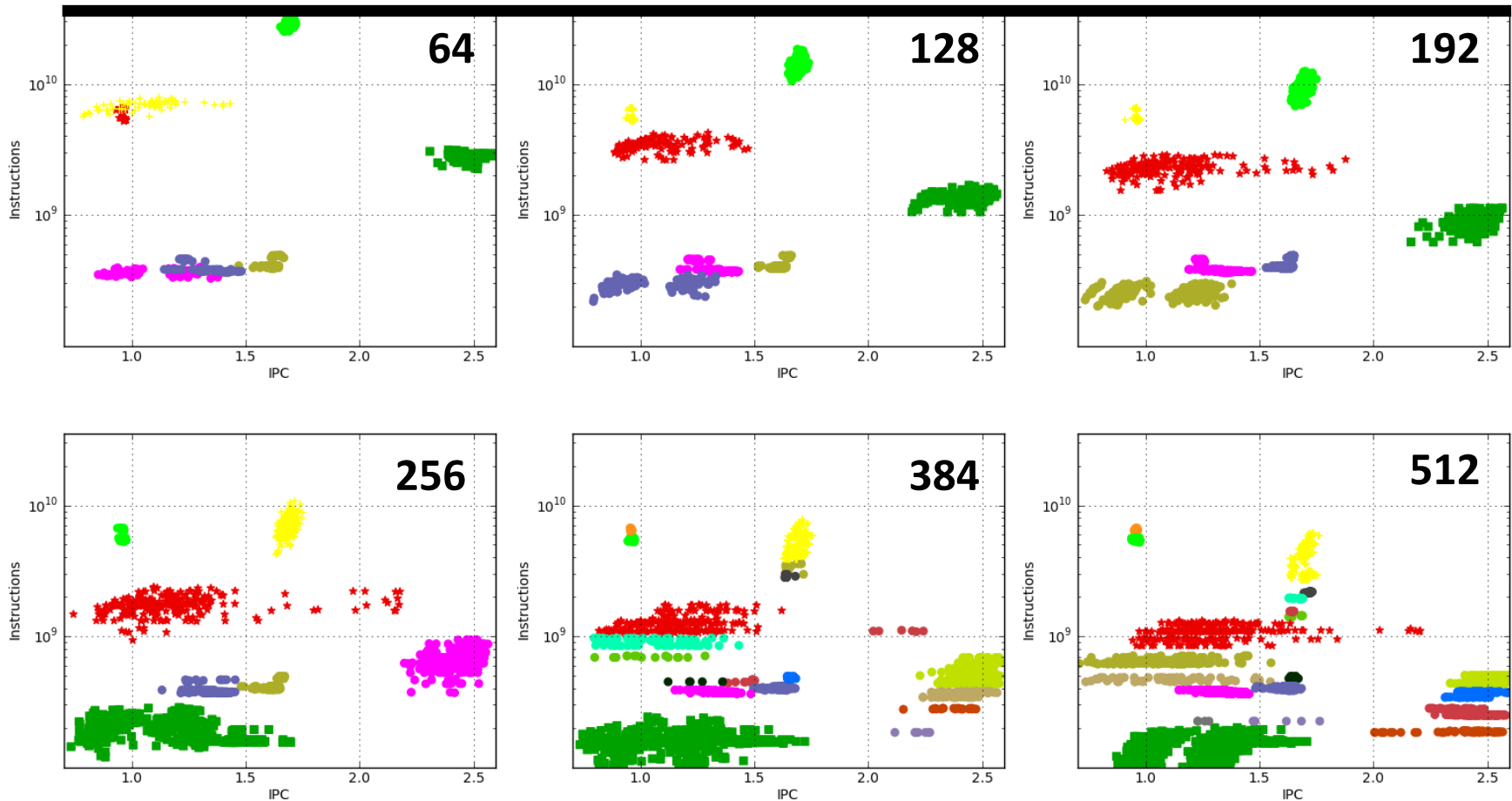


... we balance Clusters 1 & 2?



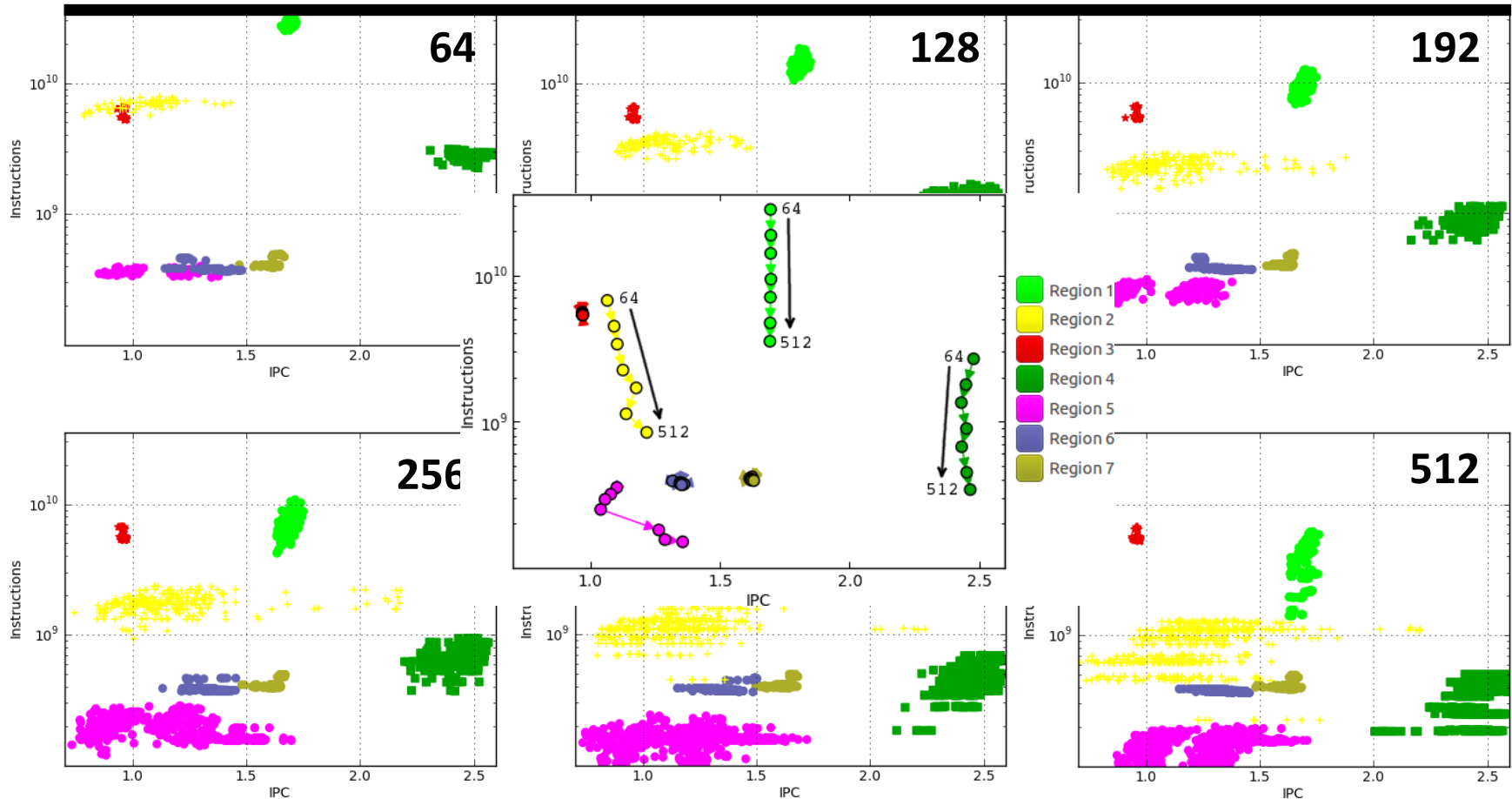
Tracking scability through clustering

- OpenMX (strong scale from 64 to 512 tasks)



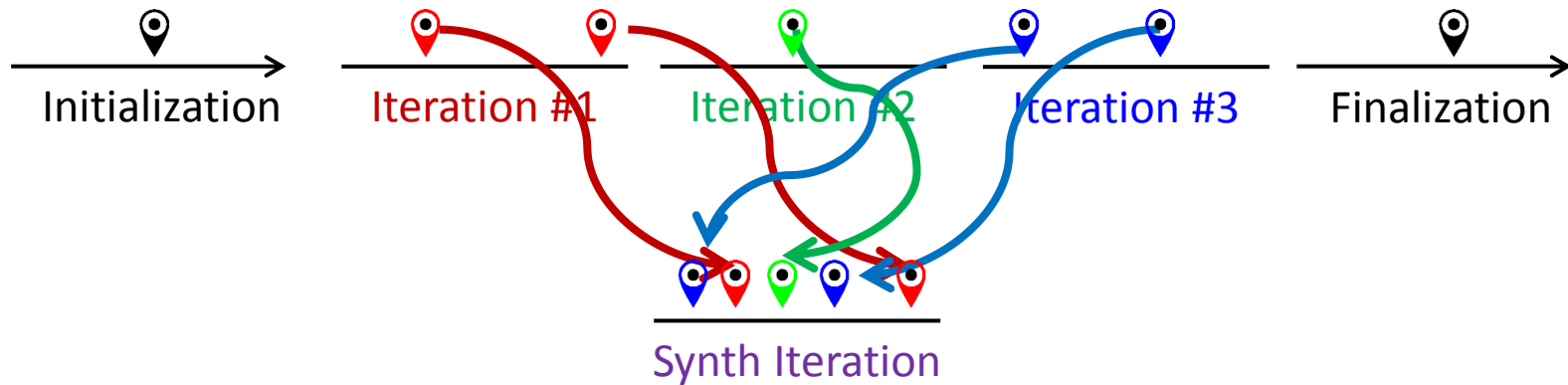
Tracking scalability through clustering

- OpenMX (strong scale from 64 to 512 tasks)



Folding

- Instantaneous metrics with minimum overhead
 - Combine instrumentation and sampling
 - Instrumentation delimits regions (routines, loops, ...)
 - Sampling exposes progression within a region
 - Captures performance counters and call-stack references

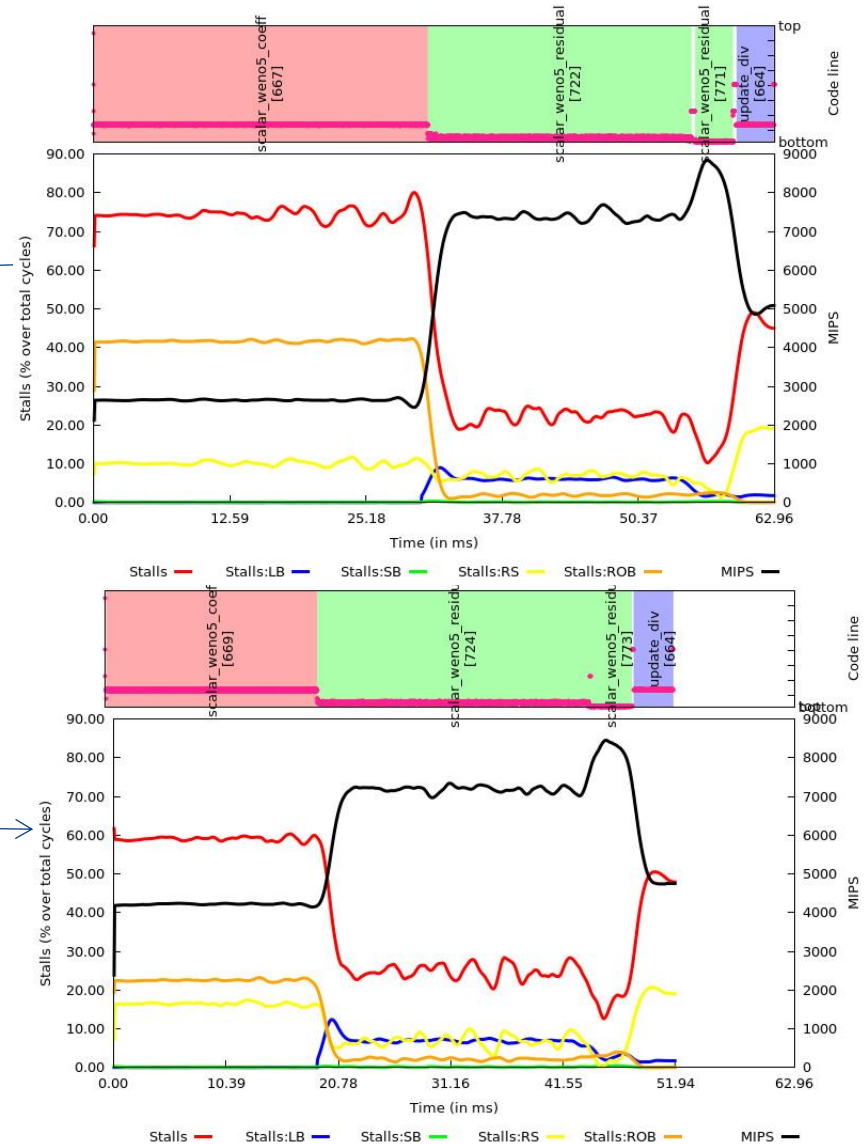


“Blind” optimization

- From folded samples of a few levels to timeline structure of “relevant” routines

Recommendation without access to source code

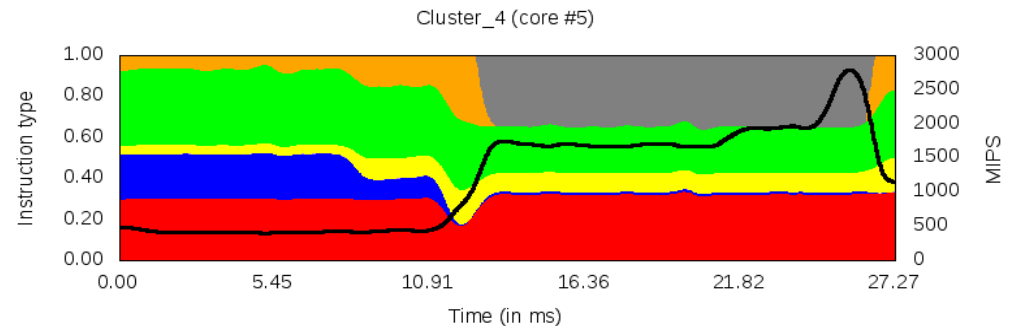
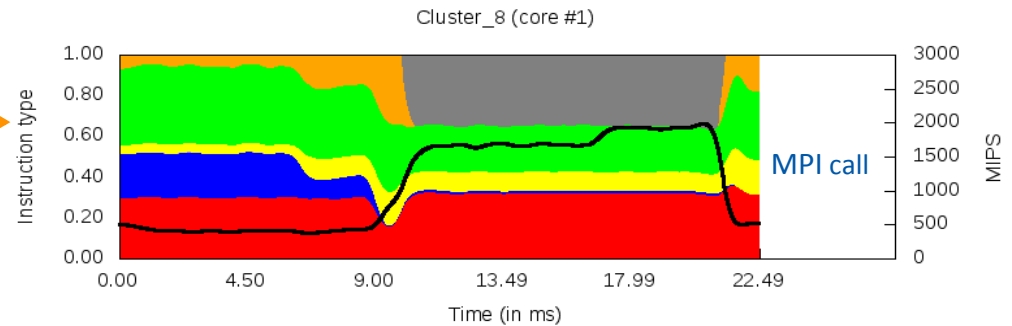
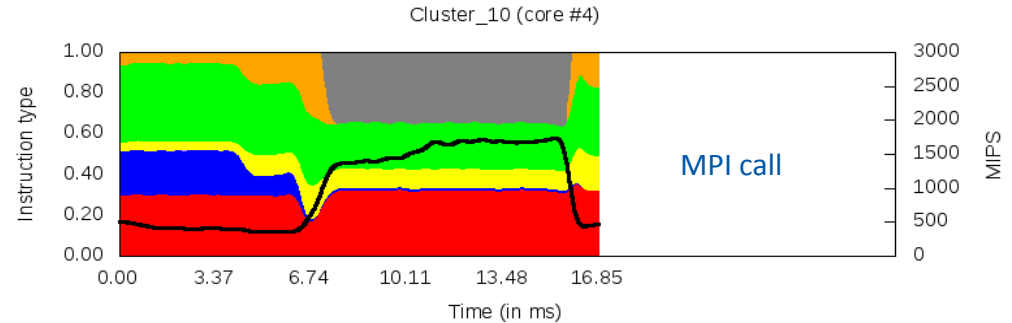
Evolution for Stall distribution model
Appl * Task * Thread * - Group_0 - Cluster_2



CG-POP multicore MN3 study

- Unbalanced MPI application
- Same code
- Different duration
- Different performance

Instruction mix model for the unbalanced CGPOP on different cores of the same hexacore chip



Methodology



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

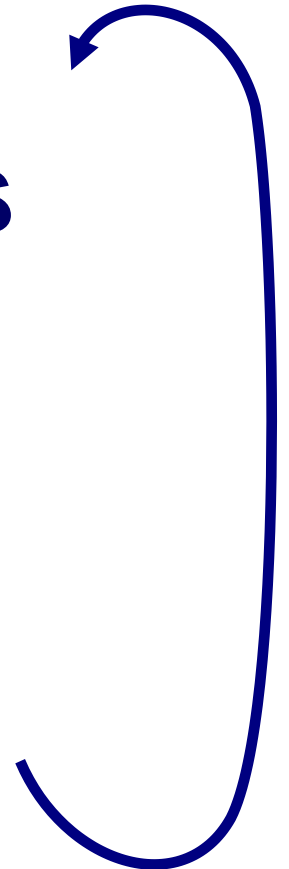
Performance analysis tools objective

Help generate hypotheses

Help validate hypotheses

Qualitatively

Quantitatively



First steps

- Parallel efficiency – percentage of time invested on computation
 - Identify sources for “inefficiency”:
 - load balance
 - Communication /synchronization
- Serial efficiency – how far from peak performance?
 - IPC, correlate with other counters
- Scalability – code replication?
 - Total #instructions
- Behavioral structure? Variability?

Paraver Tutorial:
Introduction to Paraver and Dimemas methodology

BSC Tools web site

- tools.bsc.es
 - downloads
 - Sources / Binaries
 - Linux / windows / MAC
 - documentation
 - Training guides
 - Tutorial slides
- Getting started
 - Start wxparaver
 - Help → tutorials and follow instructions
 - Follow training guides
 - Paraver introduction (MPI): Navigation and basic understanding of Paraver operation