

# Understanding applications using the BSC performance tools

---

---

Judit Gimenez (judit@bsc.es)  
German Llort(german.llort@bsc.es)

---

---

## Humans are visual creatures

---

- Films or books?
  - Two hours vs. days (months)
- Memorizing a deck of playing cards
  - Each card translated to an image (person, action, location)
- Our brain loves pattern recognition
  - What do you see on the pictures?

PROCESS

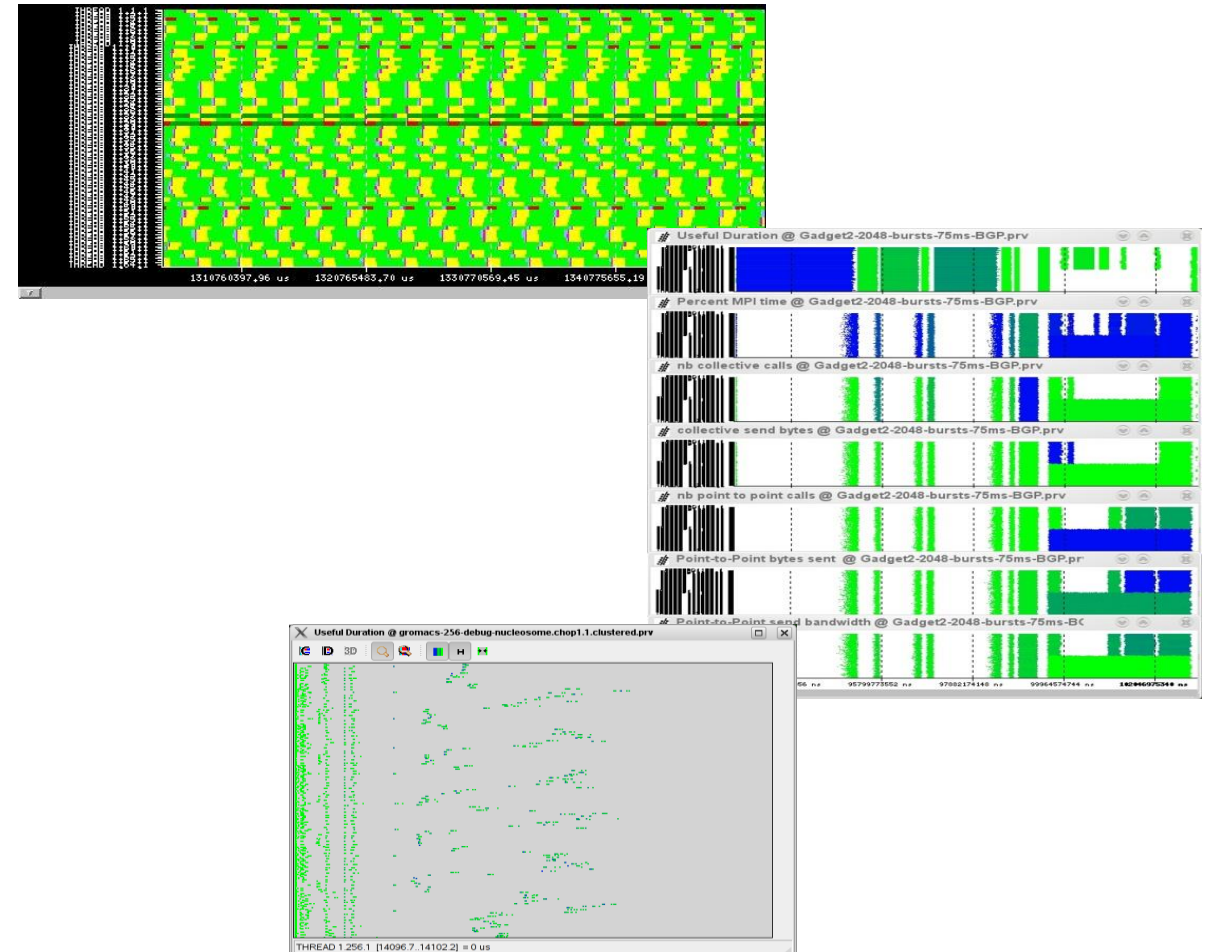
STORE

IDENTIFY



## Our Tools

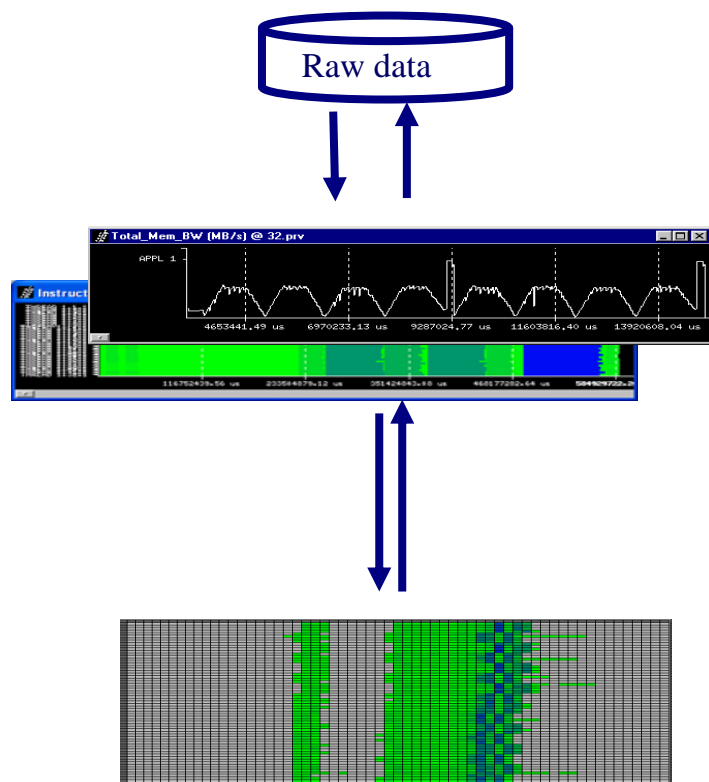
- Since 1991
- Based on traces
- Open Source
  - <http://www.bsc.es/paraver>
- Core tools:
  - Paraver (paramedir) – offline trace analysis
  - Dimemas – message passing simulator
  - Extrae – instrumentation
- Focus
  - Detail, variability, flexibility
  - Behavioral structure vs. syntactic structure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





# Paraver

## Paraver: Performance data browser



Timelines

2/3D tables  
(Statistics)

Trace visualization/analysis

+ trace manipulation

Goal = Flexibility

No semantics

Programmable

Comparative analyses

Multiple traces

Synchronize scales

## Timelines

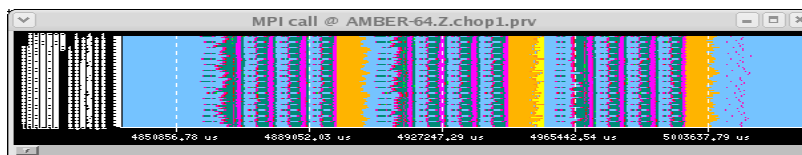
- Each window displays one view
  - Piecewise constant** function of time



$$s(t) = S_i, i \in [t_i, t_{i+1})$$

- Types of functions

- Categorical
  - State, user function, outlined routine

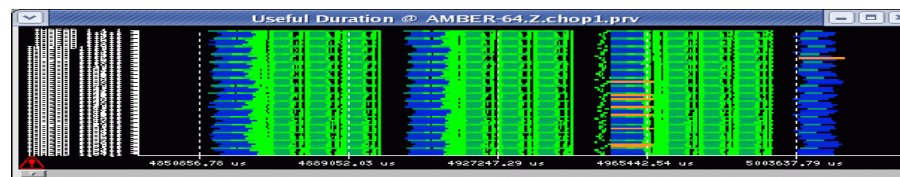


$$S_i \in [0, n] \subset N, \quad n <$$

- Logical
  - In specific user function, In MPI call, In long MPI call

$$S_i \in \{0, 1\}$$

- Numerical
  - IPC, L2 miss ratio, Duration of MPI call, duration of computation burst

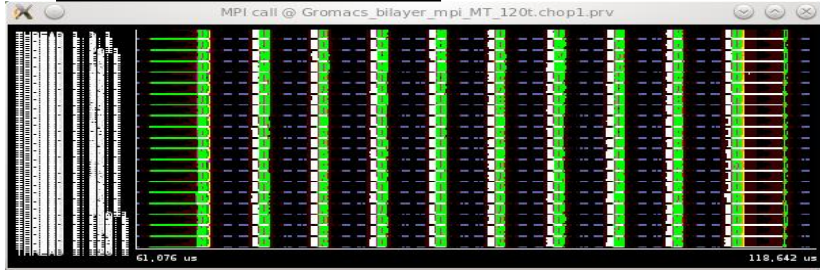


$$S_i \in R$$

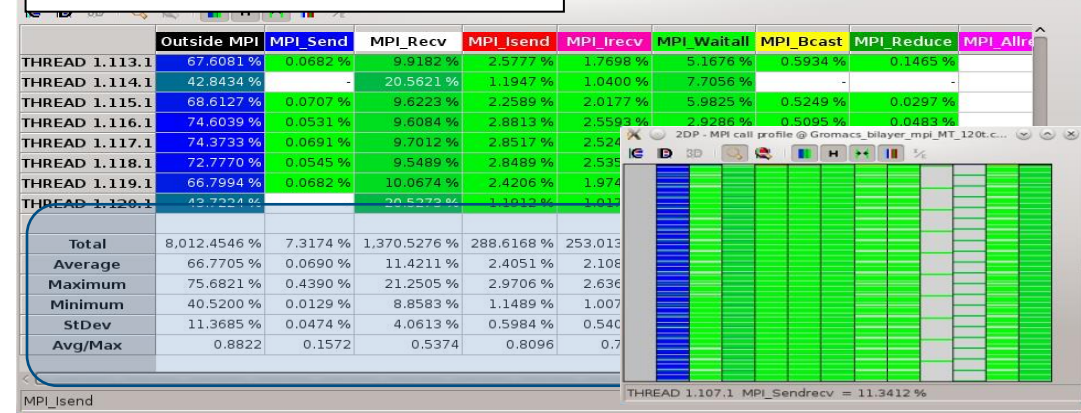
## Tables: Profiles, histograms, correlations

- From timelines to tables

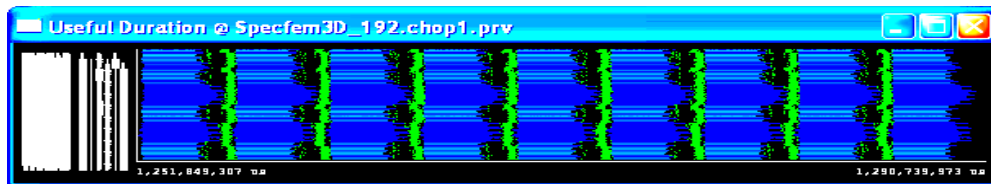
MPI calls



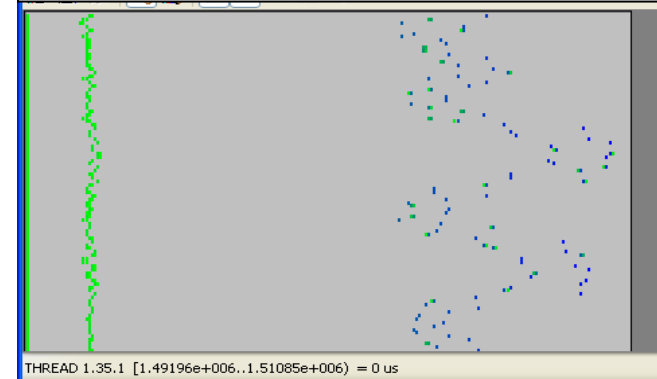
MPI calls profile



Useful Duration

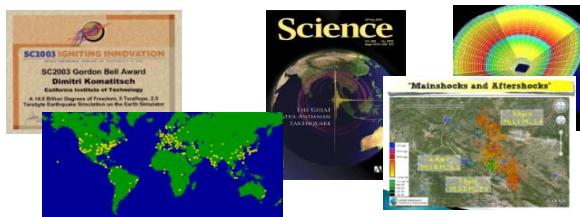


Histogram Useful Duration

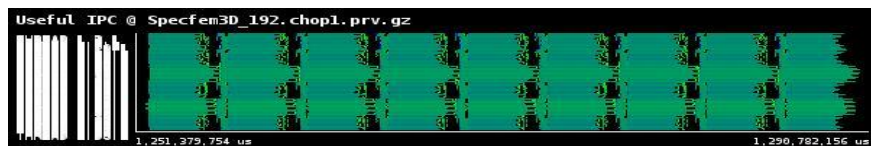
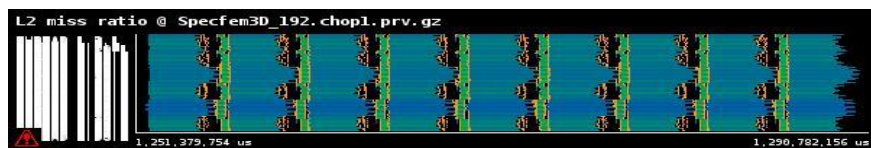
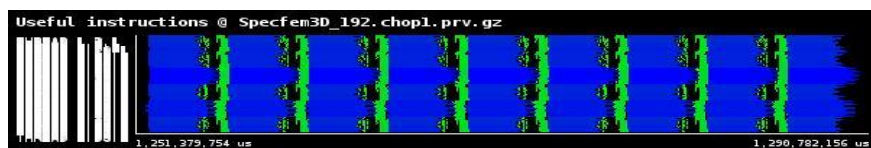
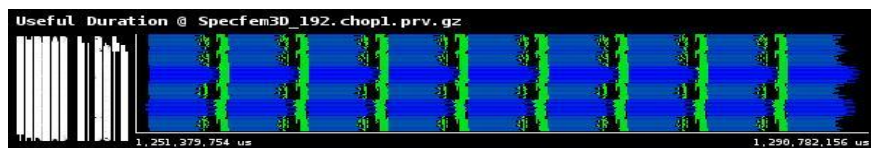




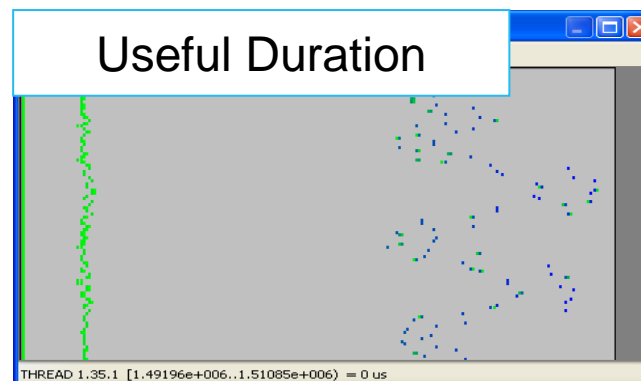
# Analyzing variability through histograms and timelines



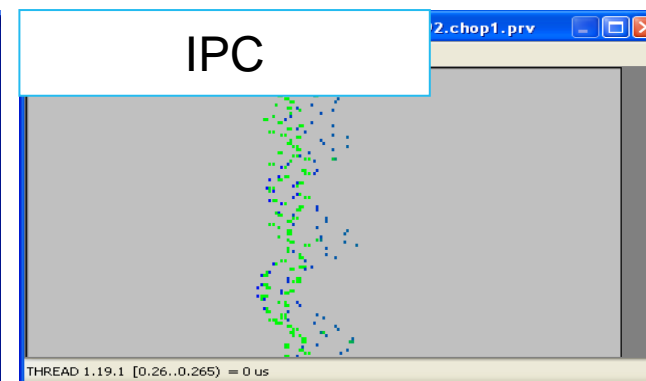
SPECFEM3D



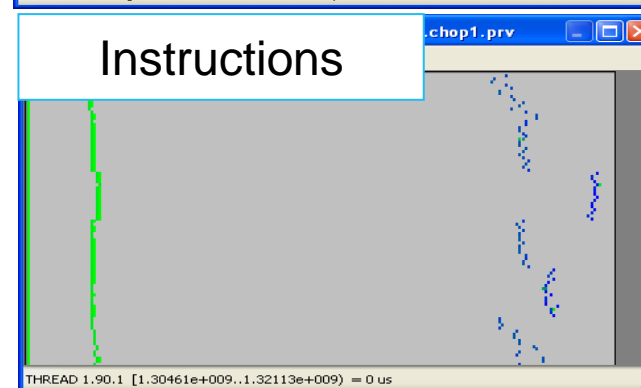
Useful Duration



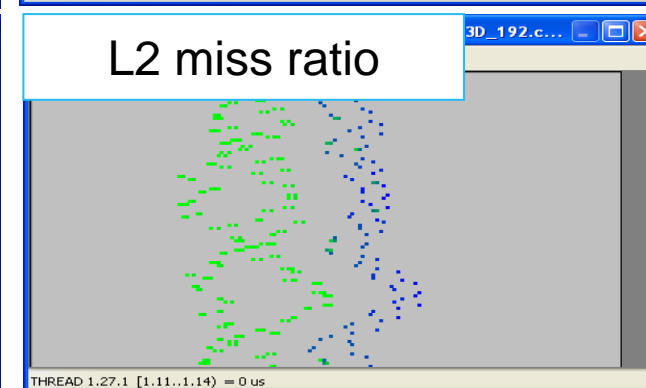
IPC



Instructions



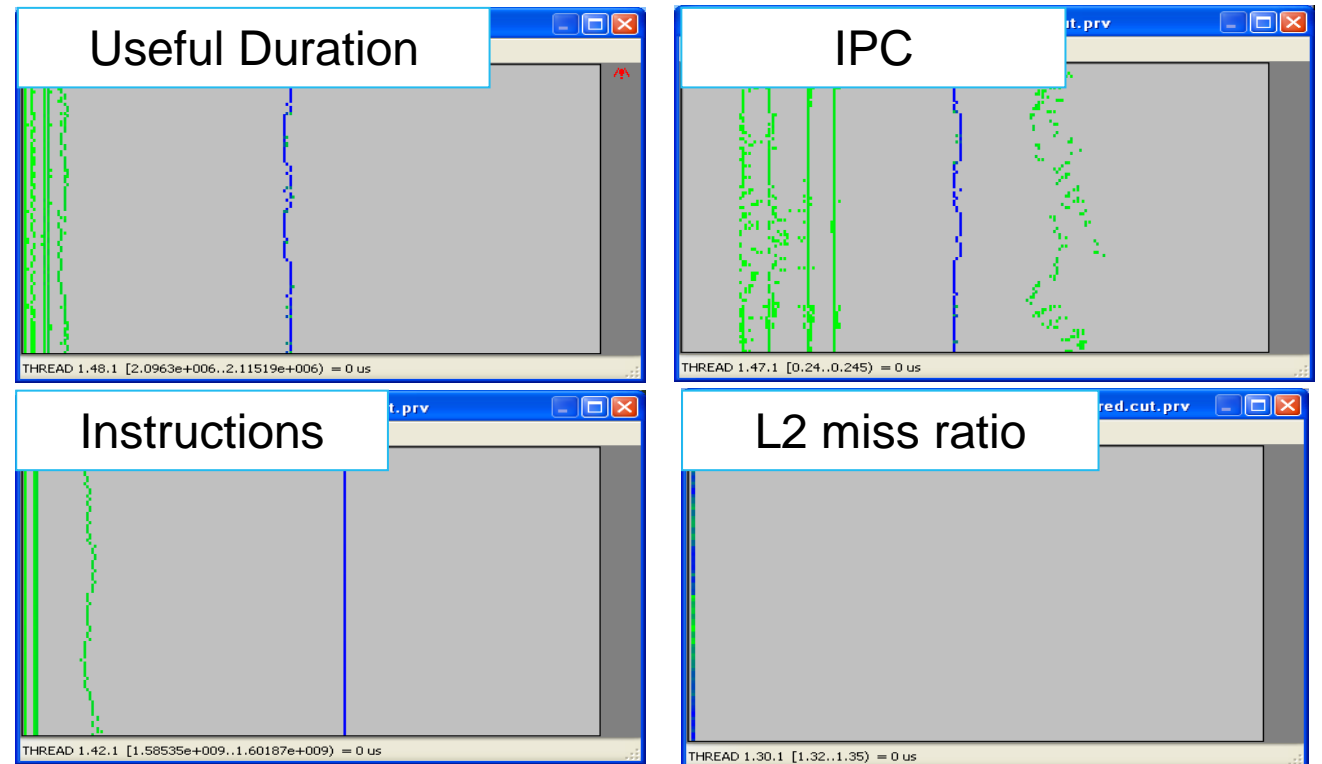
L2 miss ratio





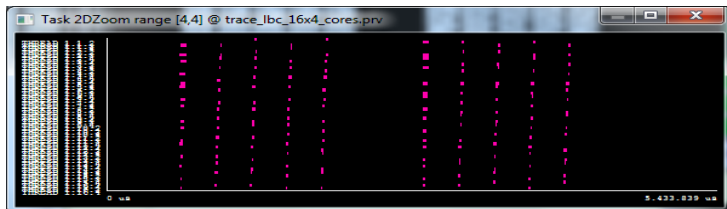
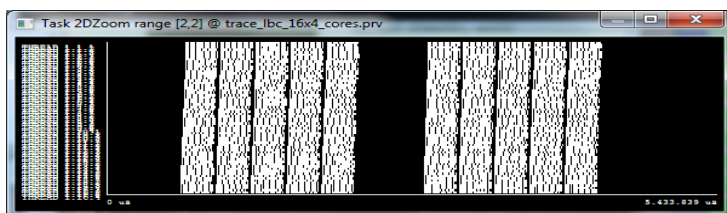
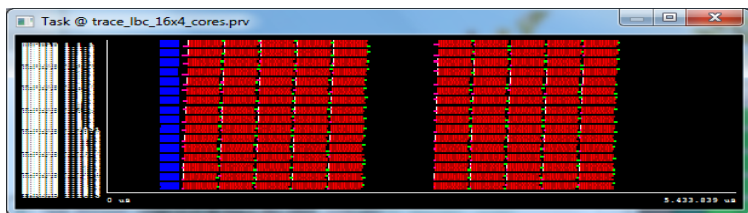
## Analyzing variability through histograms and timelines

- By the way: six months later ....

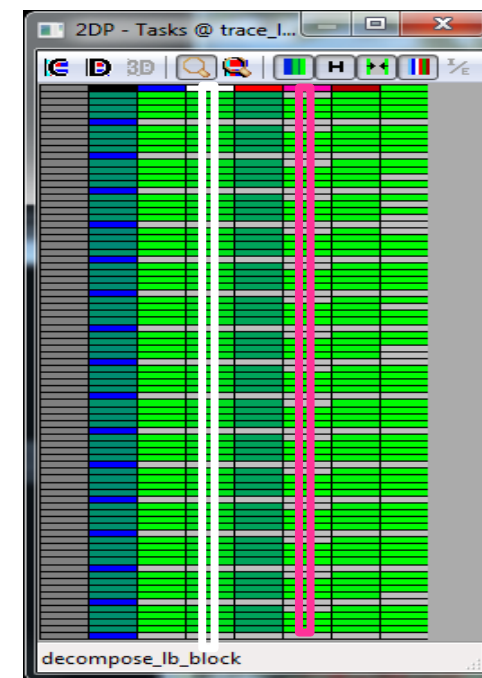


## From tables to timelines

- Where in the timeline do the values in certain table columns appear?  
ie. want to see the time distribution of a given routine?

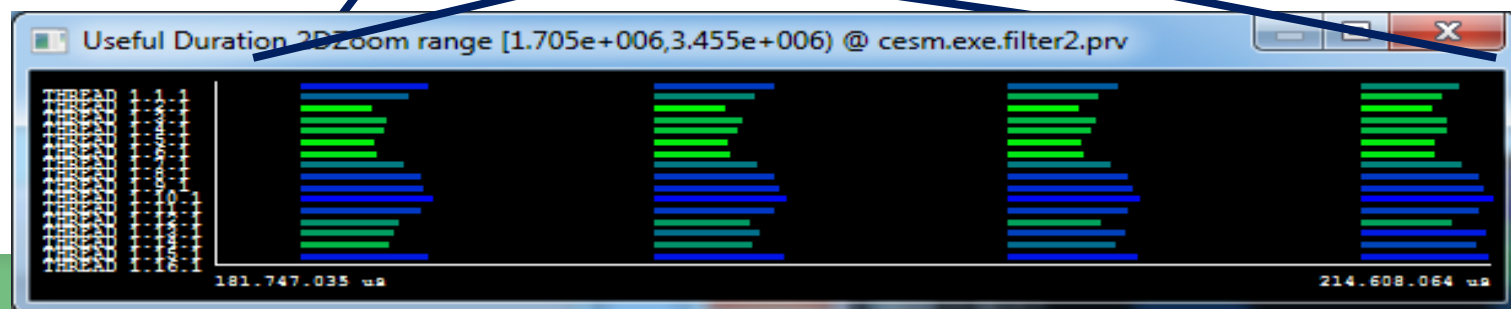
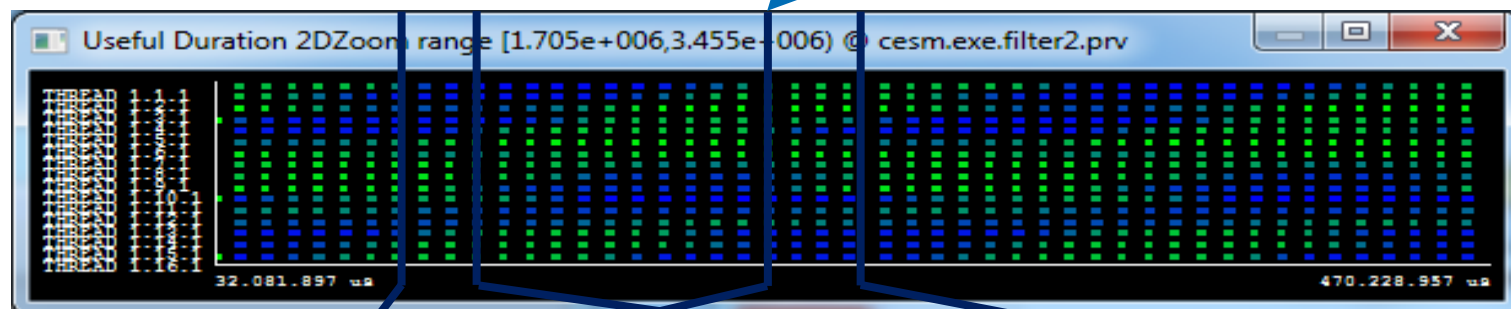
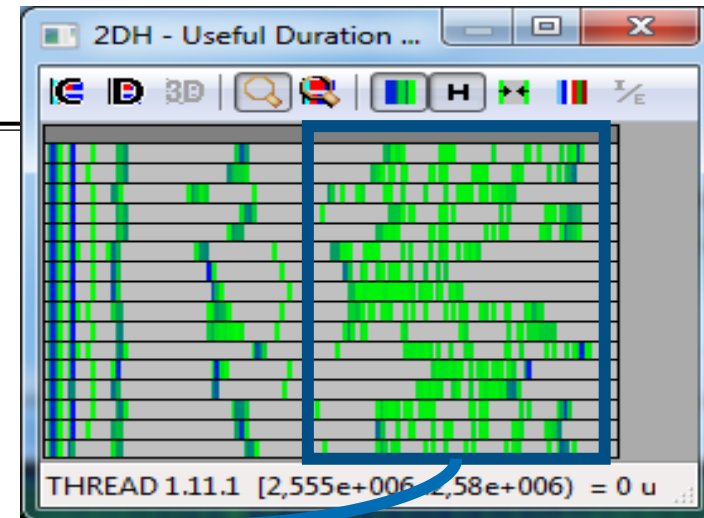


Only showing when a given value happens



## Variability ... is everywhere

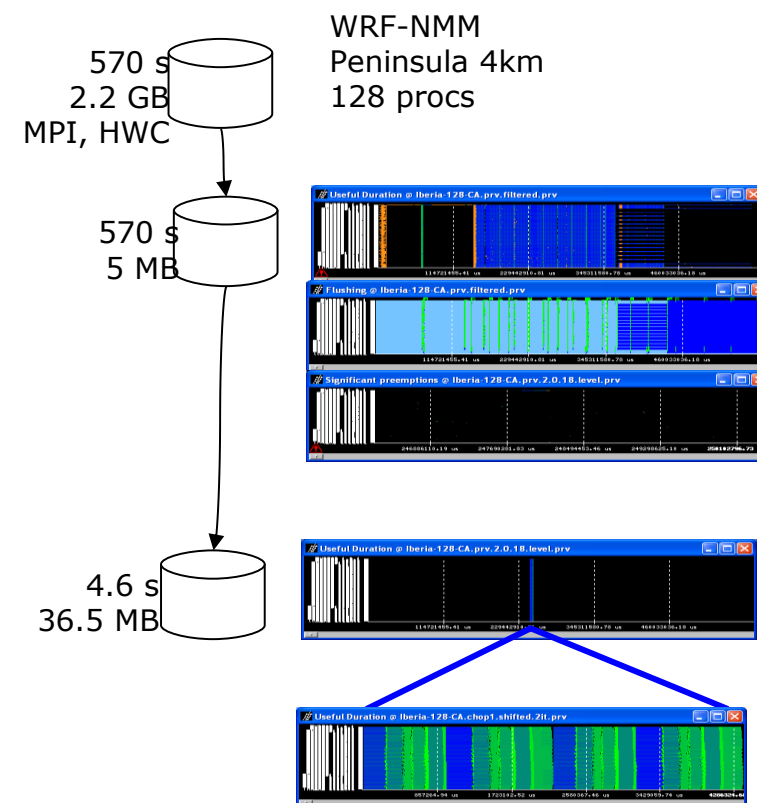
- CESM: 16 processes, 2 simulated days
- Histogram useful computation duration shows high variability
- How is it distributed?
- Dynamic imbalance
  - In space and time
  - Day and night.
  - Season ? ☺





## Trace manipulation

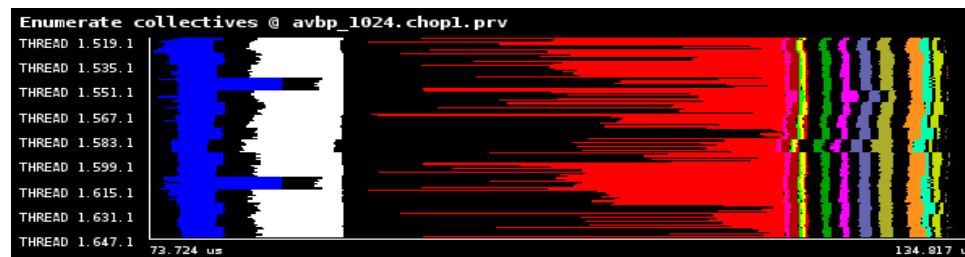
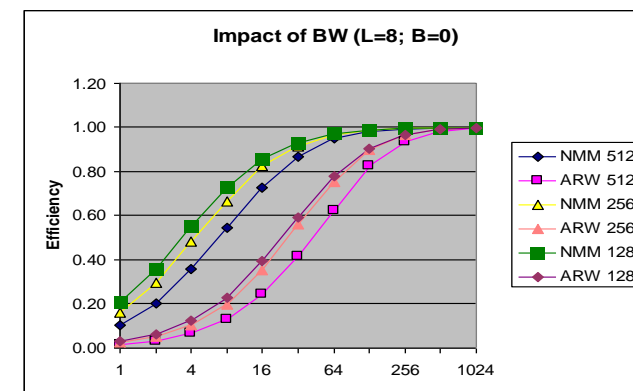
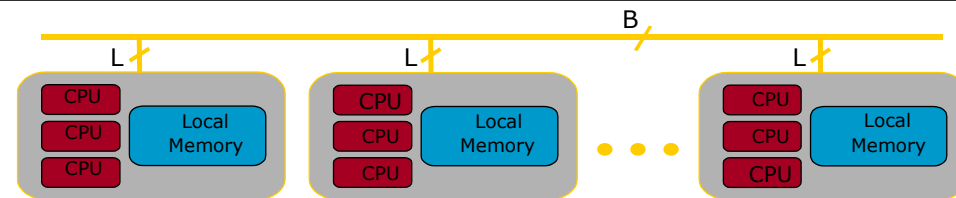
- Data handling/summarization capability
- Filtering
  - Subset of records in original trace
  - By duration, type, value,...
  - Filtered trace IS a paraver trace and can be analysed with the same cfgs (as long as needed data kept)
- Cutting
  - All records in a given time interval
  - Only some processes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



# Dimemas

# Dimemas: Coarse grain, Trace driven simulation

- Simulation: Highly non linear model
  - MPI protocols, resources 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 of an app. equally sensitive to network?
- MPI sanity check
  - Modeling nominal
- Paraver – Dimemas tandem
  - Analysis and prediction
  - What-if from selected time window



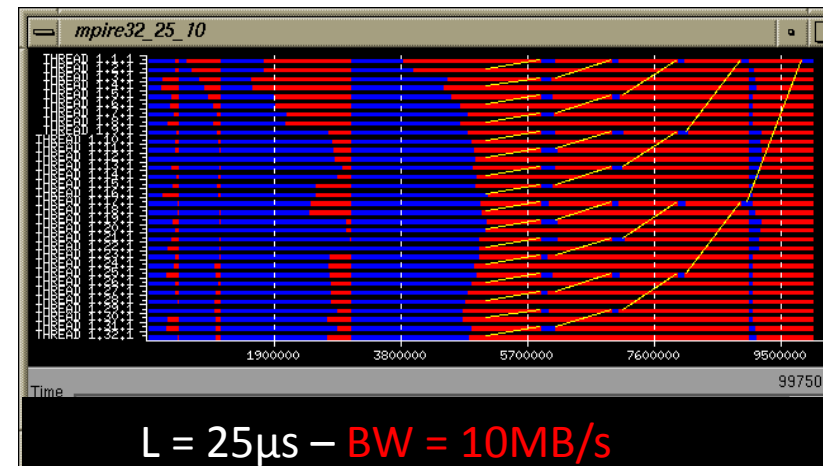
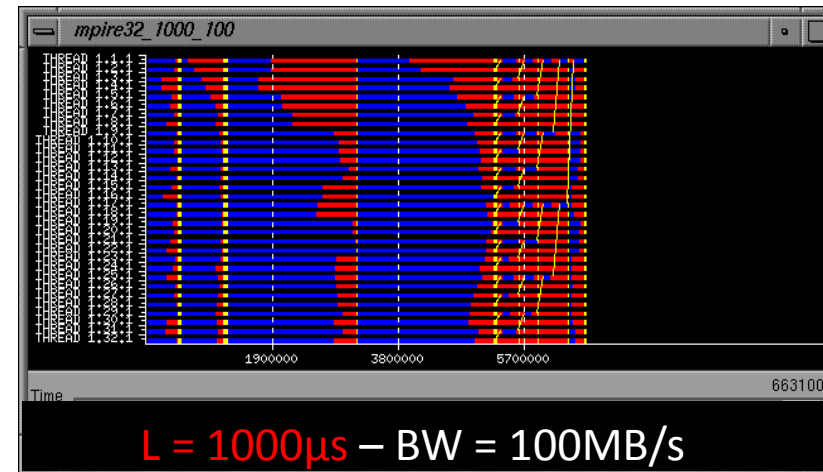
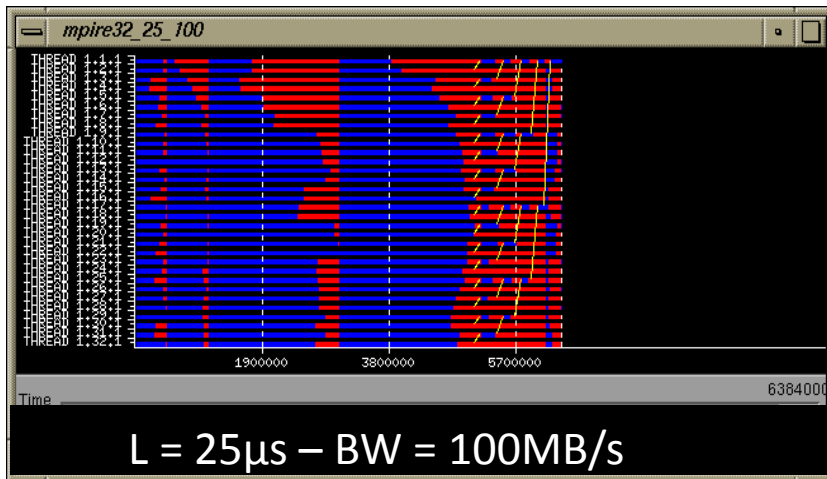
4ms

Detailed feedback on simulation (trace)



## Network sensitivity

- MPIRE 32 tasks, no network contention



All windows  
same scale

## Ideal machine

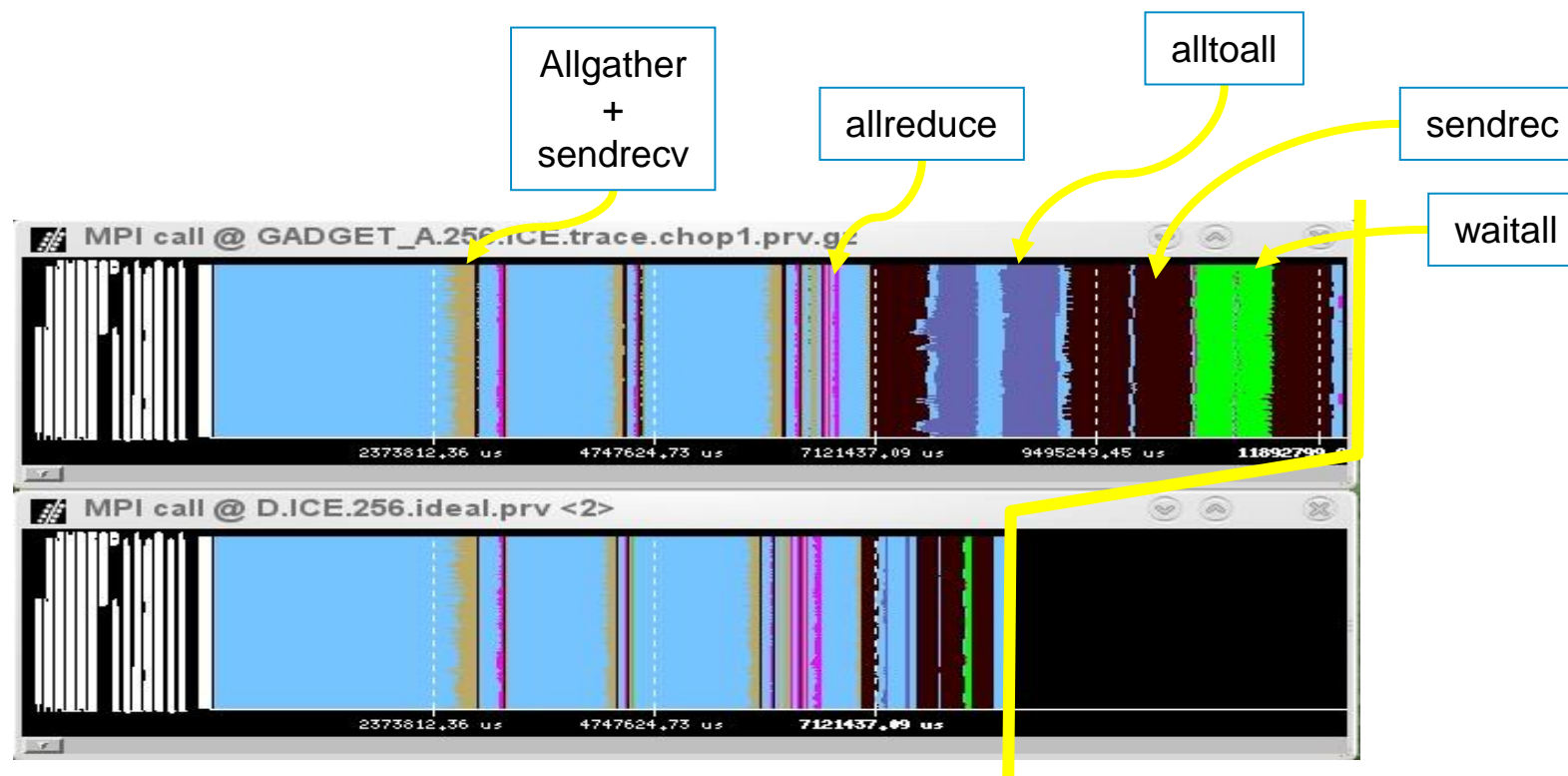
The impossible machine:  $BW = \infty$ ,  $L = 0$

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

GADGET @ Nehalem cluster  
256 processes

Real  
run

Ideal  
network



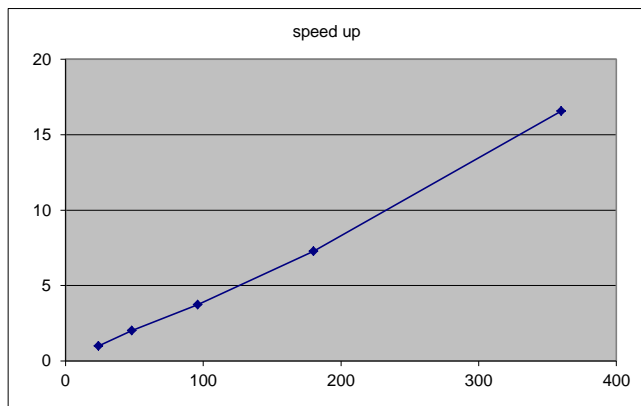
# Models



# Why scaling?

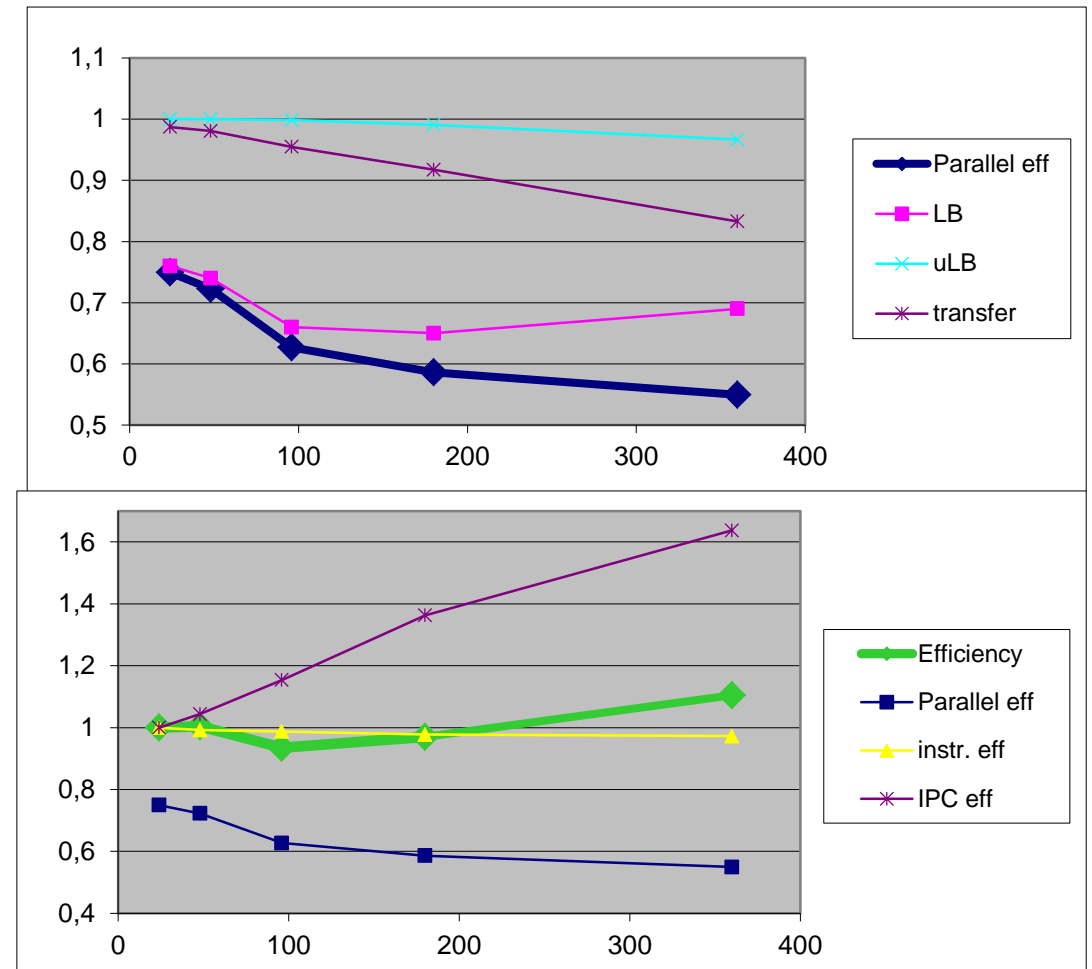
CG-POP mpi2s1D - 180x120

Good scalability !!  
Should we be happy?



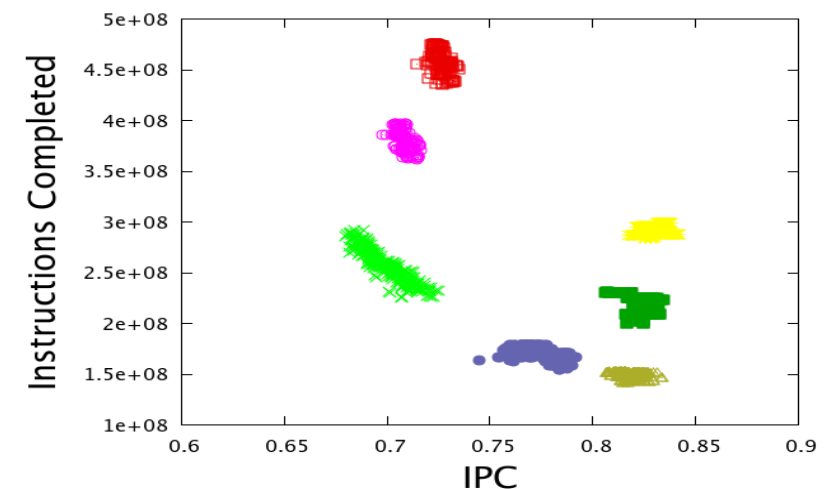
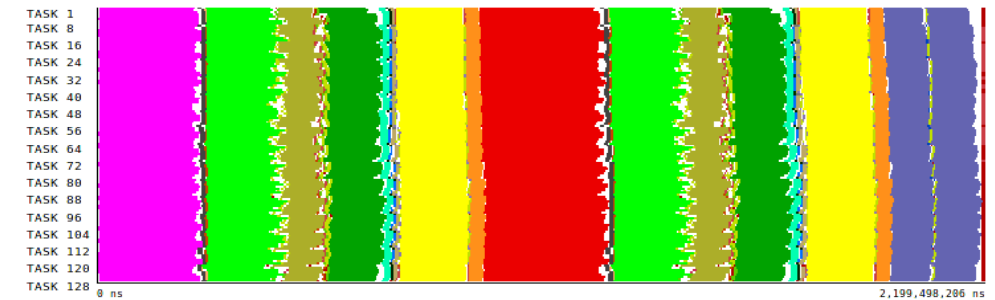
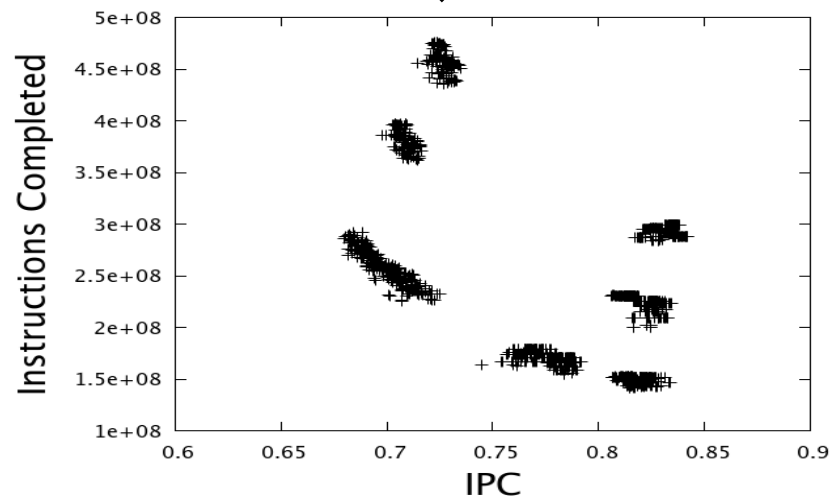
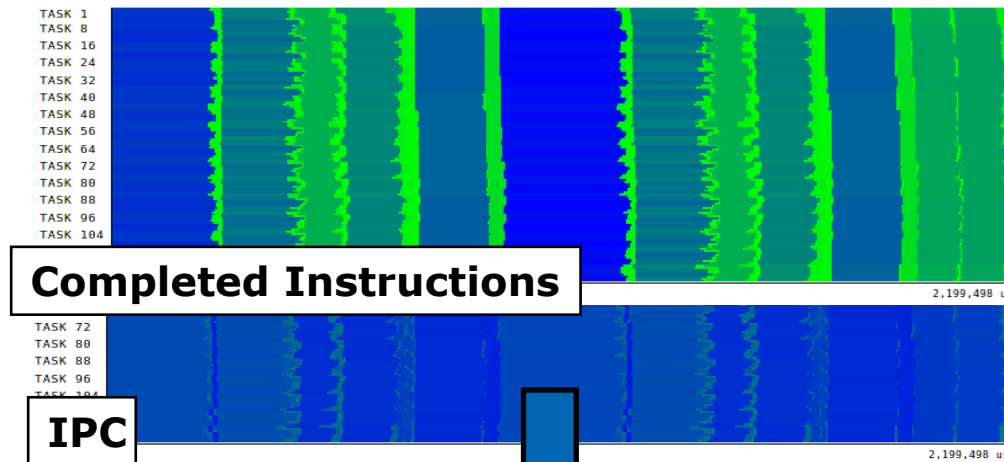
$$\eta_{\parallel} = LB * Ser * Trf$$

$$\eta = \eta_{\parallel} * \eta_{instr} * \eta_{IPC}$$



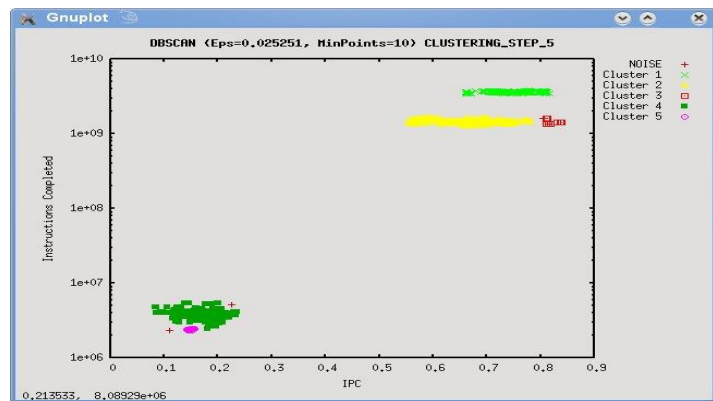
# Clustering

## Using Clustering to identify structure

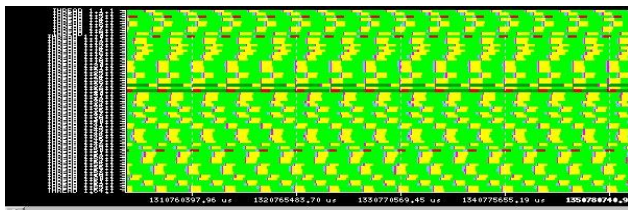




# Performance @ serial computation bursts

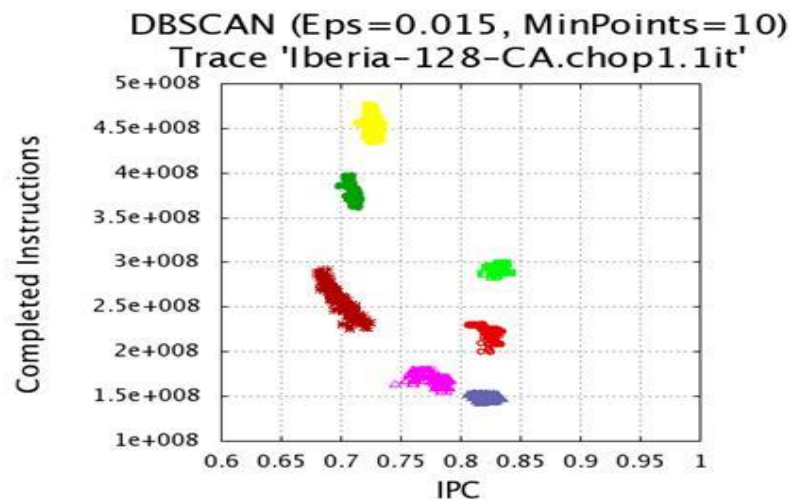


SPECfem3D

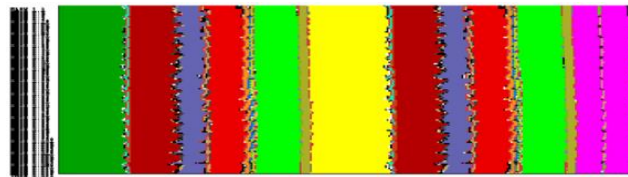


Asynchronous SPMD

Balanced #instr variability  
in IPC



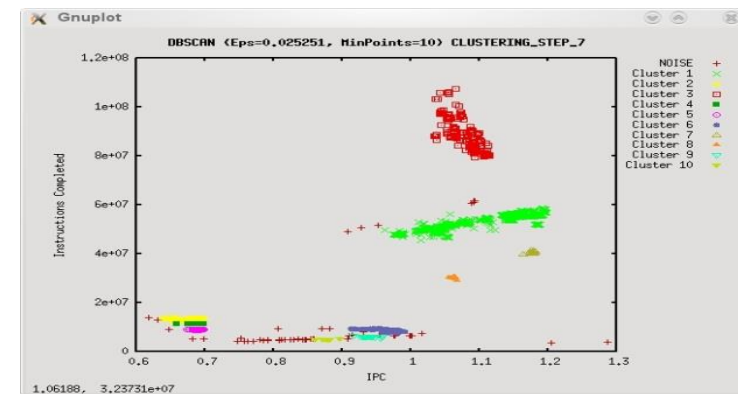
WRF 128 cores



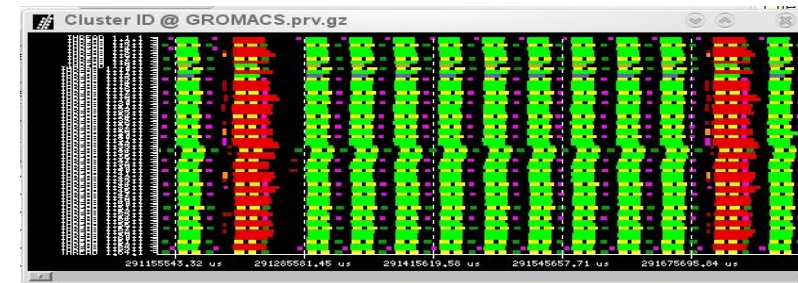
SPMD

Repeated substructure

Coupled imbalance



GROMACS

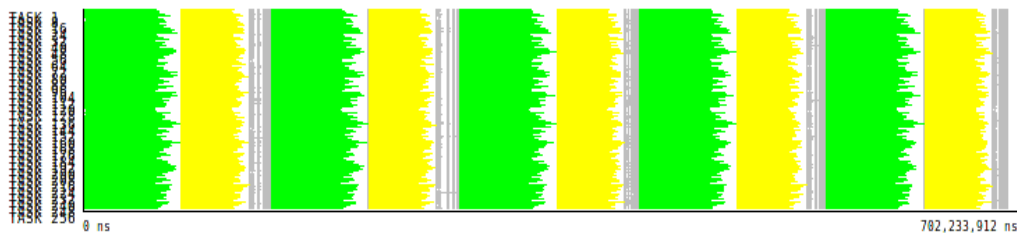


MPMD structure

Different coupled  
imbalance trends

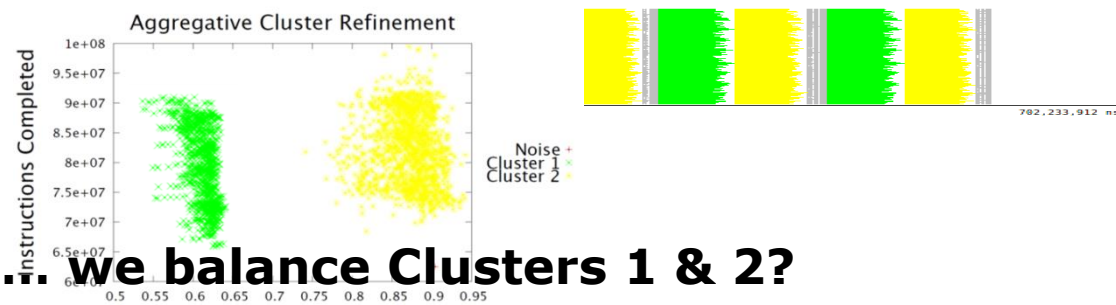
## Integrating models and analytics

What if ....



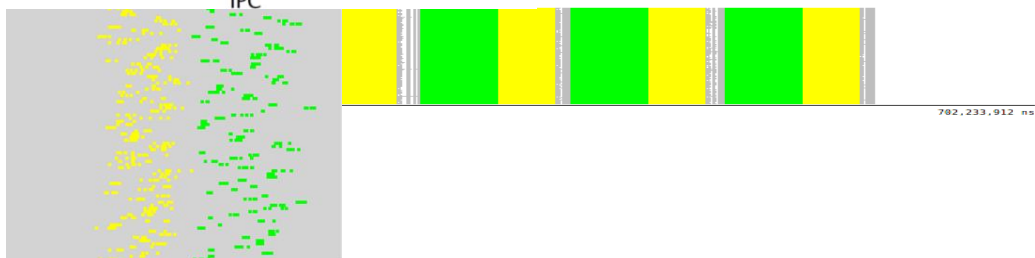
PEPC

... we increase the IPC of Cluster1?



13% gain

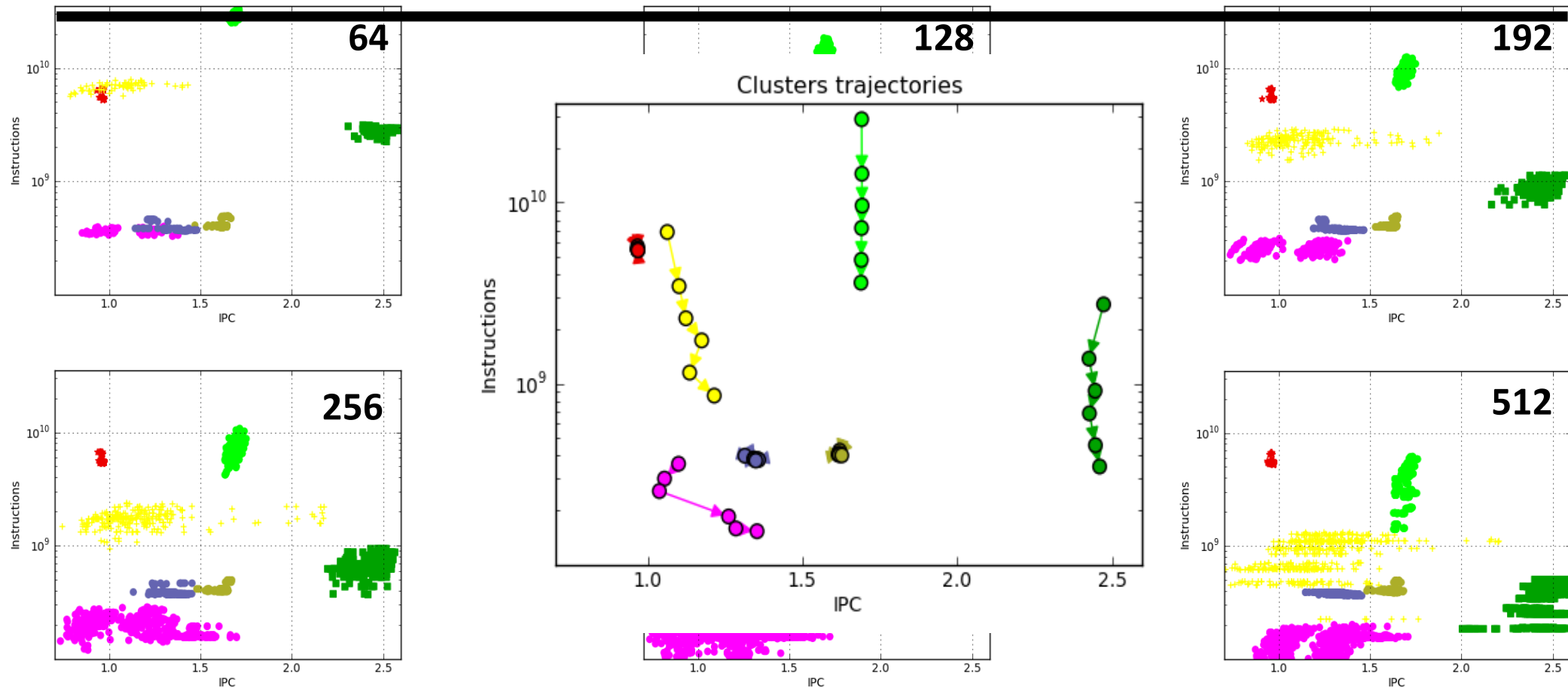
... we balance Clusters 1 & 2?



19% gain

## Tracking: scalability through clustering

OpenMX (strong scale from 64 to 512 tasks)





# Folding



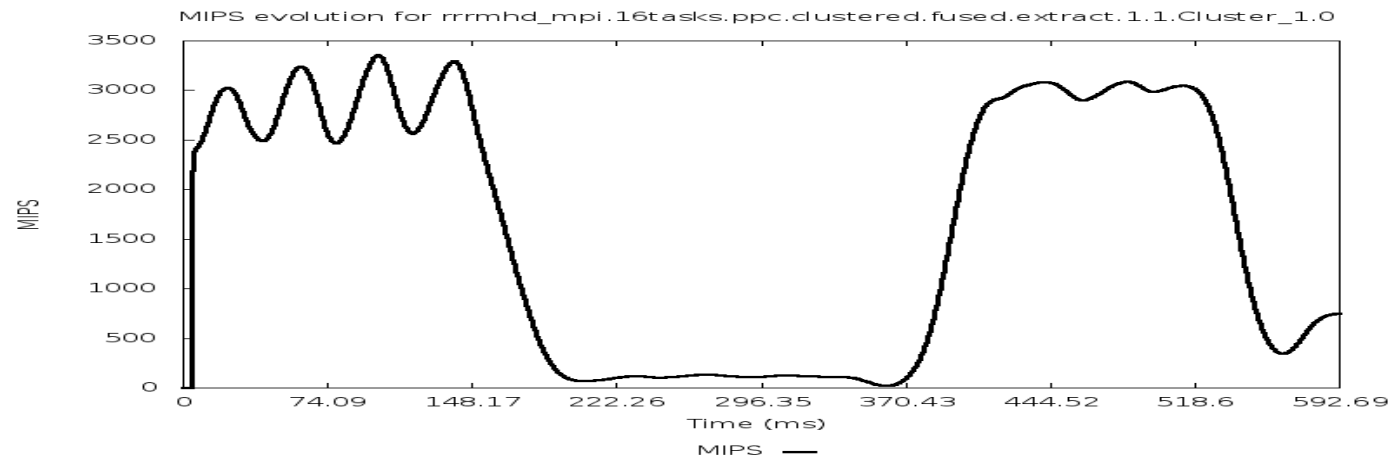
## Folding: Detailed metrics evolution

---

Performance of a sequential region = 2000 MIPS

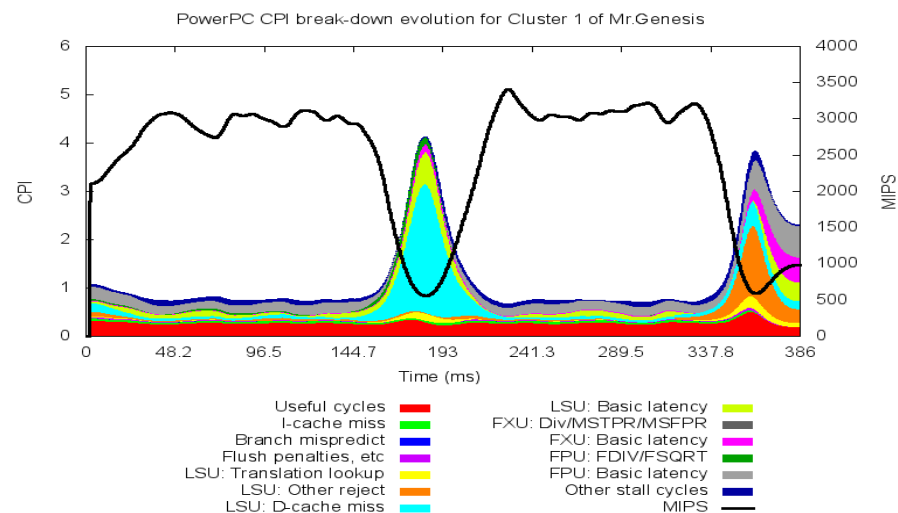
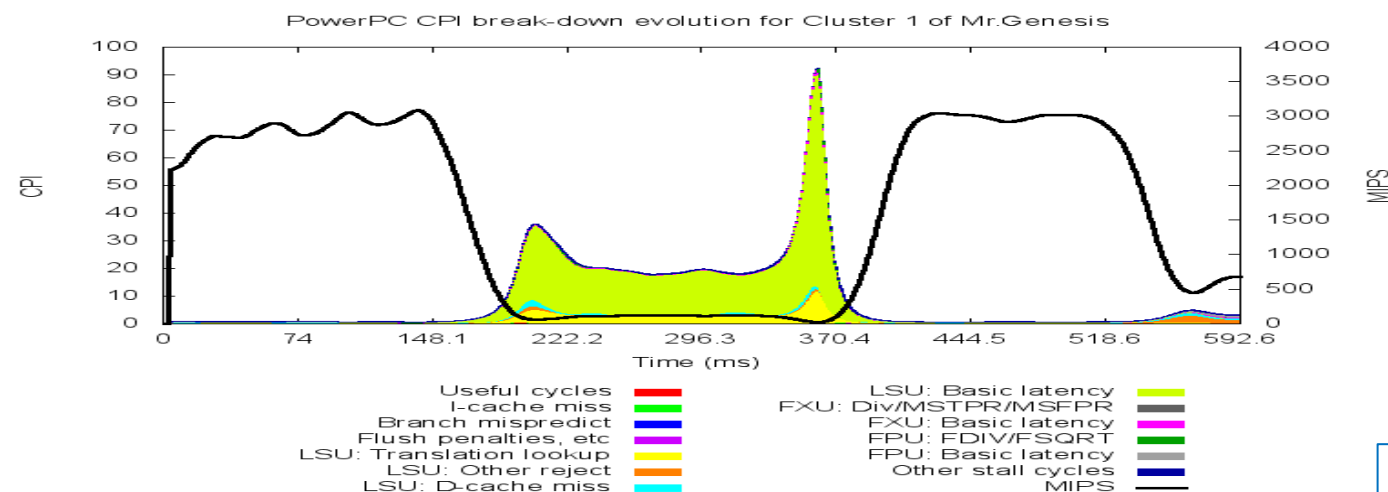
Is it good enough?

Is it easy to improve?



## Folding: Instantaneous CPI stack

### MRGENESIS

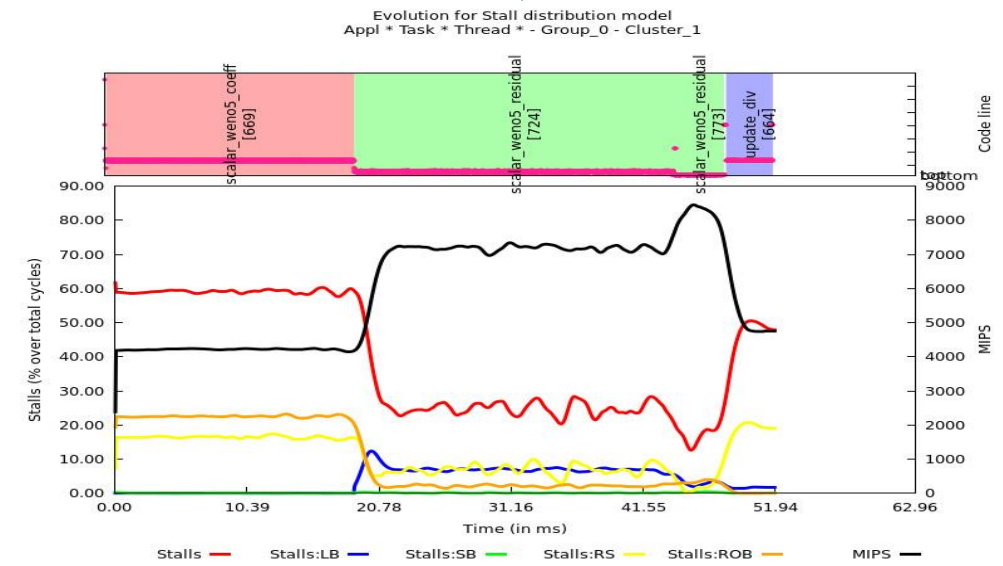
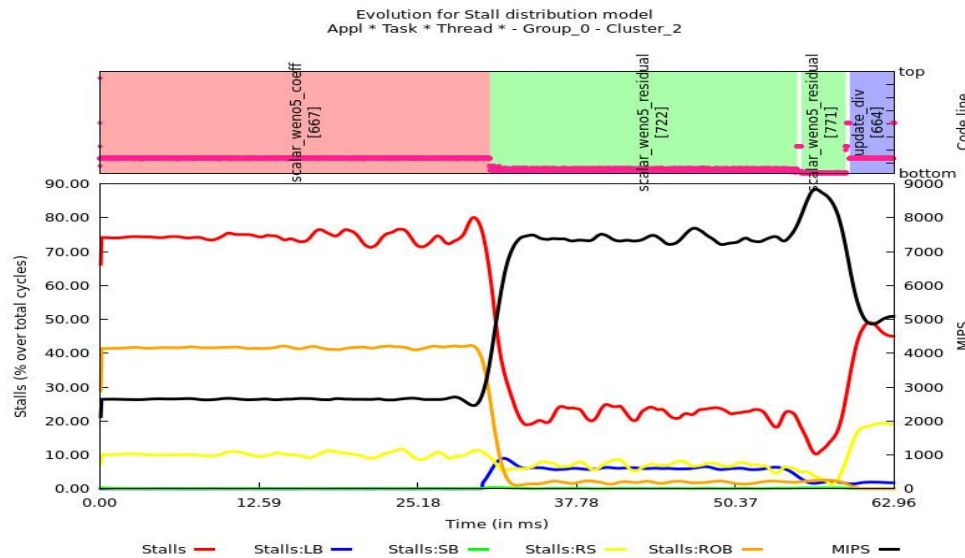


- Trivial fix.(loop interchange)
- Easy to locate?
- Next step?
- Availability of CPI stack models for production processors?
  - Provided by manufacturers?

## “Blind” optimization

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

Recommendation without access to source code



# Methodology



## 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

---

- [www.bsc.es/paraver](http://www.bsc.es/paraver)
- 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