

# Understanding applications using the BSC performance tools

---

Lau Mercadal, Germán Llort

✉ [tools@bsc.es](mailto:tools@bsc.es)

Barcelona Supercomputing Center

---

# Humans are visual creatures

---

- Painting / photo or description?
  - Our brain processes visual impressions 60,000 times faster than text
  - It takes only 13 ms for the human brain to process an image
- Memorizing a deck of playing cards
  - Each card translated to an image (person, action, location)
- Our brain loves pattern recognition
  - 90% of all communication happens visually
  - The human eye can differentiate approximately 10 million different colours

PROCESS

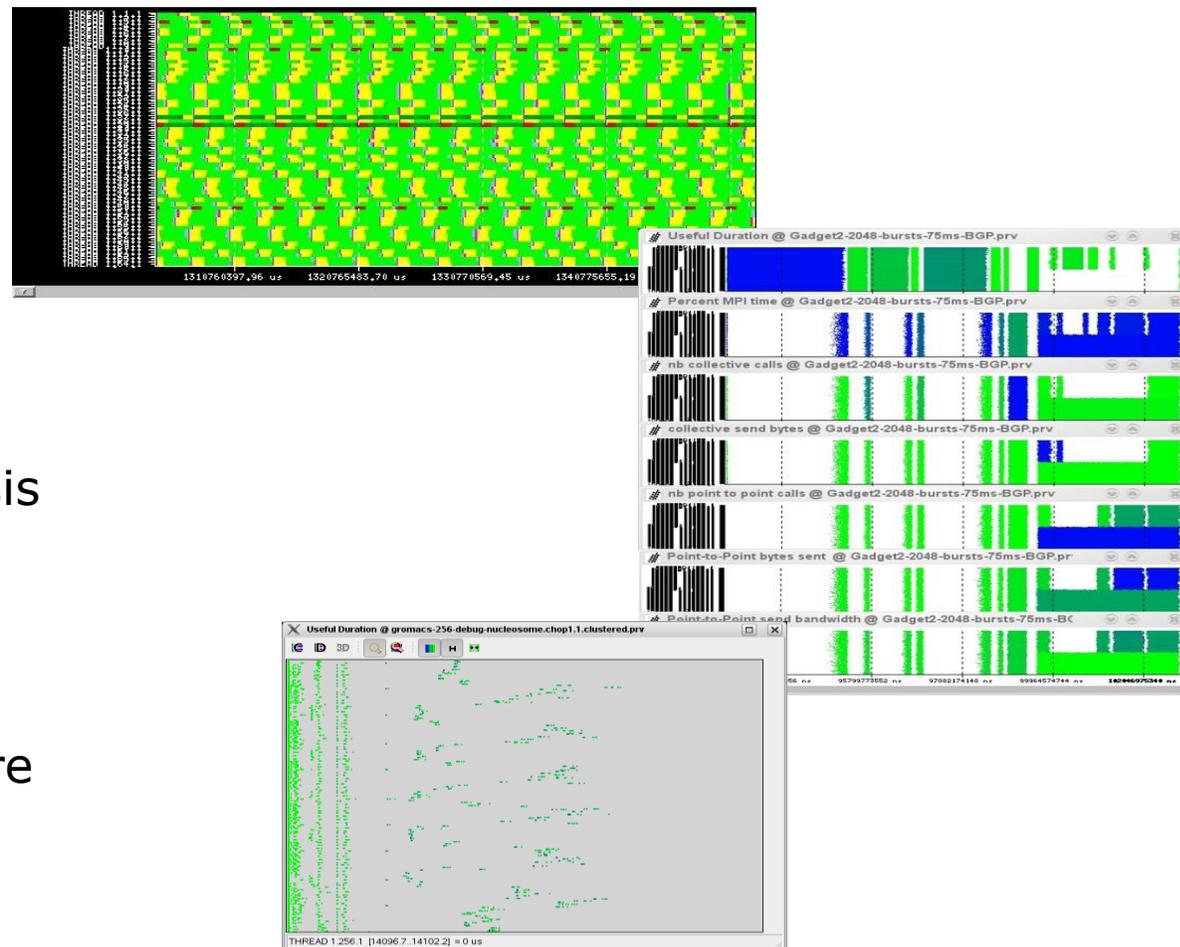
STORE

IDENTIFY



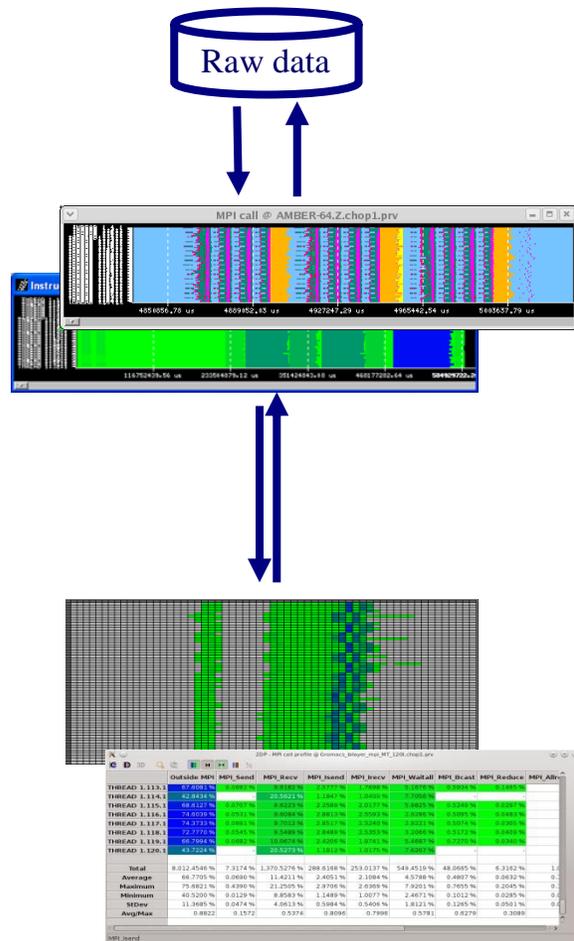
# Our Tools

- Since 1991
- Based on traces
- Open Source
  - <https://tools.bsc.es>
- Core tools:
  - **Extrae** – instrumentation
  - **Paraver** (paramedir) – offline trace analysis
  - **Dimemas** – message passing simulator
- Focus
  - Detail, variability, flexibility
  - Behavioural 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

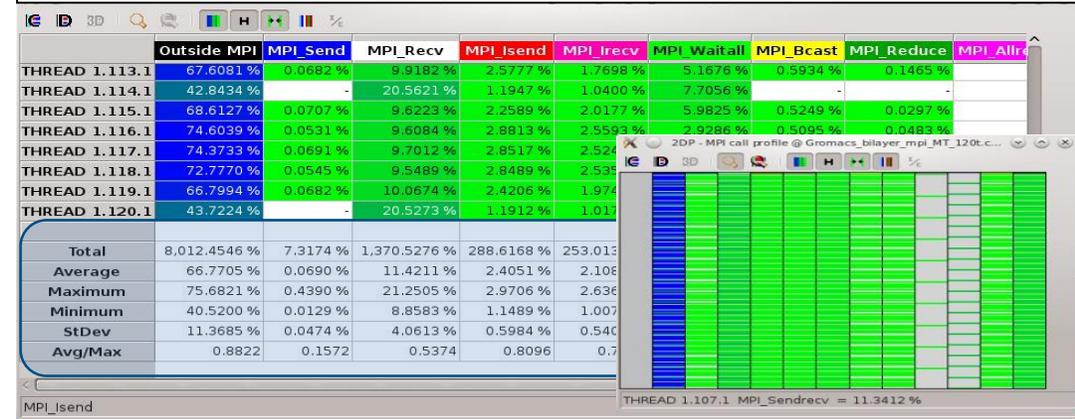
Synchronize scales

# From timelines to tables

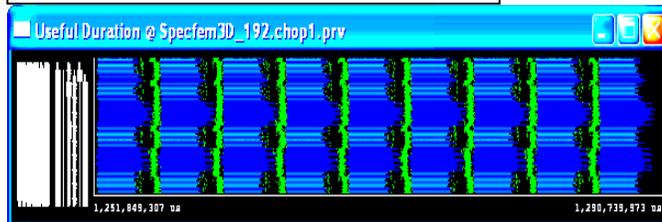
## MPI calls



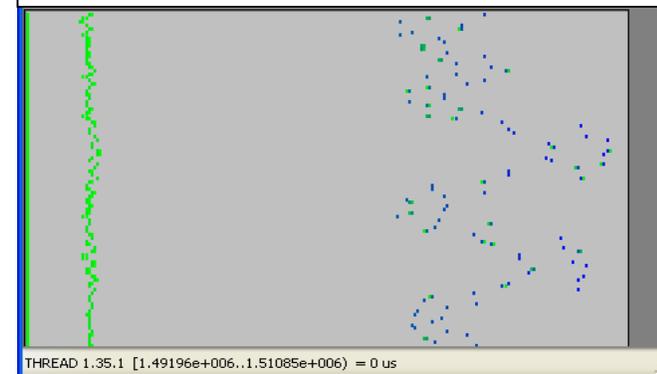
## MPI calls profile



## Useful Duration



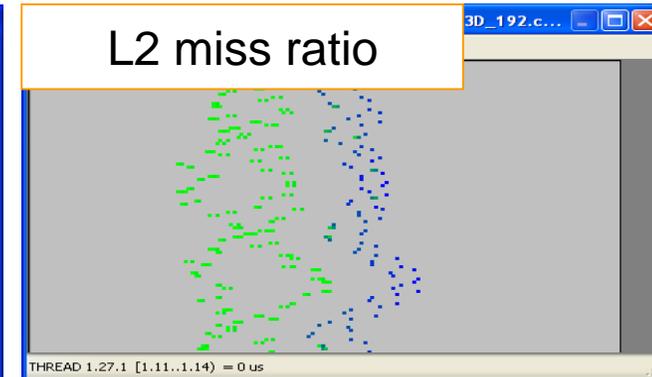
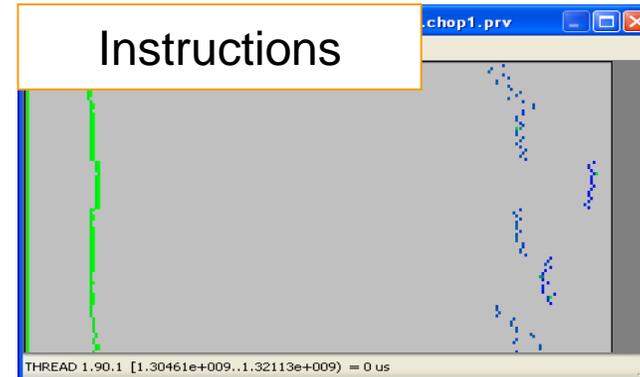
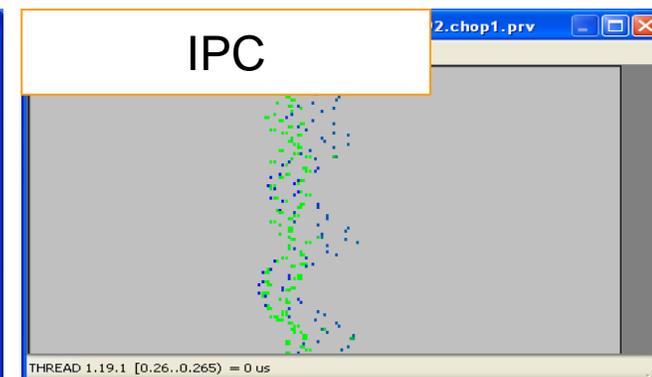
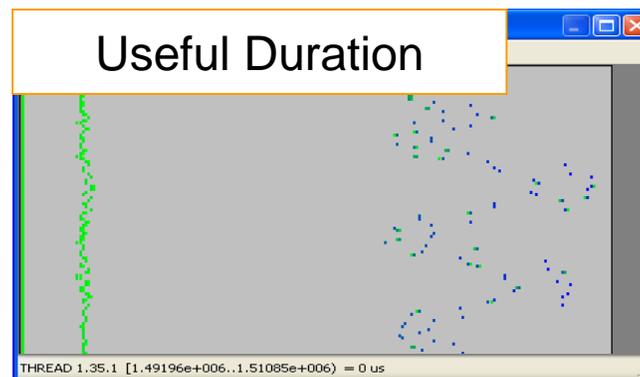
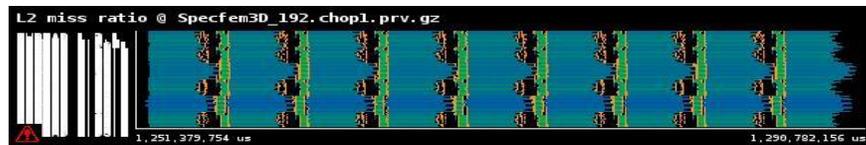
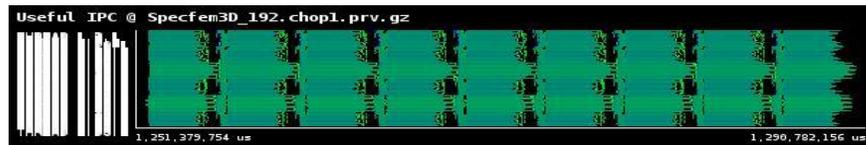
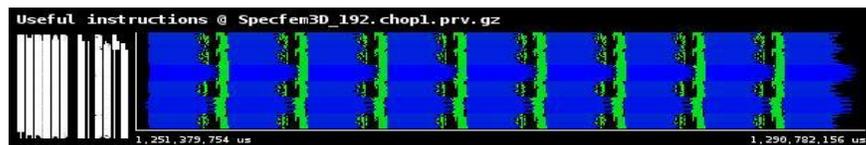
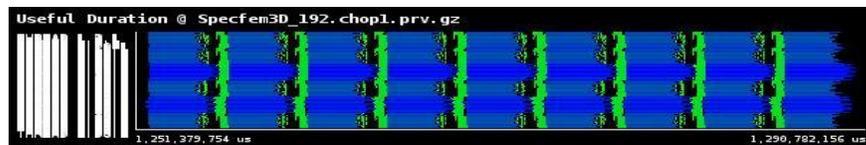
## Histogram Useful Duration



# Analyzing variability through histograms and timelines

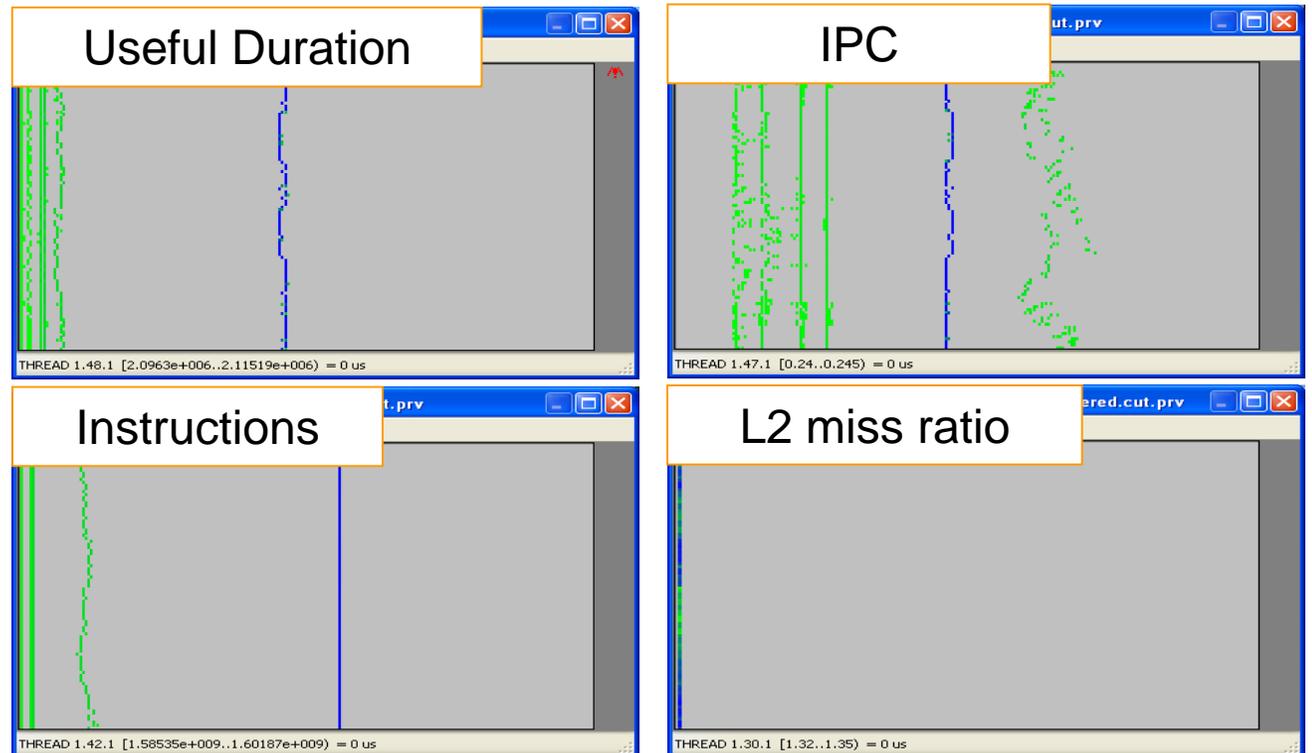


SPECFEM3D



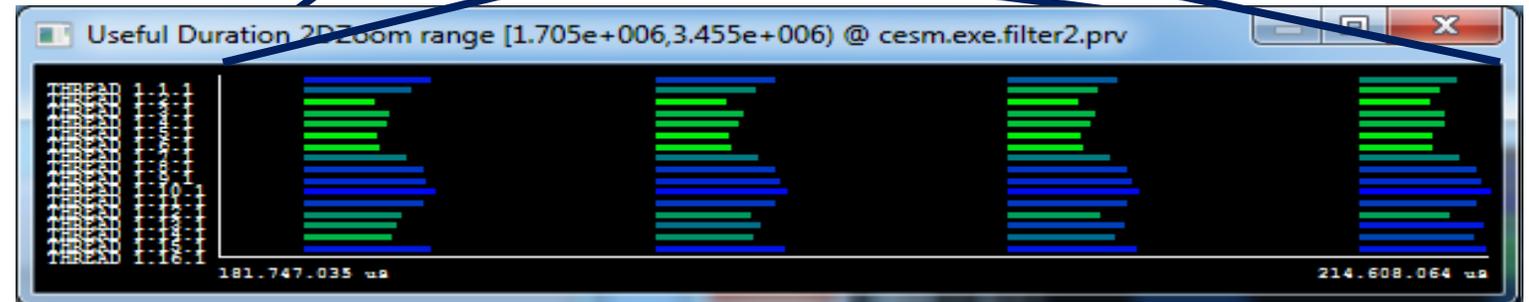
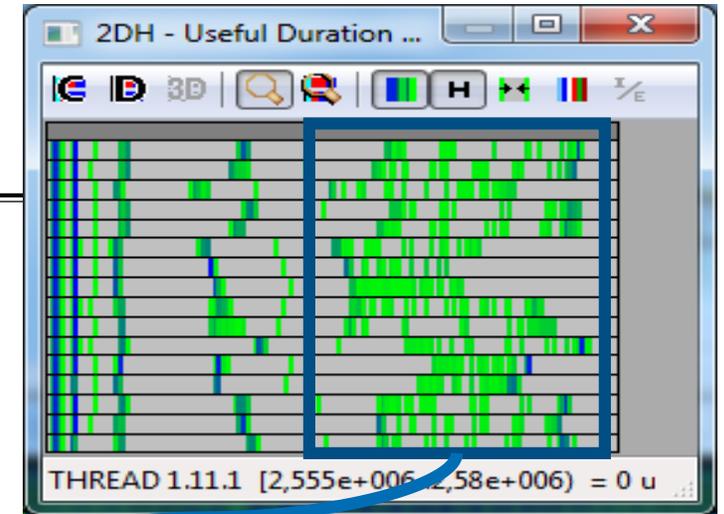
# Analyzing variability through histograms and timelines

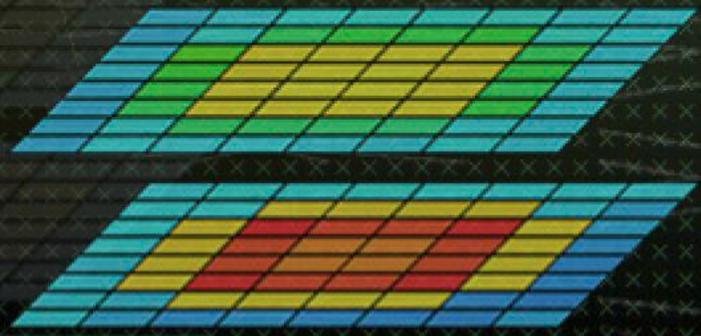
- By the way: six months later ...



# 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 ? ☺

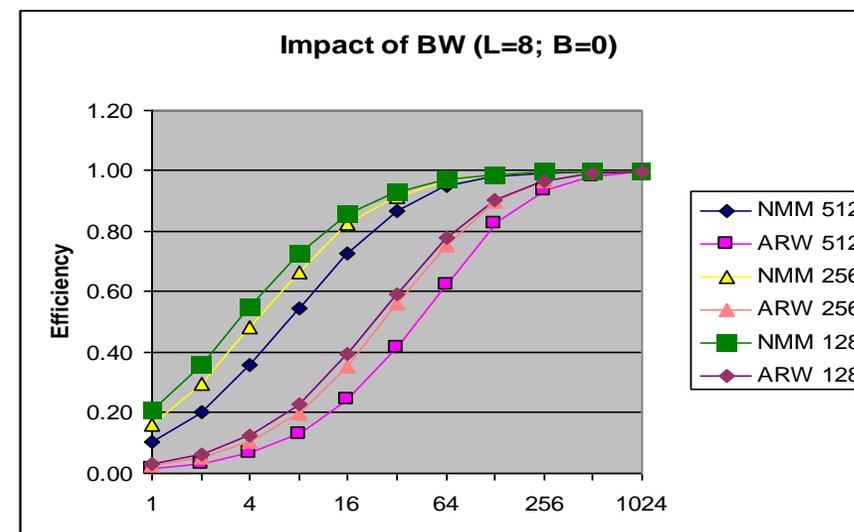
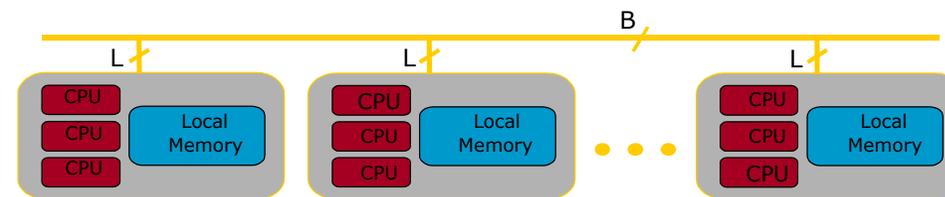




# Dimemas

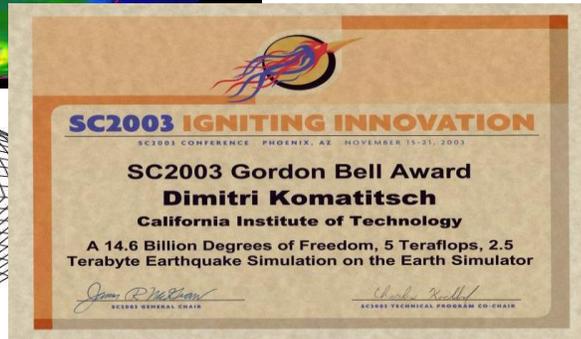
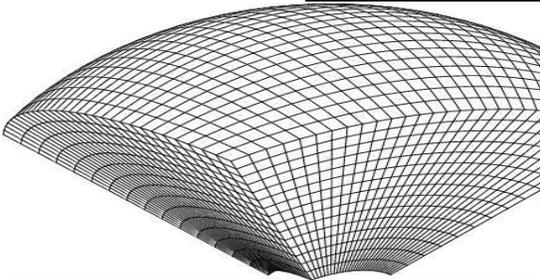
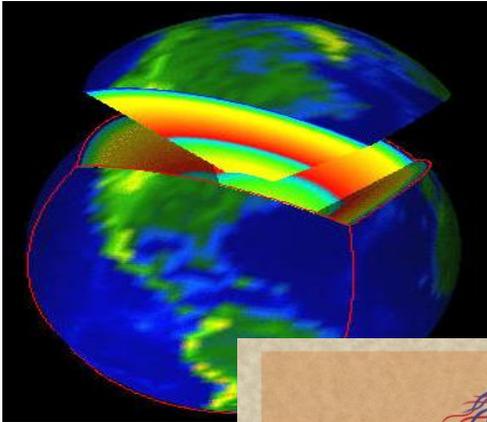
# Dimemas: Coarse grain, Trace driven simulation

- 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
  - Modelling nominal
- Paraver – Dimemas tandem
  - Analysis and prediction
  - What-if from selected time window

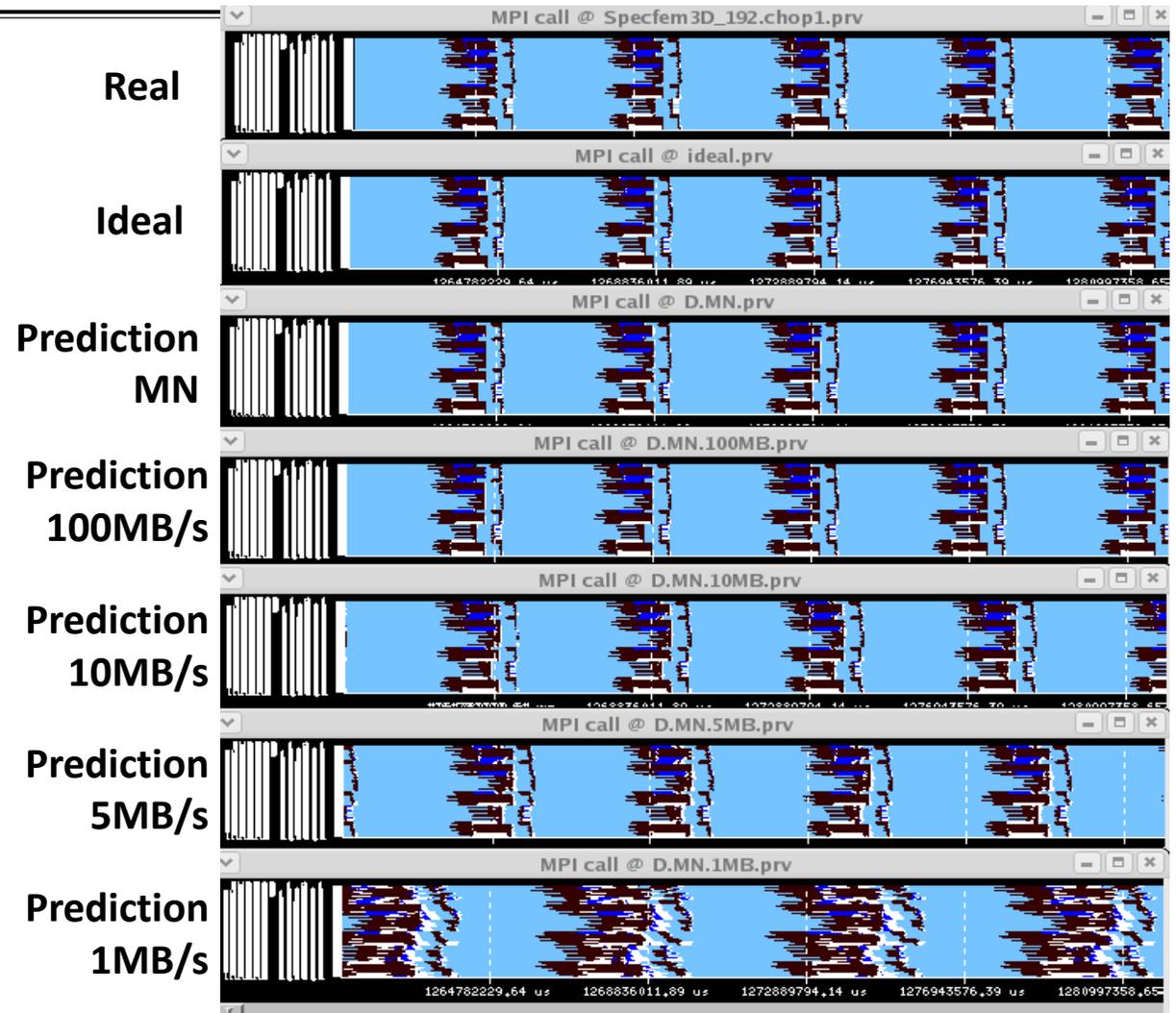


# What if we had asynchronous communications?

- SPECFEM3D

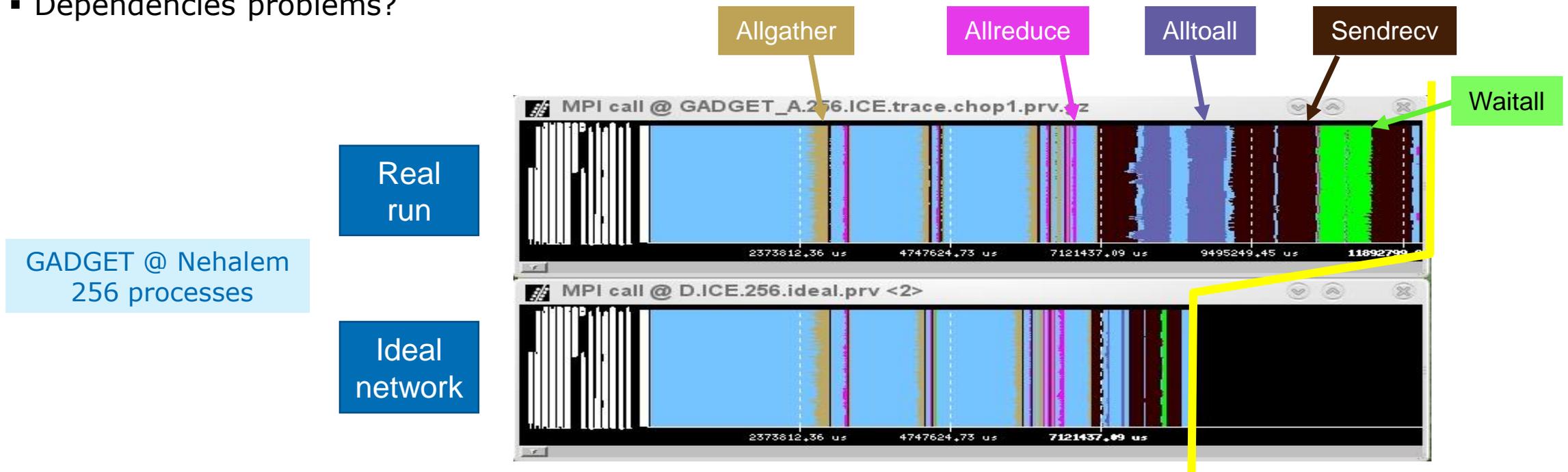


Courtesy Dimitri Komatitsch



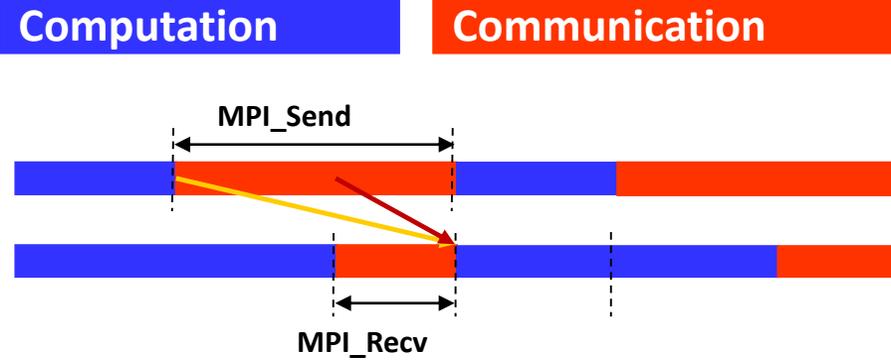
# Ideal machine

- The impossible machine:  $BW = \infty, L = 0$ 
  - Actually describes/characterizes Intrinsic application behavior
    - Load balance problems?
    - Dependencies problems?

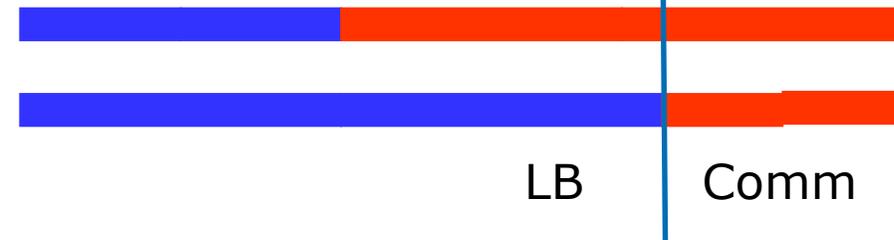


# Models

# Parallel efficiency model



Do not blame MPI



- Parallel efficiency = LB eff \* Comm eff

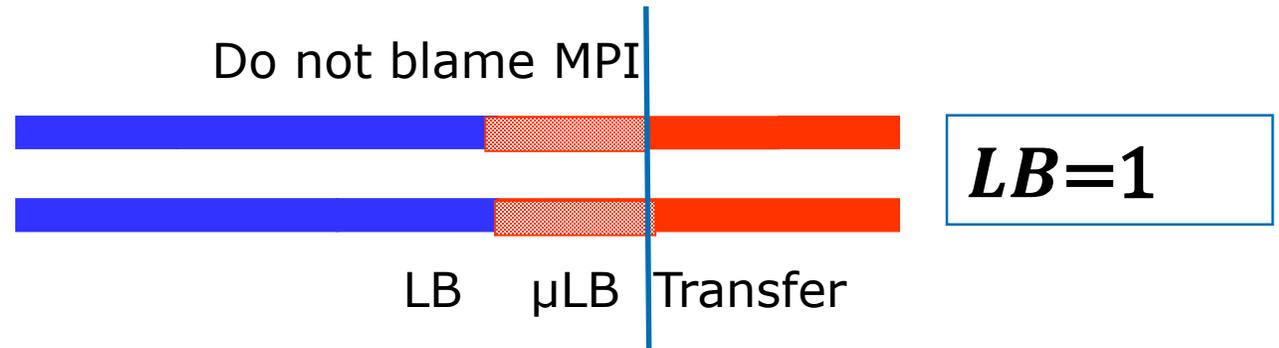
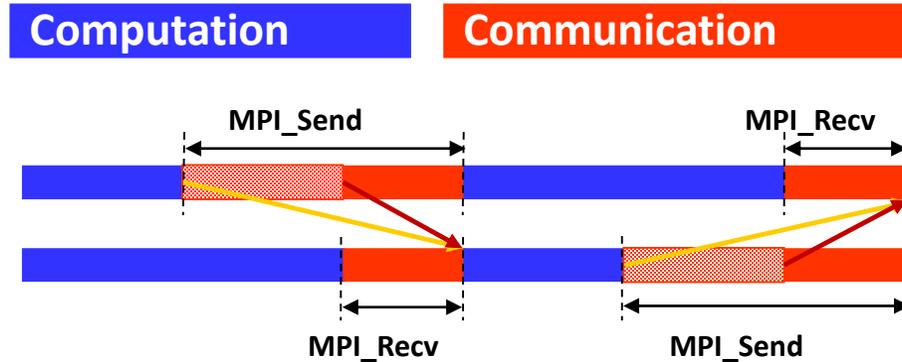
2DP - MPI call profile @ trace\_24h\_atmos\_symbols.cho...

	Outside MPI	MPI_Recv	MPI_Isend	MPI_Irecv
THREAD 1.130.1	87,93 %	2,31 %	0,01 %	0,02 %
THREAD 1.131.1	88,16 %	9,09 %	0,00 %	0,02 %
THREAD 1.132.1	88,18 %	9,09 %	0,00 %	0,02 %
THREAD 1.133.1	88,18 %	9,09 %	0,00 %	0,02 %
<b>Total</b>	9,309,74 %	306,53 %	1.411,18 %	3,83 %
<b>Average</b>	69,00 %	2,30 %	10,69 %	0,03 %
<b>Maximum</b>	88,18 %	67,62 %	54,97 %	
<b>Minimum</b>	30,67 %	0,00 %	0,00 %	
<b>StDev</b>	15,27 %	6,06 %	21,40 %	0,00 %
<b>Avg/Max</b>	0,77	0,03	0,19	0,81

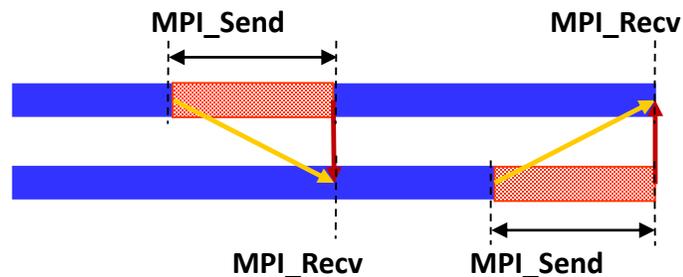
Annotations in the table:

- $\eta$  points to the 'Maximum' row.
- CommEff points to the 'Maximum' row.
- LB points to the 'Avg/Max' row.

# Parallel efficiency refinement: $LB * \mu LB * \text{Transfer}$



- Serializations / dependences ( $\mu LB$ )
- Dimemas ideal network  $\rightarrow$  Transfer (efficiency) = 1

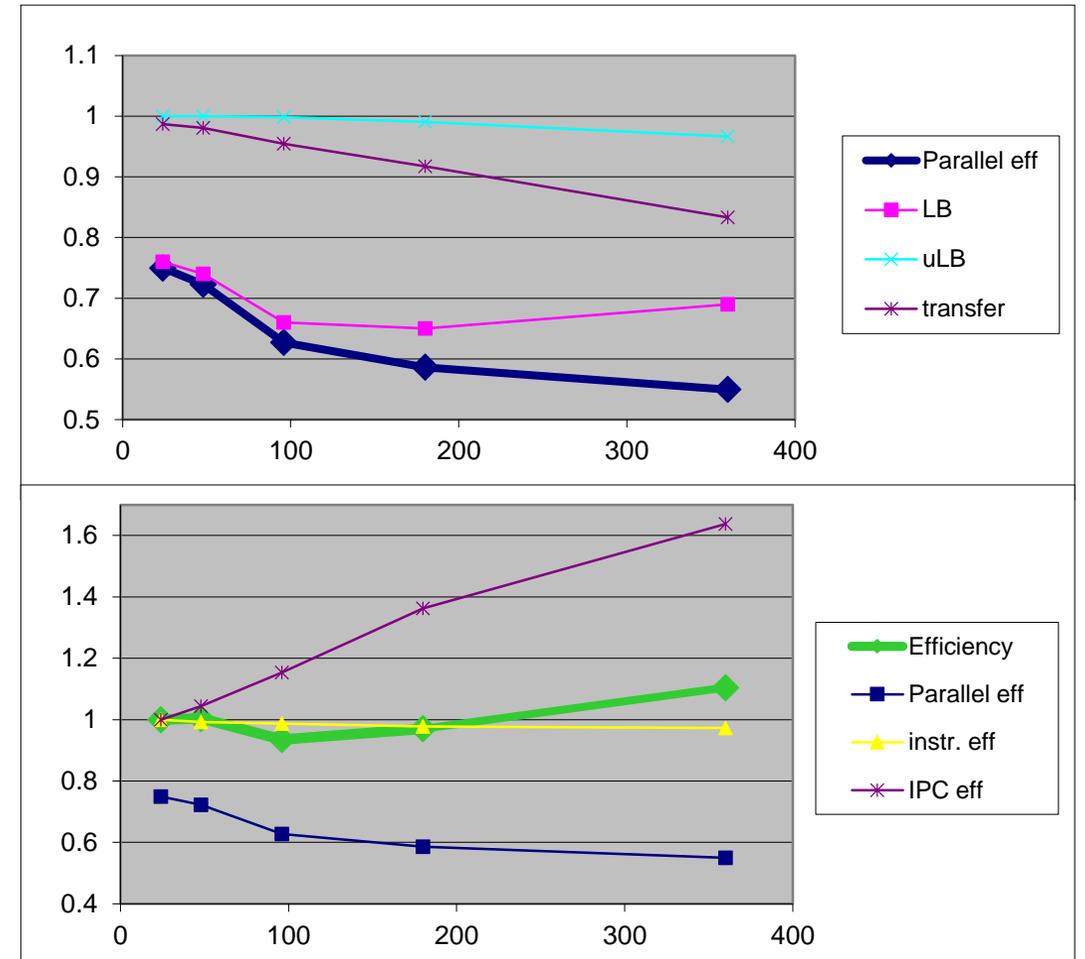
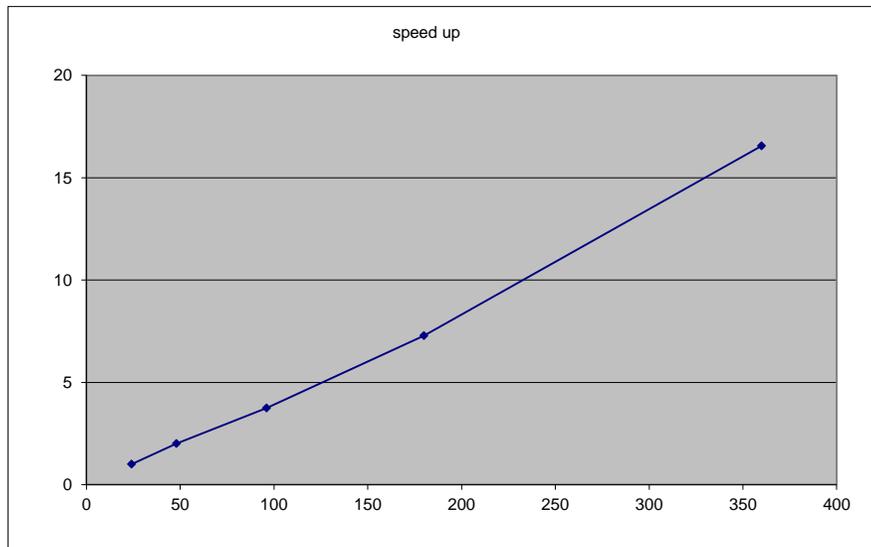


# Why scaling?

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

CG-POP mpi2s1D - 180x120

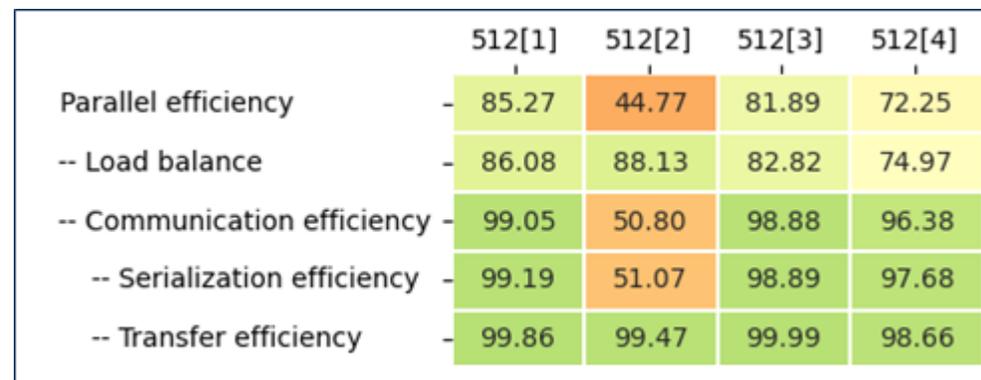
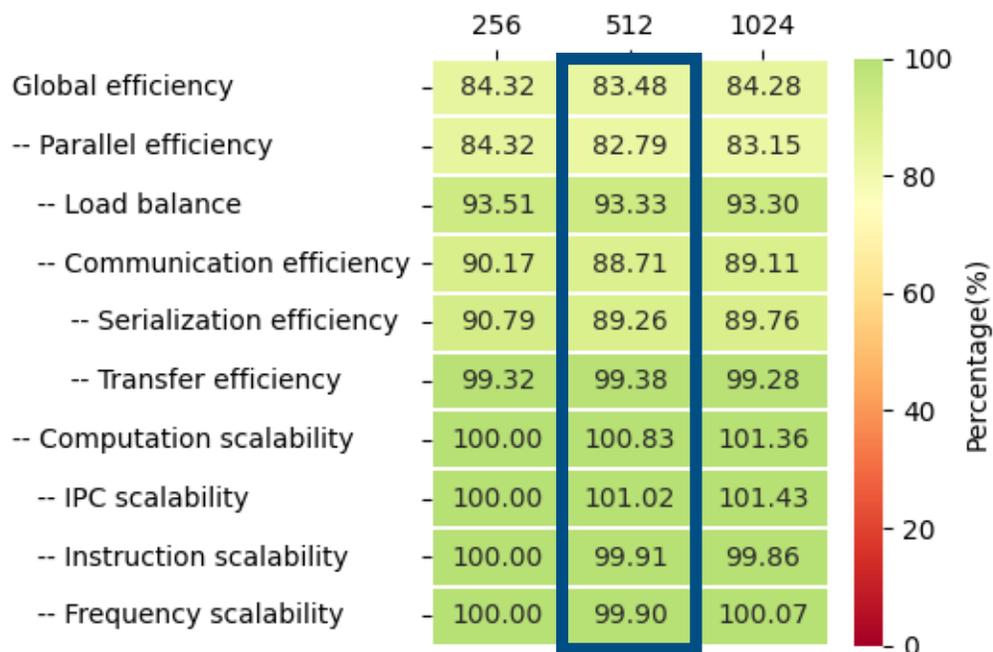
Good scalability !!  
Should we be happy?



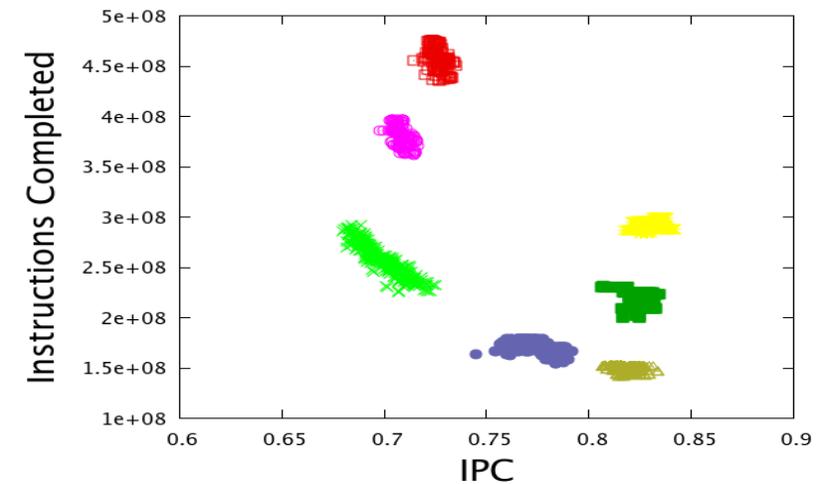
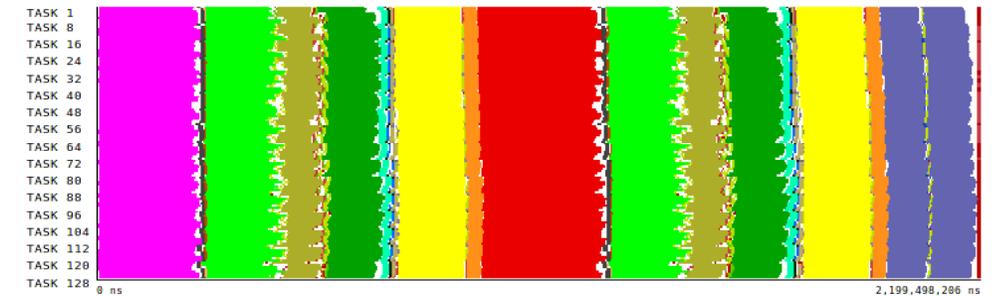
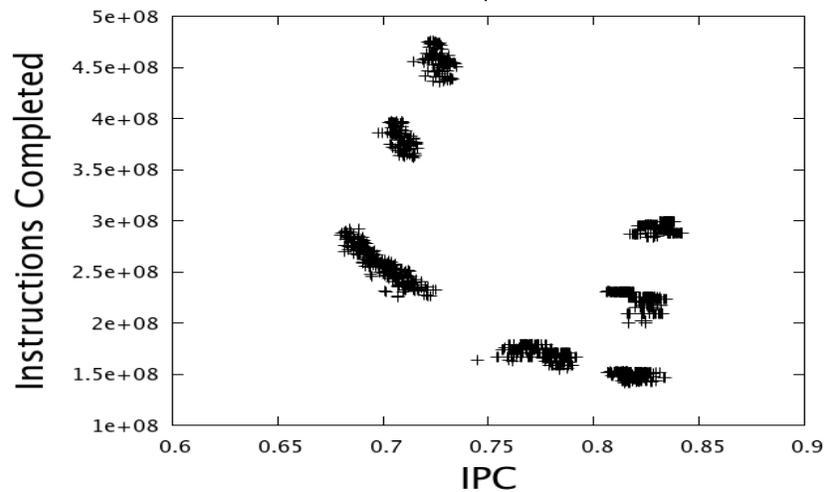
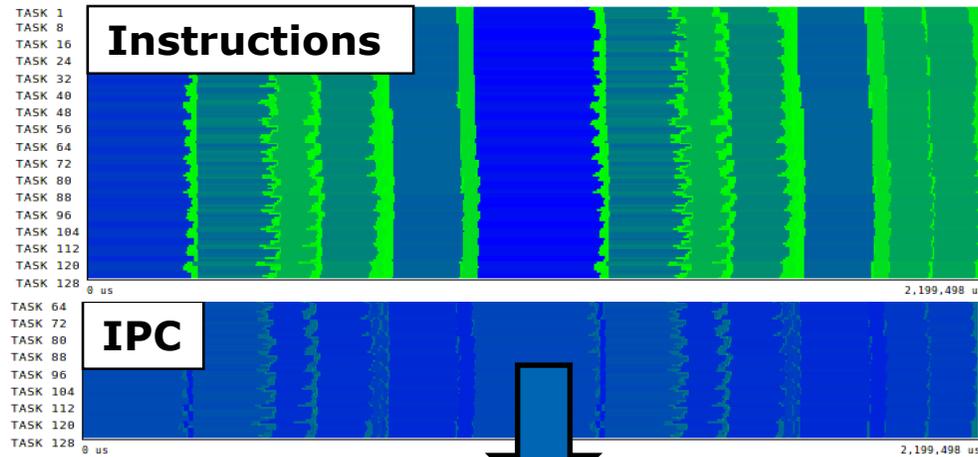
# Analytics

## Basic Analysis

- Scripts that automatically compute the efficiency model from a Paraver trace (or many traces for scalability studies)
- Dig down from global to detailed efficiencies

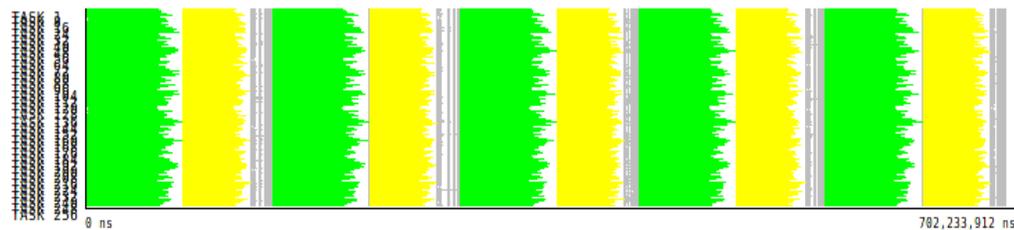


# Using Clustering to identify structure

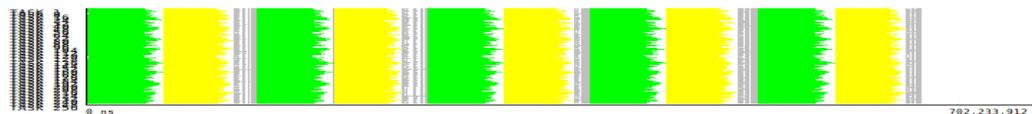
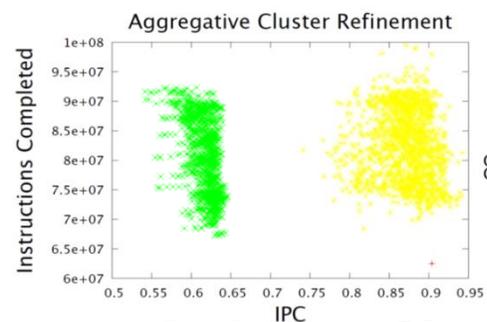


# Integrating models and analytics

What if ....

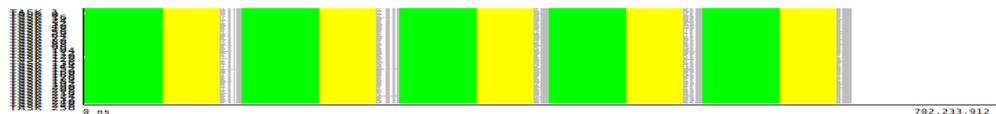
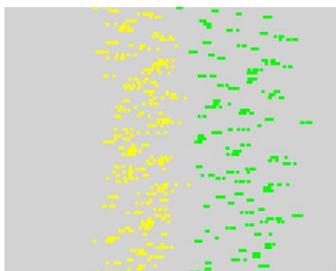


... we increase the IPC of Cluster1?



13% gain

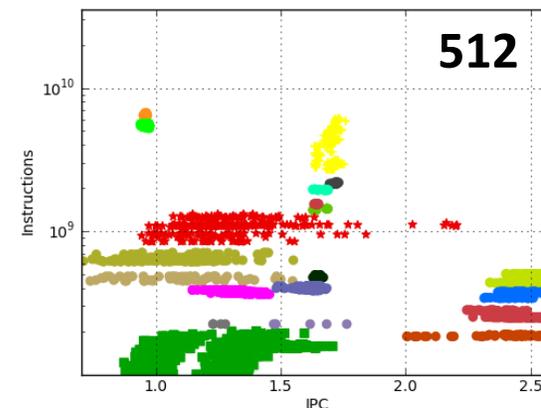
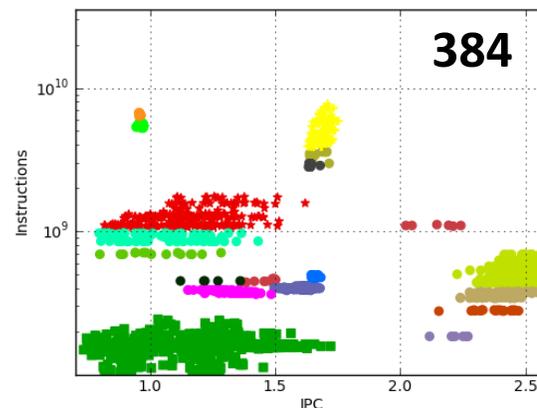
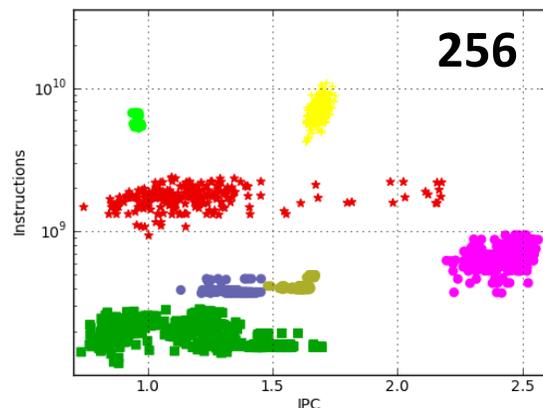
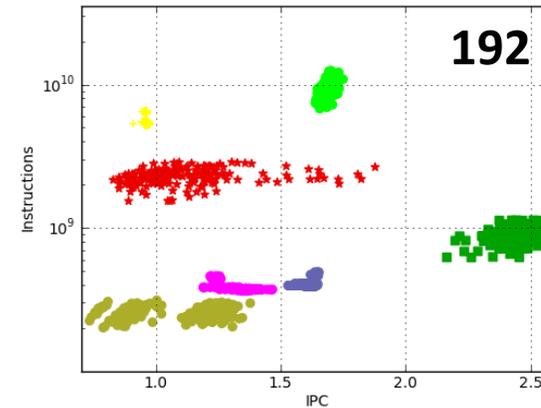
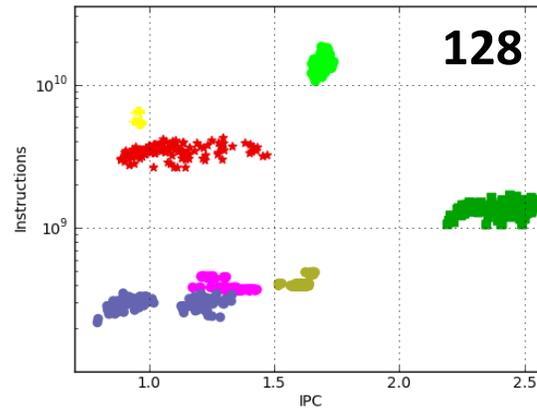
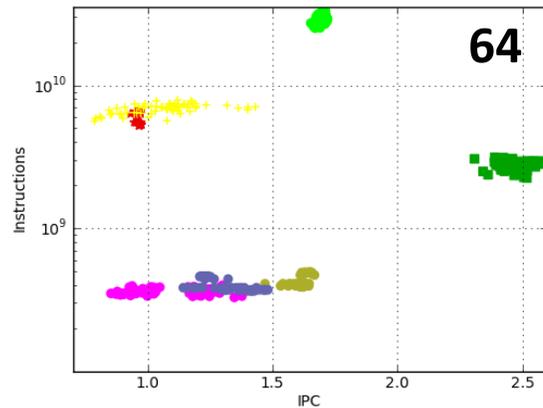
... we balance Clusters 1 & 2?



19% gain

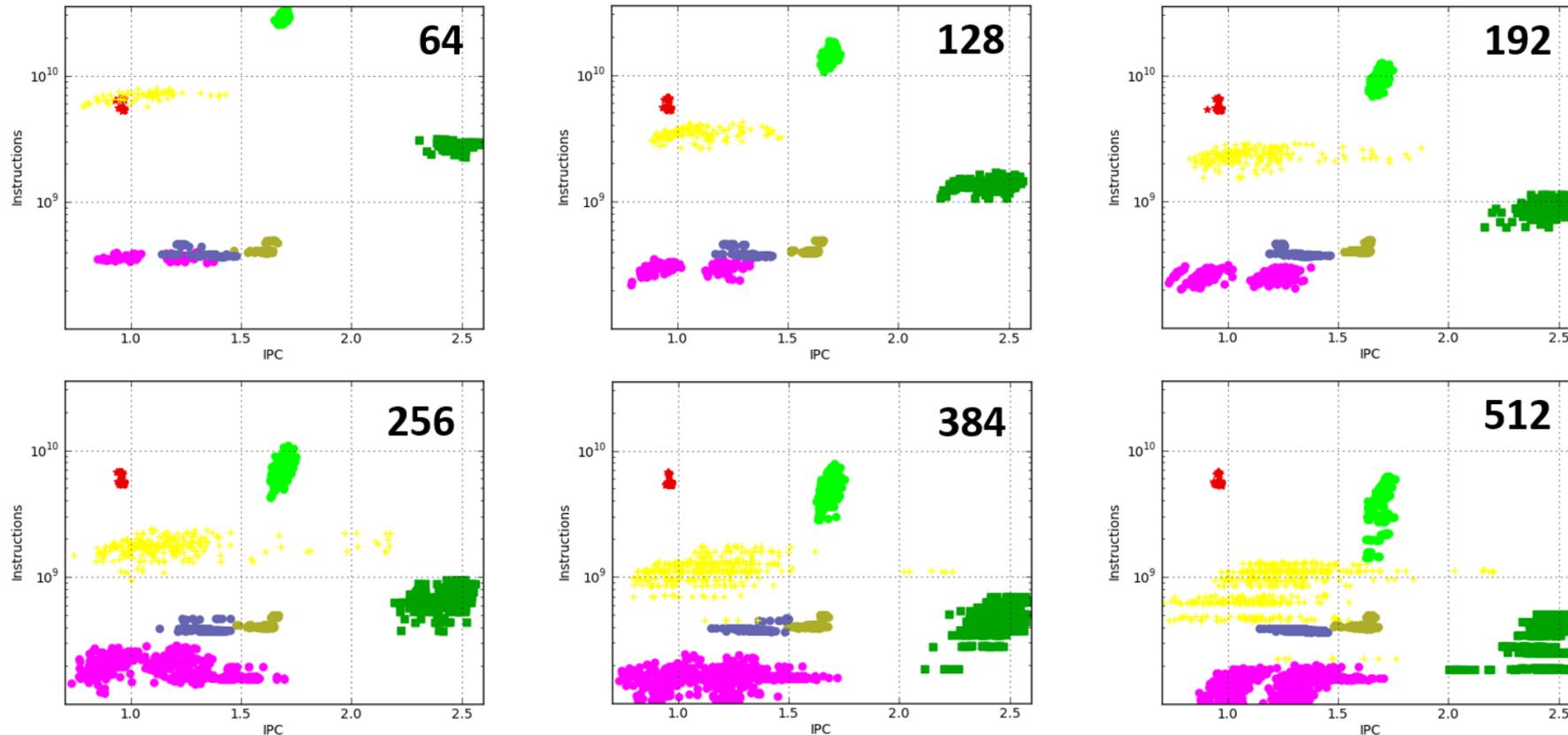
# Tracking scalability through clustering

- Study the scalability of the computing regions
  - Increasing the scale from 64 to 512 tasks



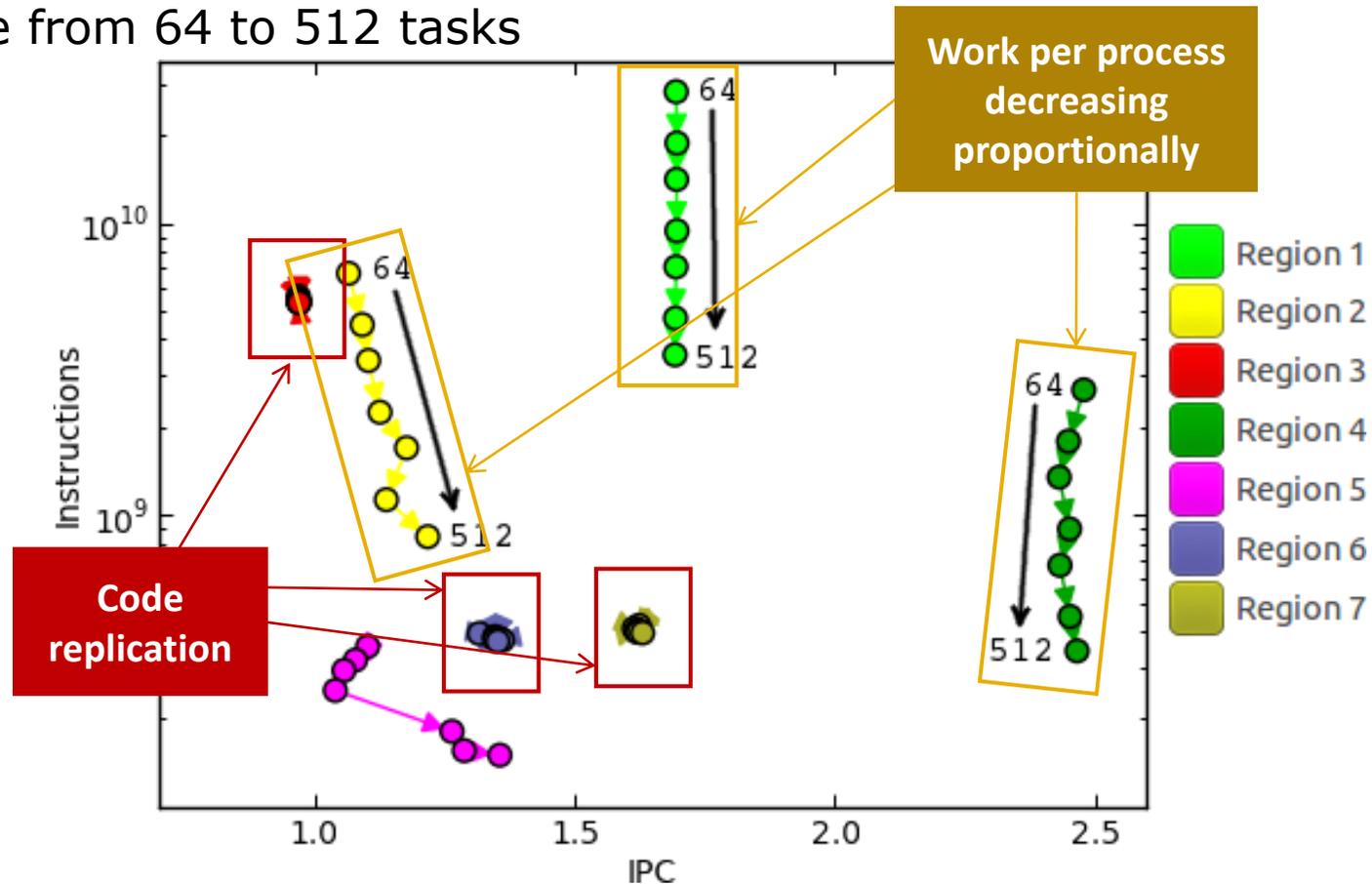
# Tracking scalability through clustering

- Study the scalability of the computing regions
  - Increasing the scale from 64 to 512 tasks



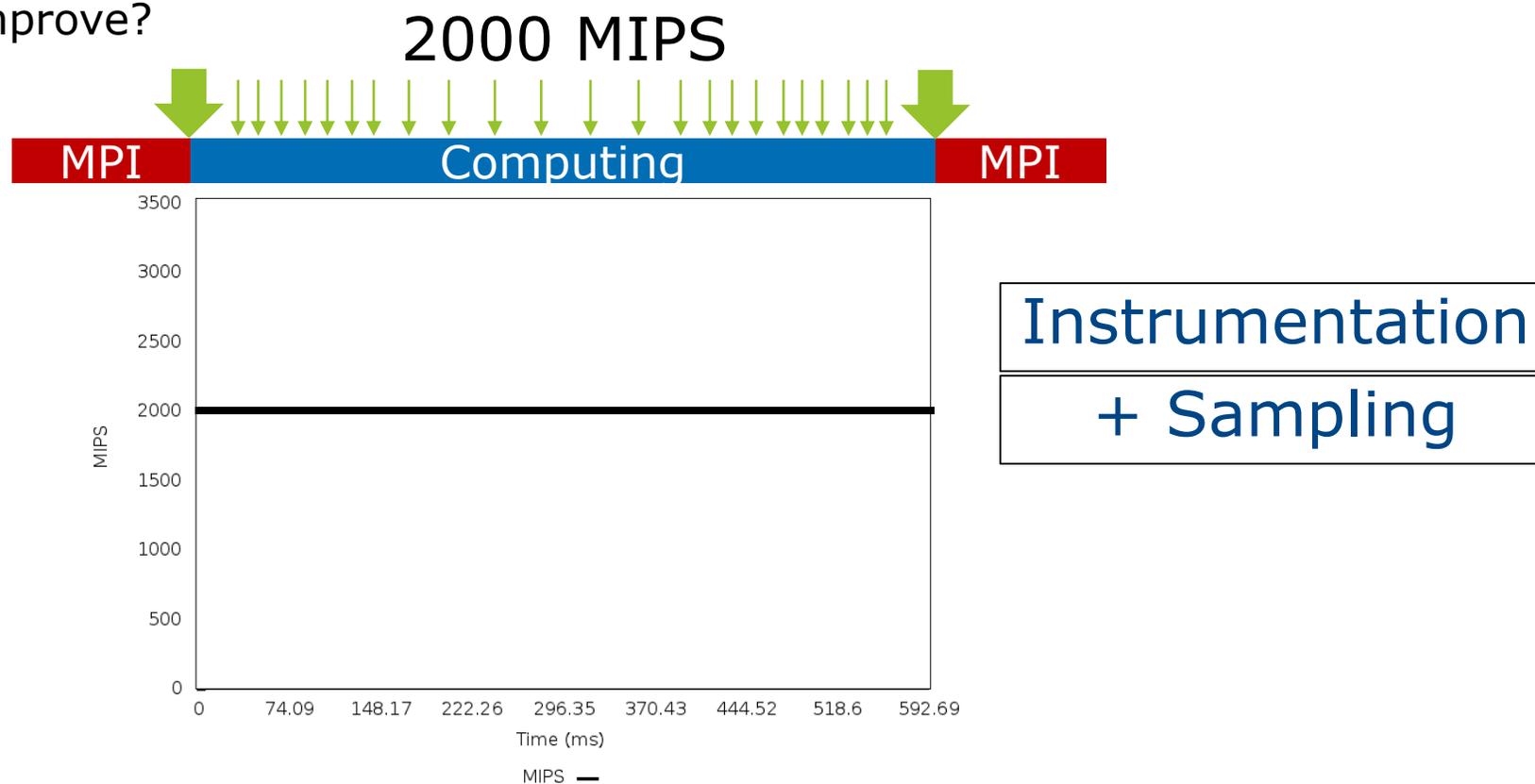
# Tracking scalability through clustering

- Study the scalability of the computing regions
  - Increasing the scale from 64 to 512 tasks

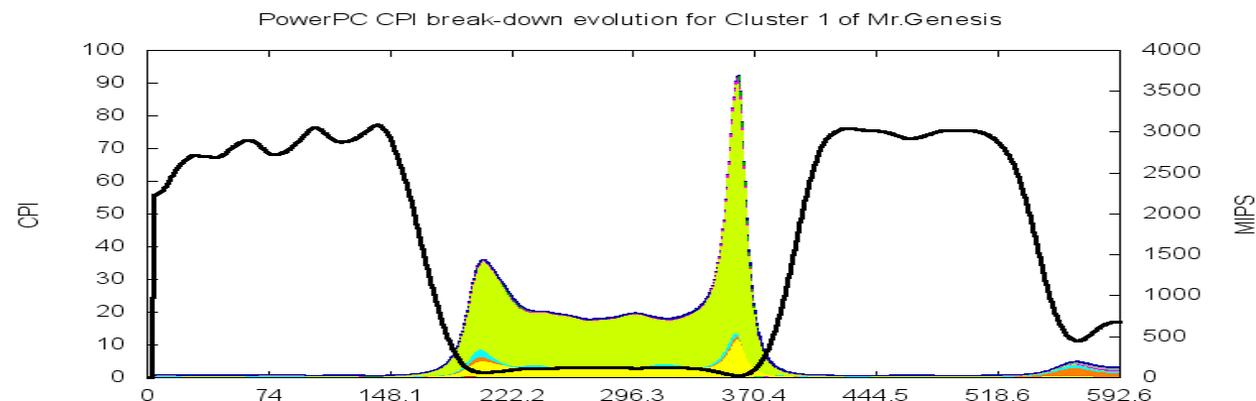


## Folding to increase details

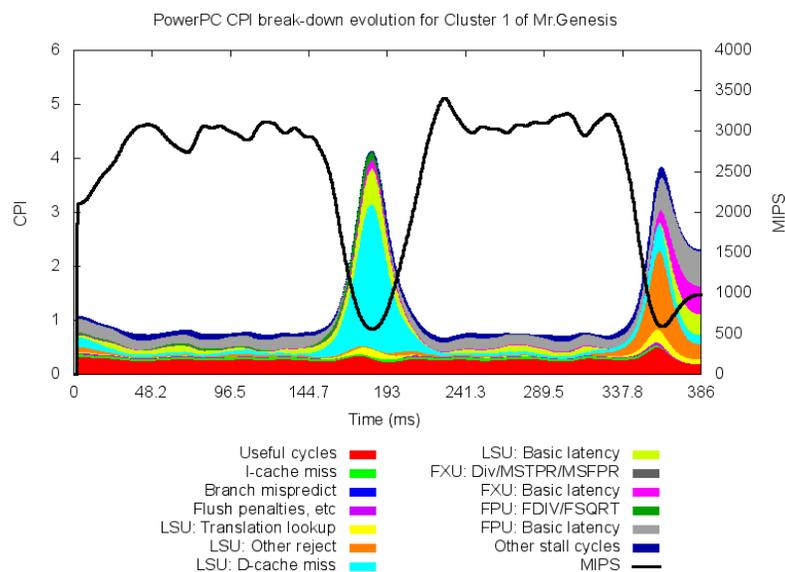
- What is the performance of a serial region?
  - Is it good enough?
  - Is it easy to improve?



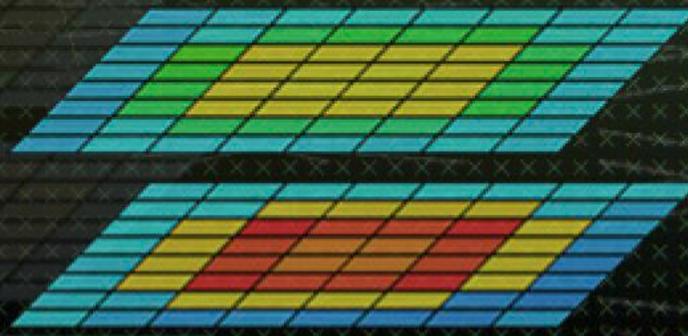
# Folding: CPI and HWC stack models



MR-GENESIS



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



# 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
- Behavioural structure? Variability?

**Tutorial:  
Introduction to Paraver &  
Dimemas methodology**

# BSC Tools web site

- <https://tools.bsc.es>

- Downloads

- Sources / Binaries



- Documentation

- Training guides
    - Tutorial slides

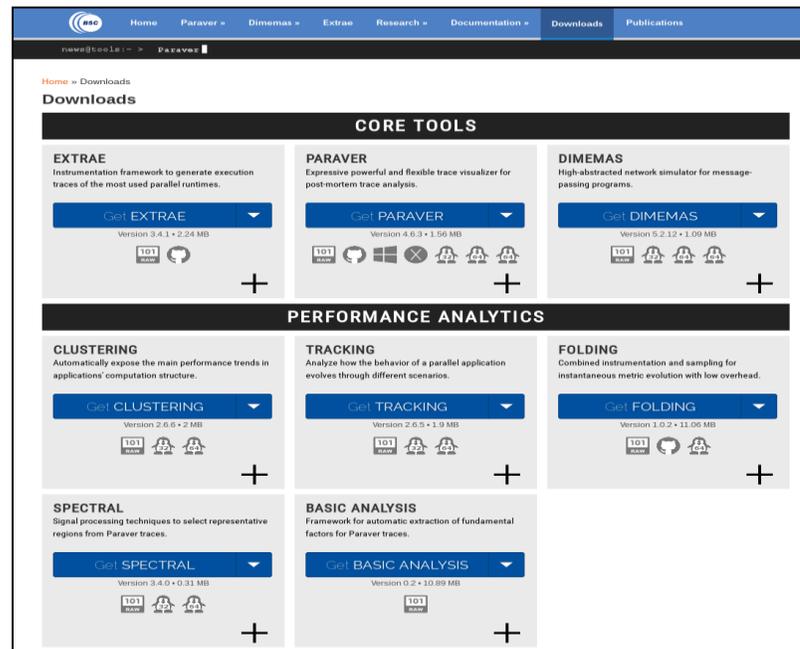
- Getting started

- Start wxparaver

- Help → Tutorials

- Follow training guides

- Paraver introduction (MPI): Navigation & Basic Understanding of Paraver operation



---

## Takeaway:

- The importance of understanding  
→ **Keep asking questions**
  - Use your brain  
→ **Use visual tools**
  - The devil is in the details  
→ **Do not miss them**
  - Don't over-theorize about your code  
→ **Look at it open-minded**
-