# Understanding applications
# using the
# BSC performance tools

Germán Llort, Lau Mercadal

✉ tools@bsc.es

Barcelona Supercomputing Center

# Humans are visual creatures

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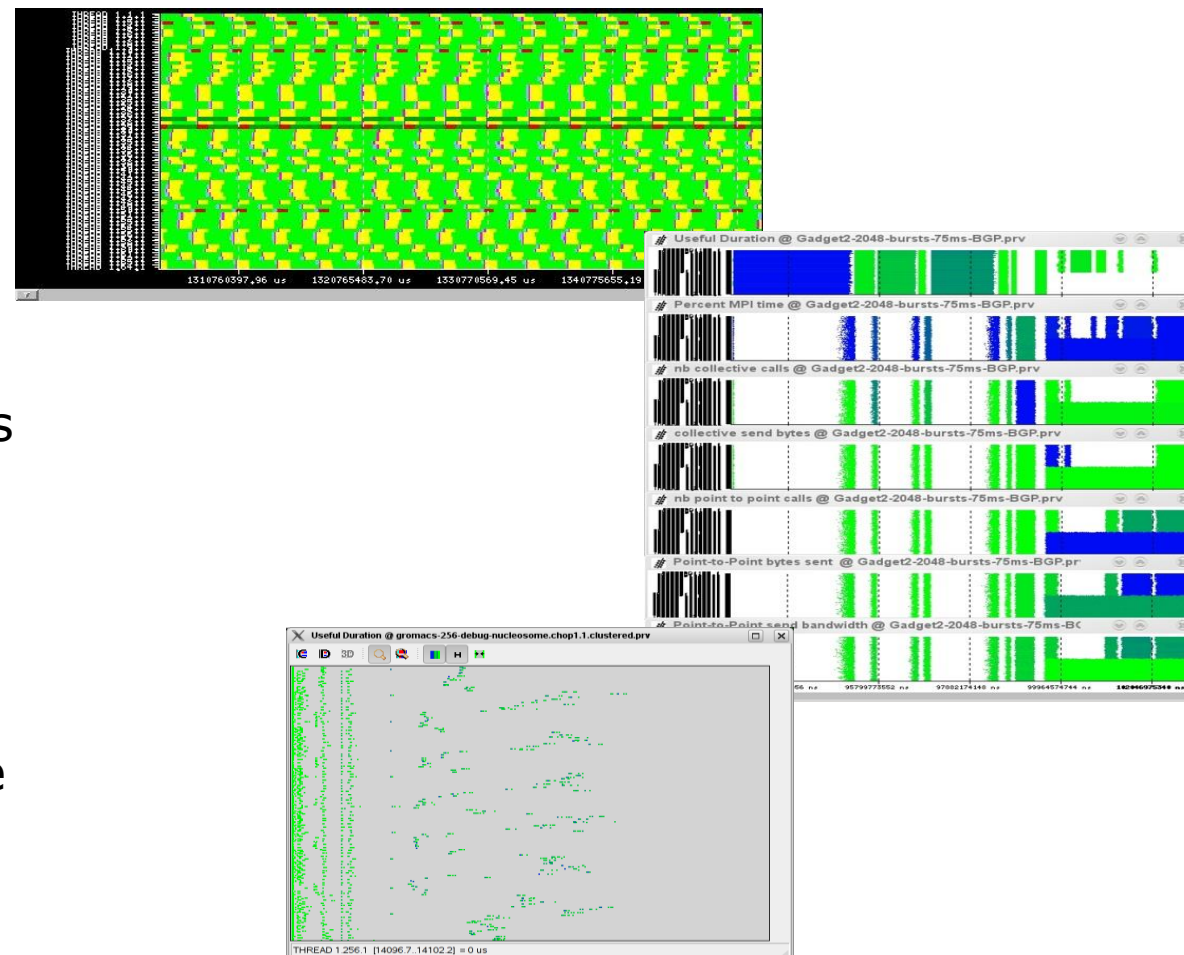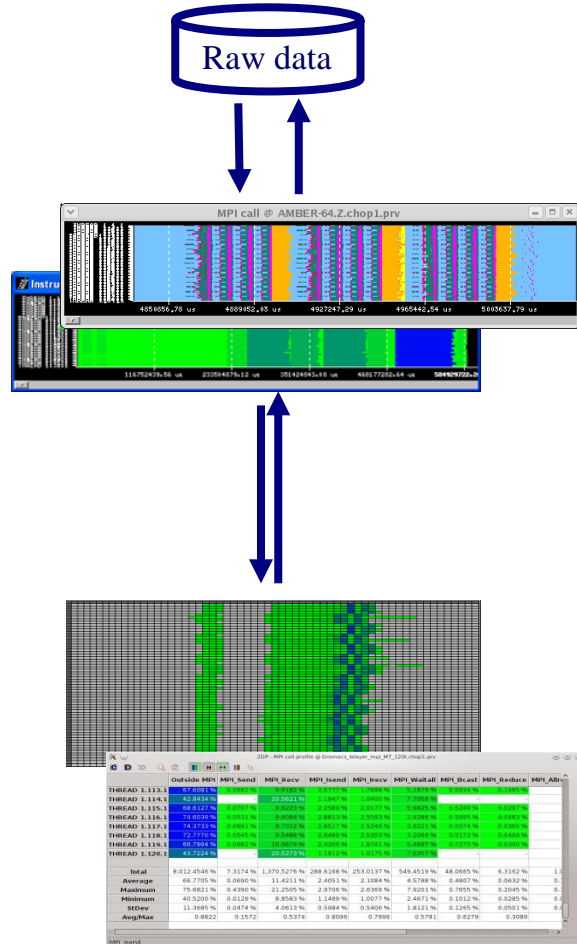  - What do you see on the pictures?

# Our Tools

- Since 1991
- Based on traces
- Open Source
  - https://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

# Paraver: Performance data browser



Timelines

2/3D tables
(Statistics)

Trace visualization/analysis

+ trace manipulation

Goal = Flexibility

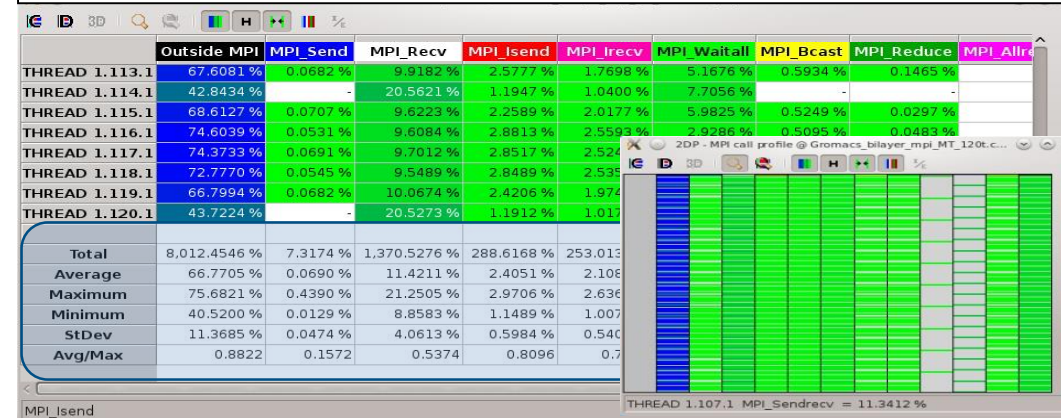No semantics

Programmable

Comparative analyses

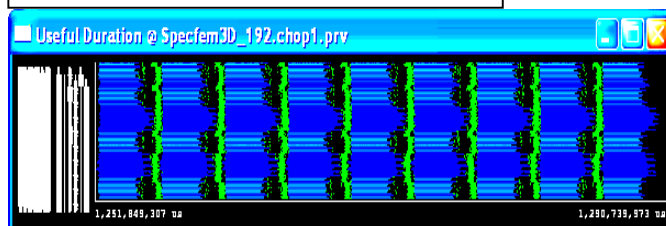Multiple traces

Synchronize scales
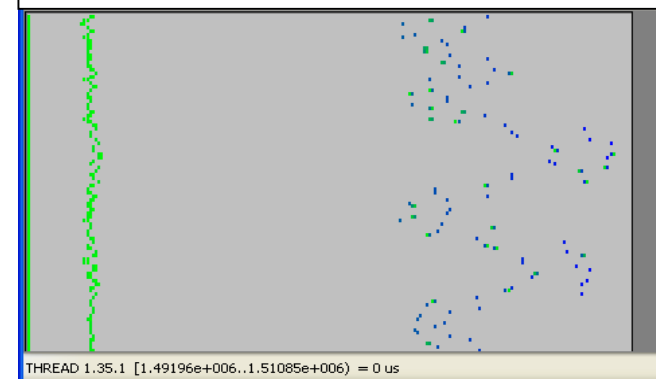
# From timelines to tables
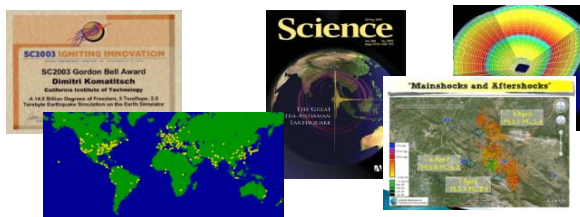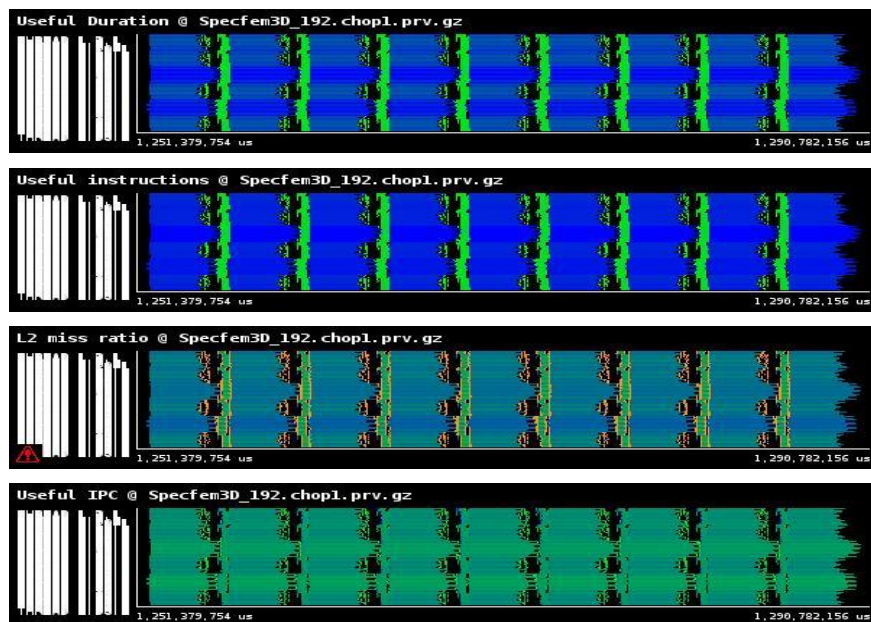


MPI calls

MPI calls profile

Useful Duration

Histogram Useful Duration
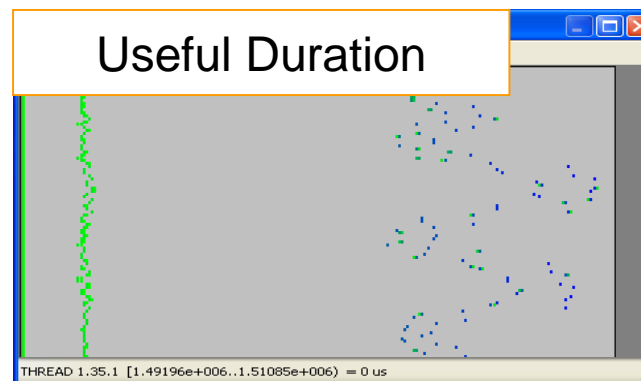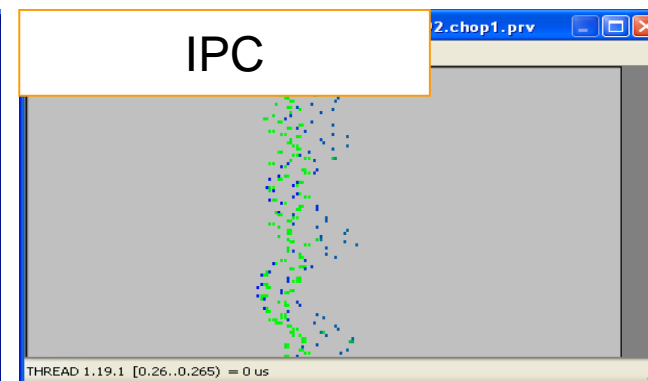
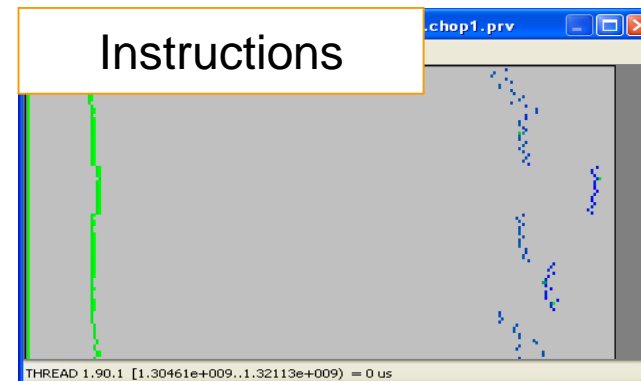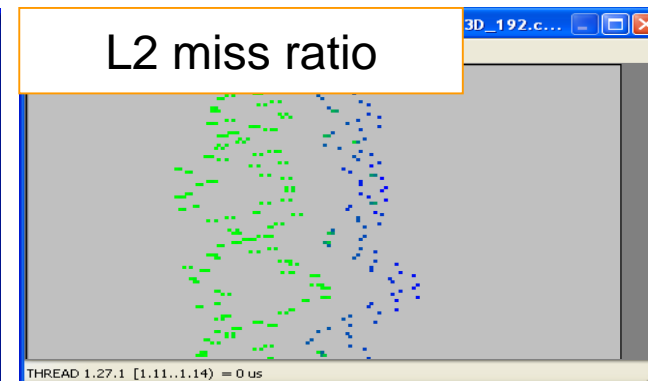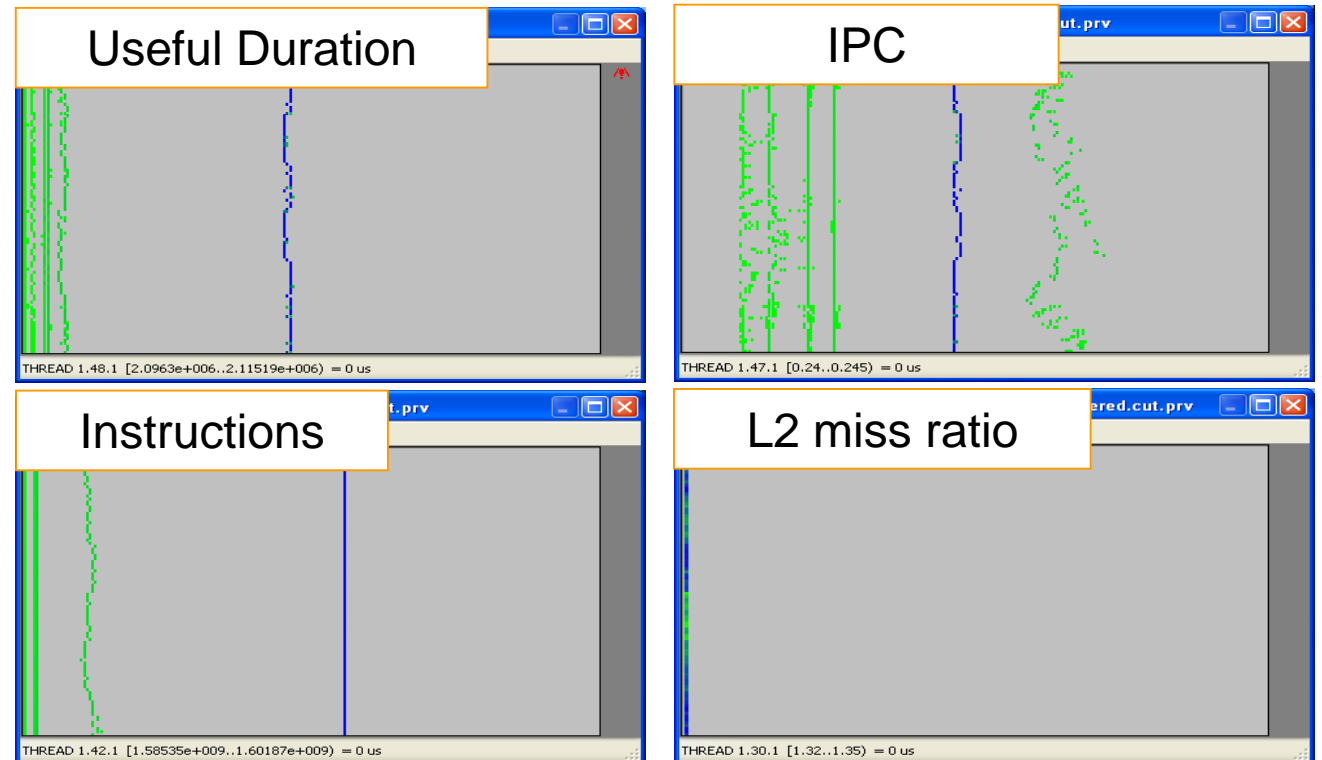# Analyzing variability through histograms and timelines
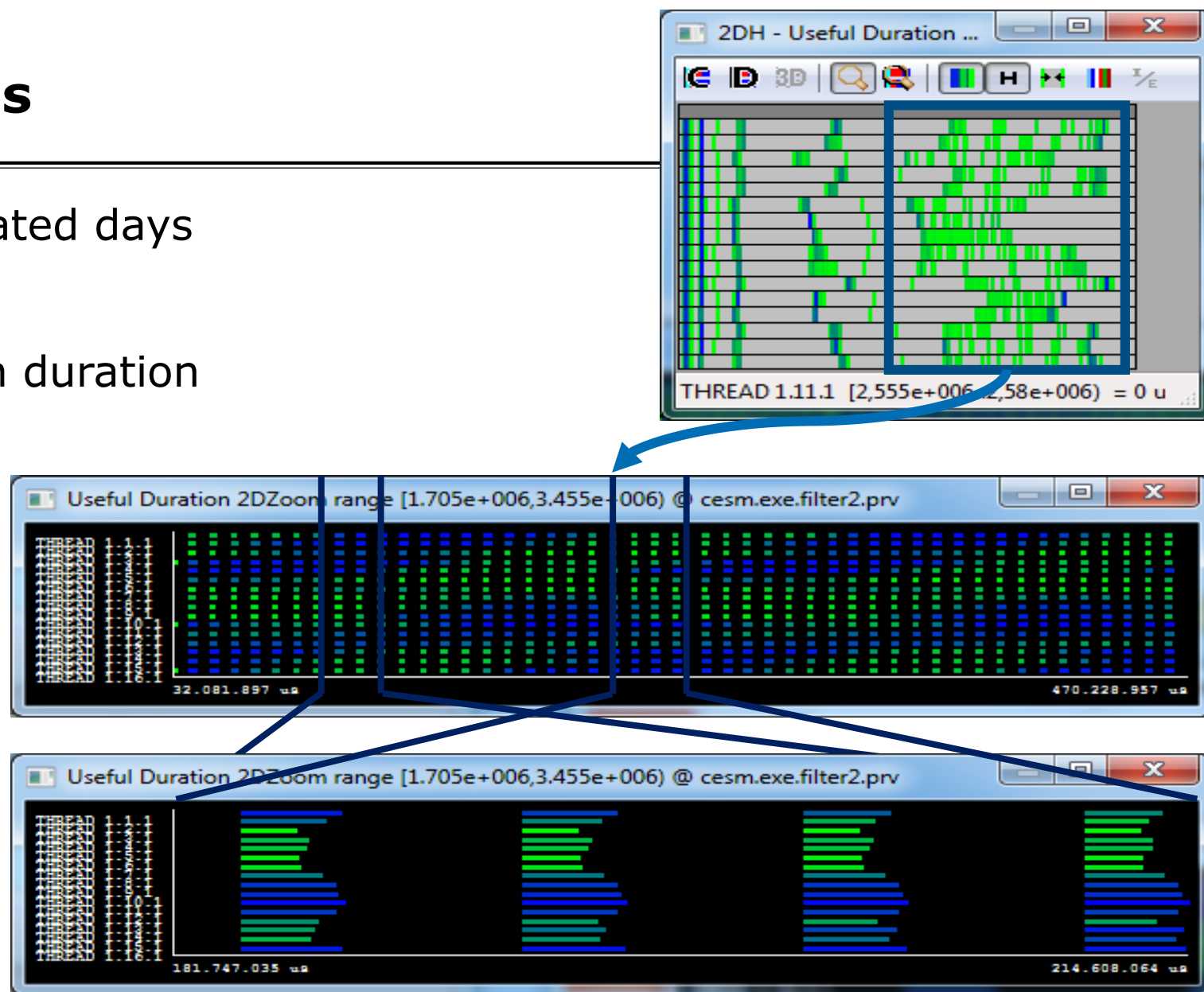
SPECFEM3D

# Analyzing variability through histograms and timelines

- By the way: six months later ….
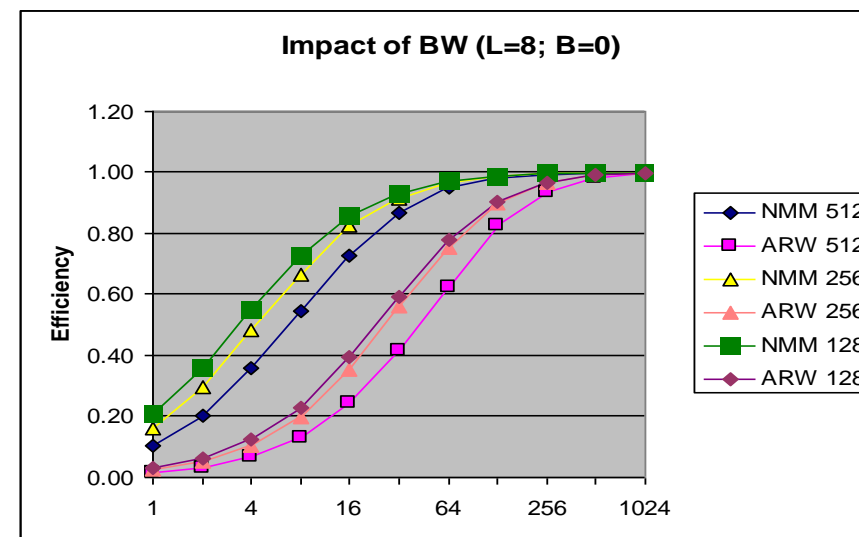
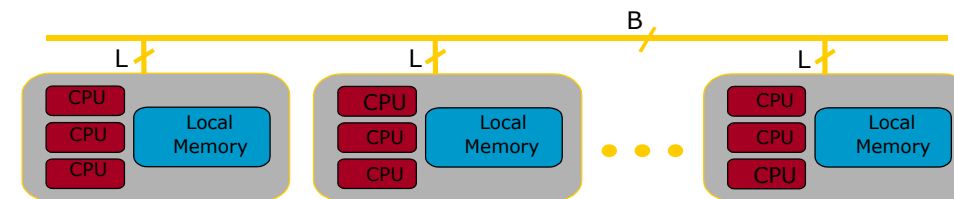# From tables to timelines



- CESM: 16 processes, 2 simulated days

- Histogram useful computation duration shows high variability
  - How is it distributed?



- Dynamic imbalance
  - In space and time
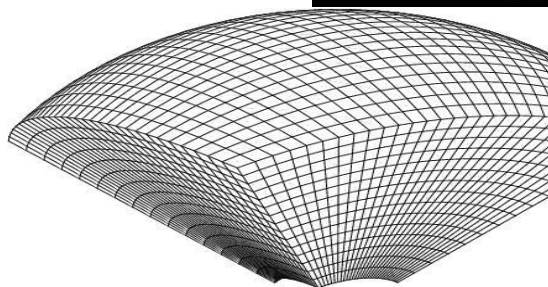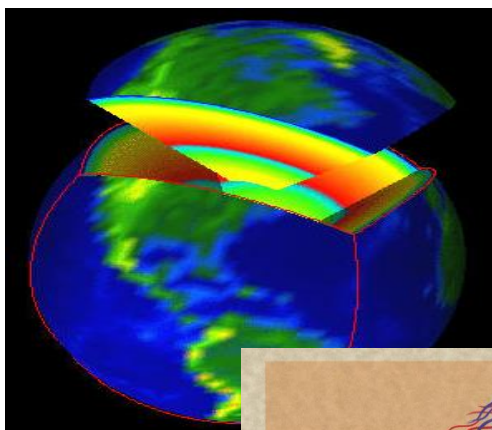  - Day and night.
  - Season ? ☺

# 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
  - Modeling nominal

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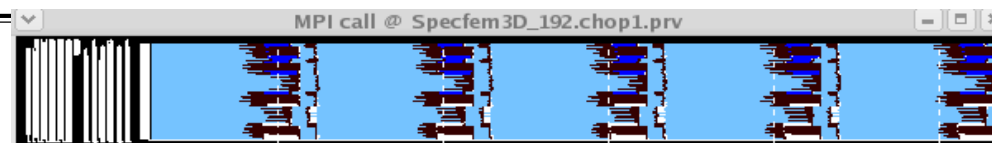  - What-if from selected time window

# What if we had asynchronous communications?

- SPECFEM3D



Courtesy Dimitri Komatitsch



Real

Ideal

Prediction
MN

Prediction
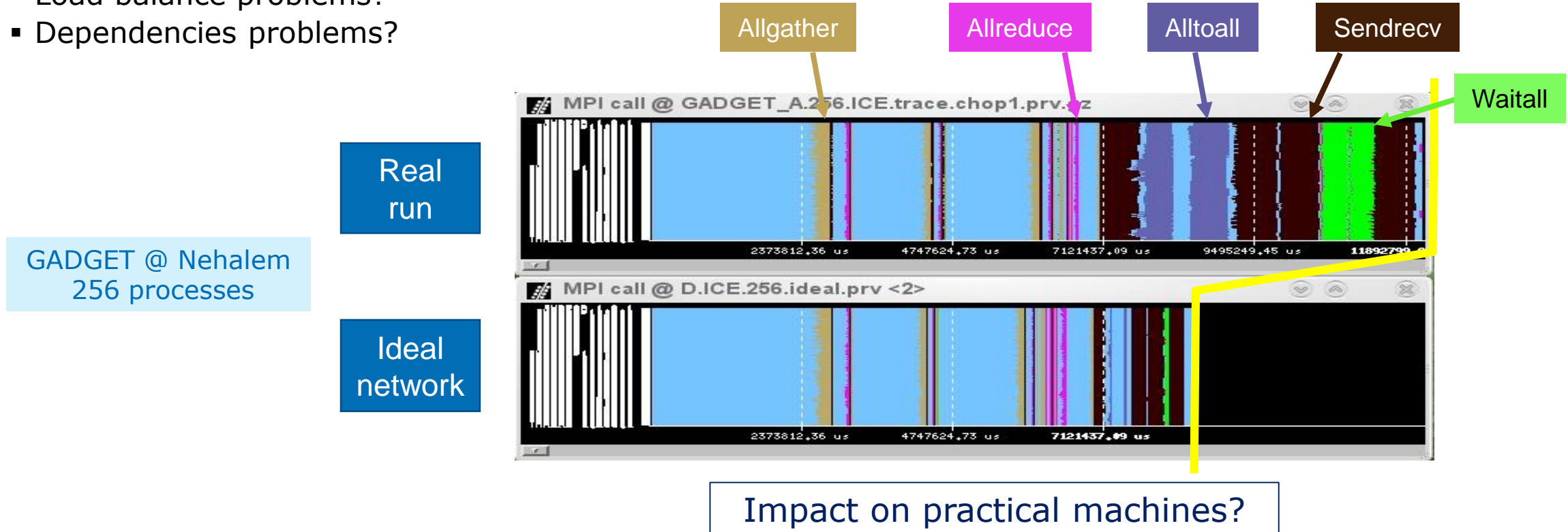100MB/s

Prediction
10MB/s
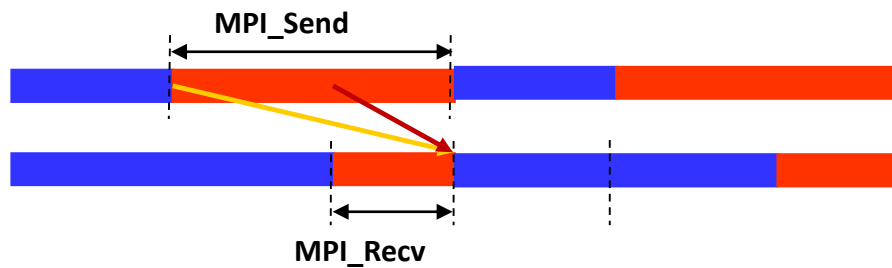
Prediction
5MB/s

Prediction
1MB/s

# Ideal machine

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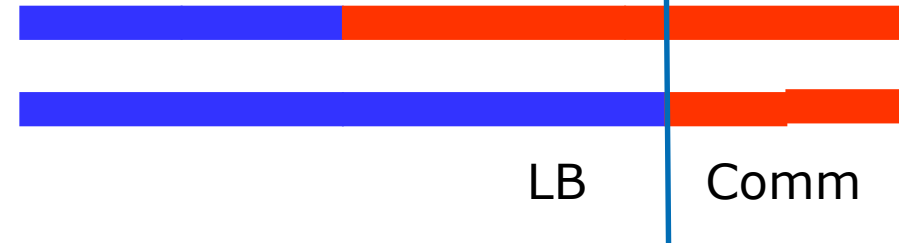    - Load balance problems?
    - Dependencies problems?



GADGET @ Nehalem
256 processes

Real run

Ideal network

Allgather    Allreduce    Alltoall    Sendrecv    Waitall

Impact on practical machines?

# Models

# Parallel efficiency model



- Parallel efficiency = LB eff * Comm eff

# Parallel efficiency refinement: LB * µLB * Transfer



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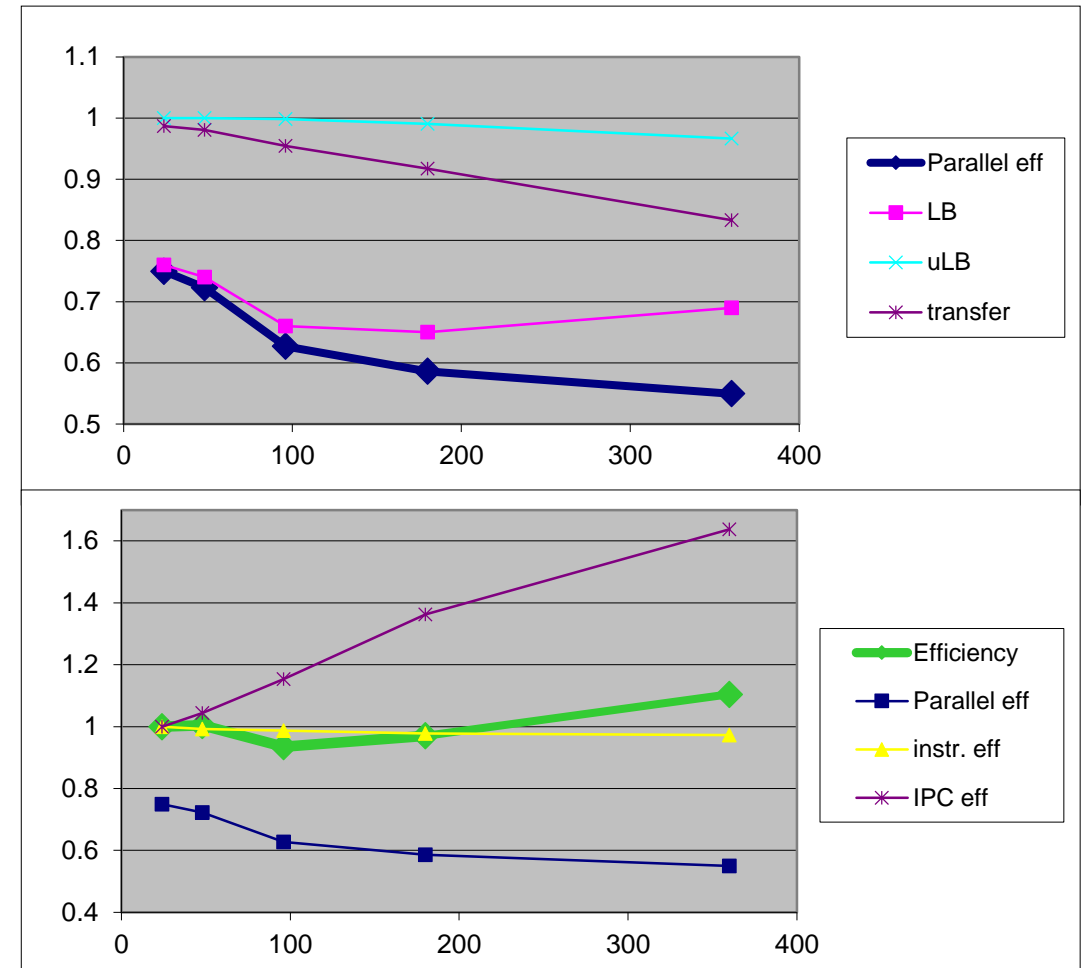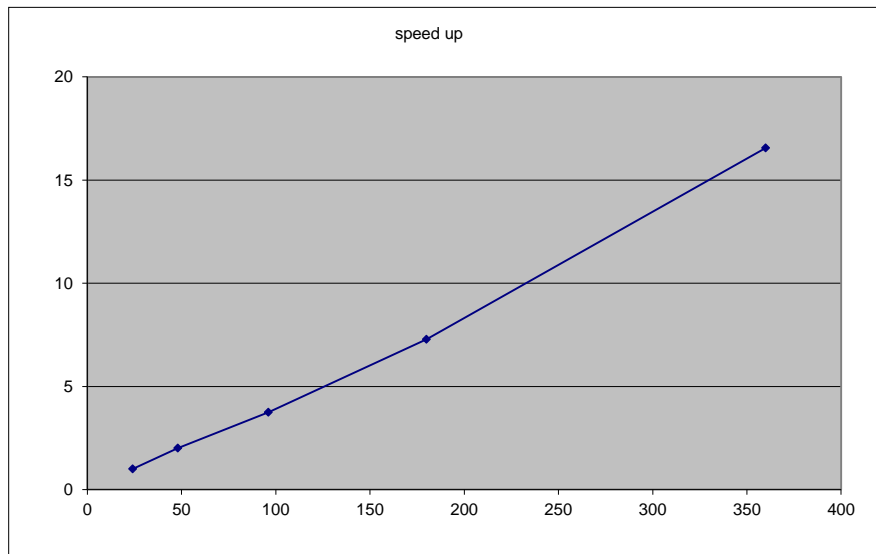- Dimemas ideal network → Transfer (efficiency) = 1

# Why scaling?

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

CG-POP mpi2s1D - 180x120

Good scalability !!
Should we be happy?



speed up

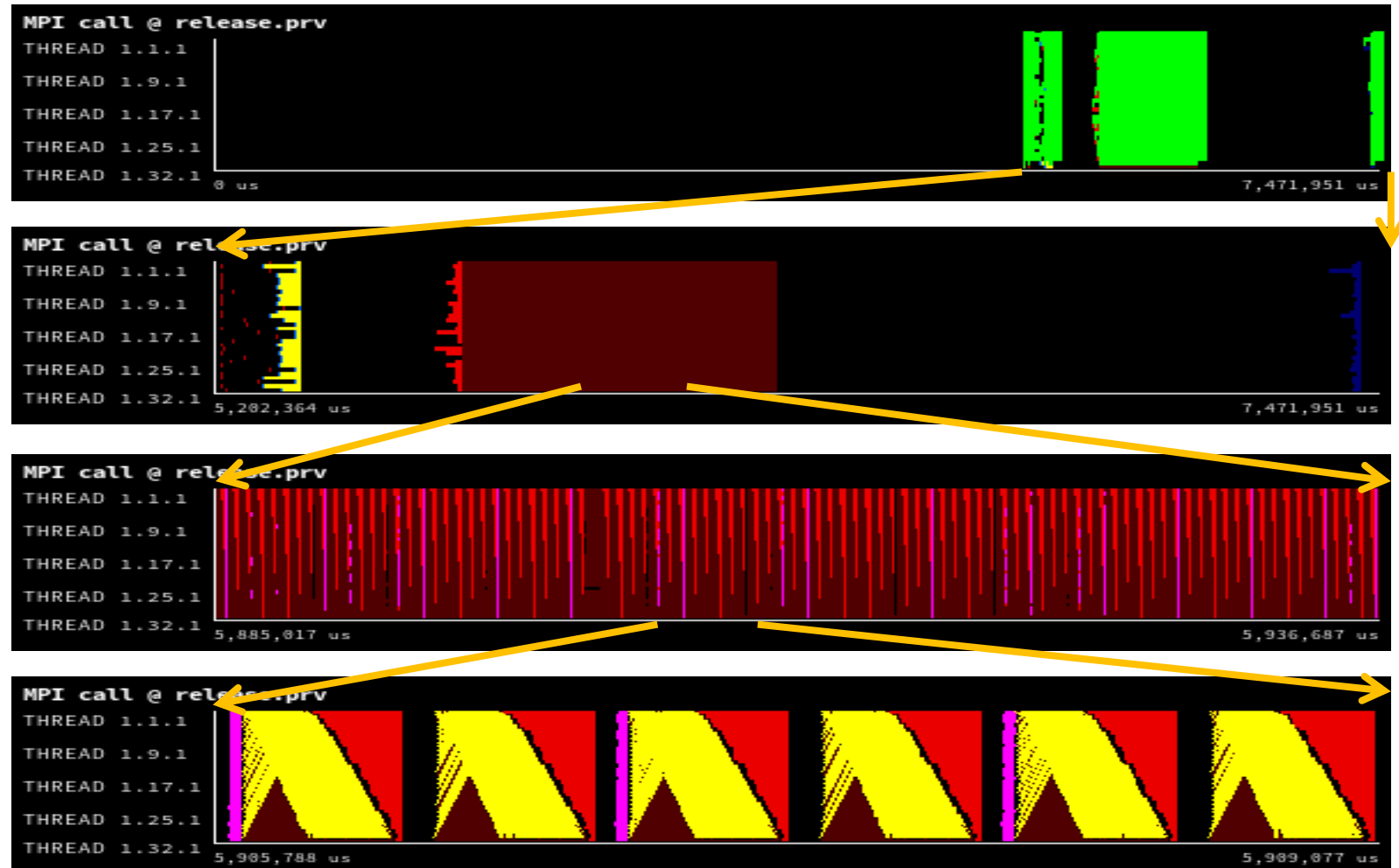# Why efficient?

Parallel efficiency =93.28
Communication = 93.84

Parallel efficiency = 77.93
Communication = 79.79

Parallel efficiency = 28.84
Communication = 30.42

# Clustering

# Using Clustering to identify structure

# Integrating models and analytics

**What if ….**



PEPC

**… we increase the IPC of Cluster1?**



**13% gain**

**… we balance Clusters 1 & 2?**



**19% gain**

# 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
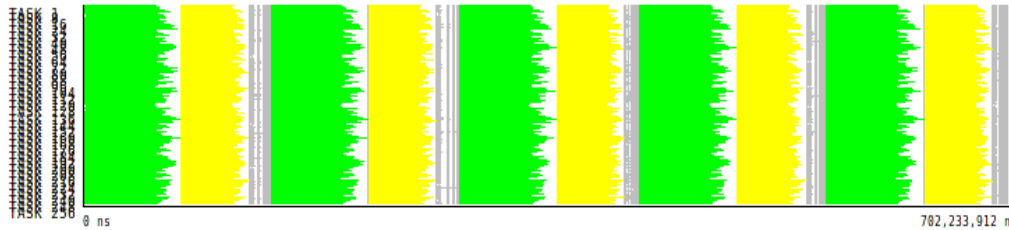
- Scalability – code replication?
  - Total instructions

- Behavioral structure? Variability?

**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
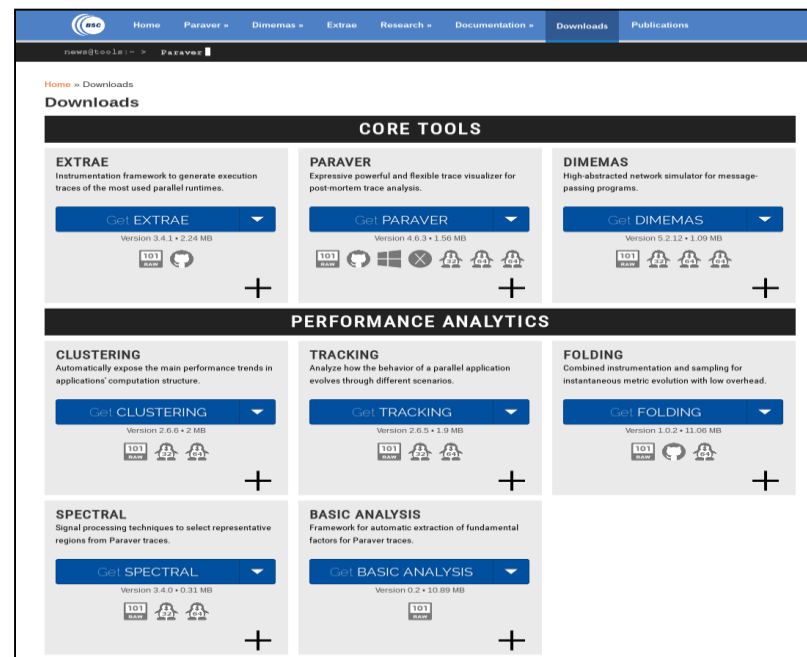  ▪ Follow training guides
    ▪ Paraver introduction (MPI): Navigation & Basic Understanding of Paraver operation

# Paraver Demo

# 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 |

⚠️ **Higher Parallel Efficiency ≠ Faster!**