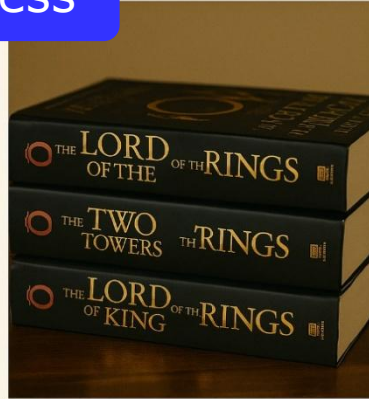# Understanding applications with BSC Tools

Lau Mercadal, Sandra Méndez, Germán Llort

✉ tools@bsc.es

Barcelona Supercomputing Center

Performance Optimisation and Productivity 3 (101143931)

Process

Store

Identify

# BSC Tools

- Since 1991
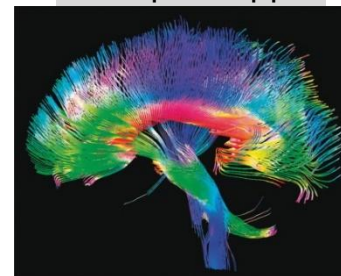
- Based on traces

- Open Source
  - https://tools.bsc.es

- Core tools:
  - **Extrae** – instrumentation
  - **Paraver** – offline trace analysis
  - **Dimemas** – message passing simulator

- Focus
  - Detail, variability, flexibility
  - Behavioural structure vs. syntactic structure
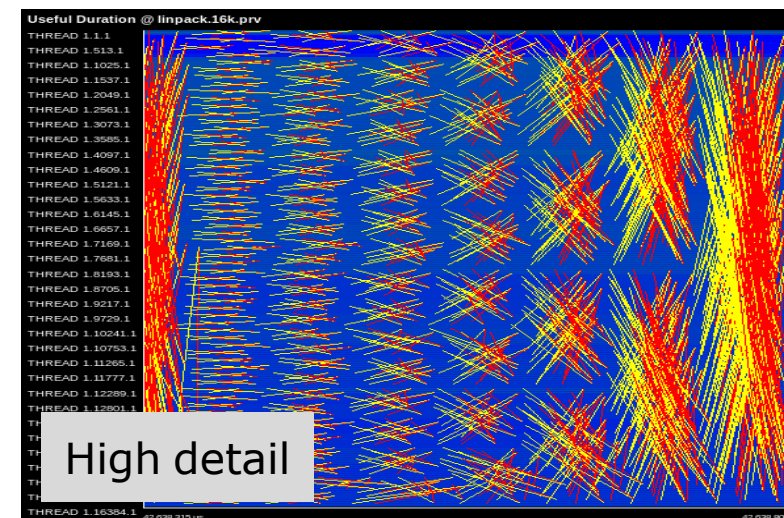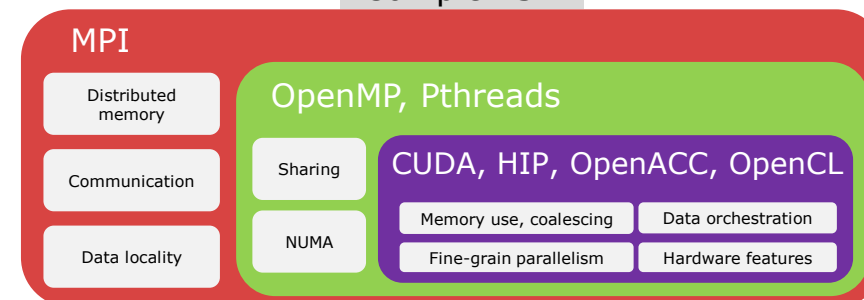  - Intelligence: Performance Analytics

Complex Apps

Complex Hw

Complex Sw



MPI

| Distributed memory |
| Communication |
| Data locality |

OpenMP, Pthreads

Sharing

NUMA

CUDA, HIP, OpenACC, OpenCL

| Memory use, coalescing | Data orchestration |
| Fine-grain parallelism | Hardware features |



Useful Duration @ linpack.16k.prv
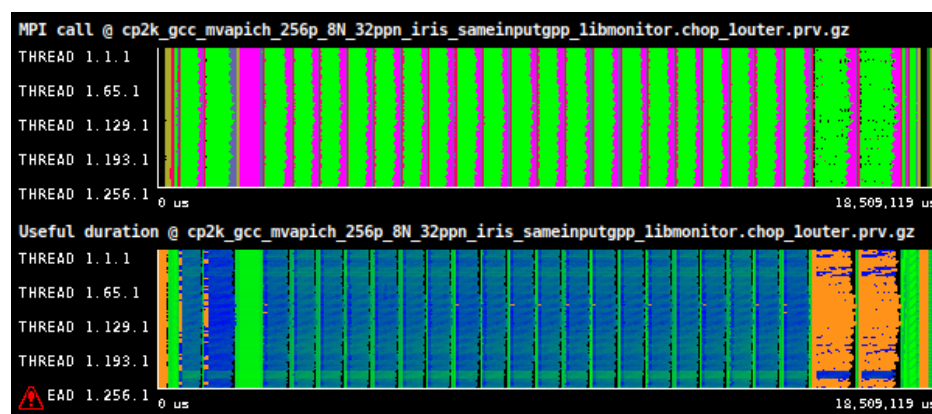
High detail

# Paraver

# Paraver – Performance data browser

```
1:151:1:147:1:294672917:294676549:1
2:151:1:147:1:294672917:50000001:0
3:151:1:147:1:294623693:294672917:158:1:154:1:294670636:297419889:4223536:8
```
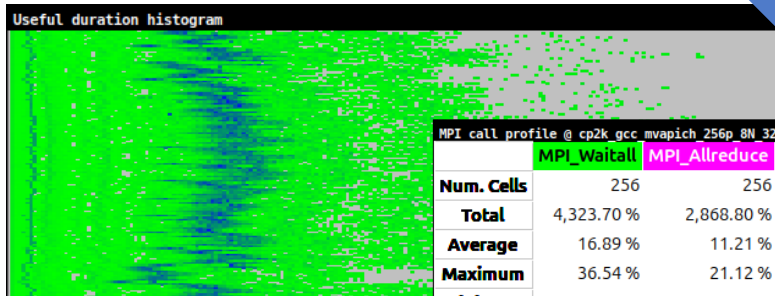
Trace visualization/analysis    Trace manipulation

Timelines



2/3D Histograms

2/3D tables (statistics)

Goal = Flexibility
No semantics
Programmable

Comparative analyses
Multiple traces
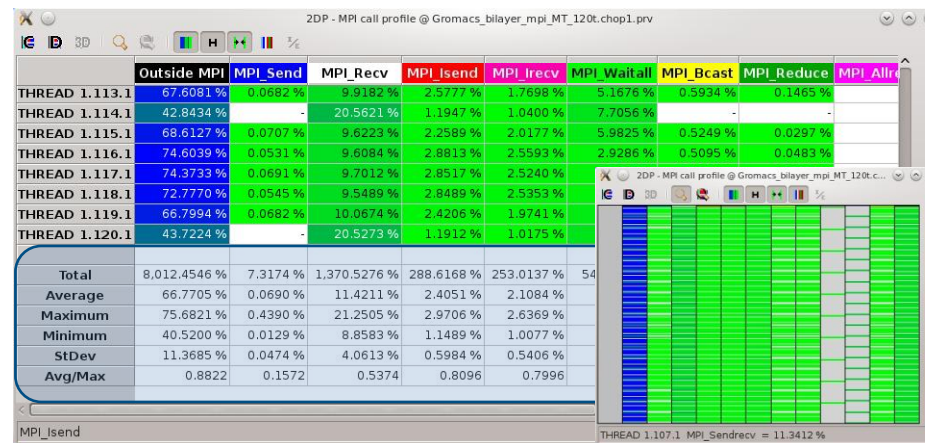Synchronize scales

# From timelines to tables

- Categorical metrics → colours
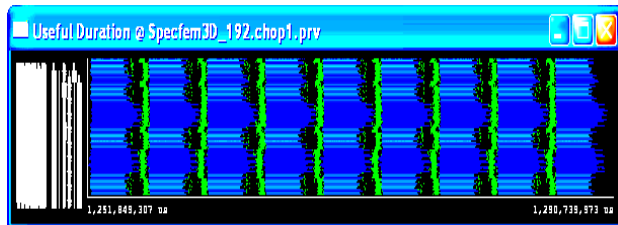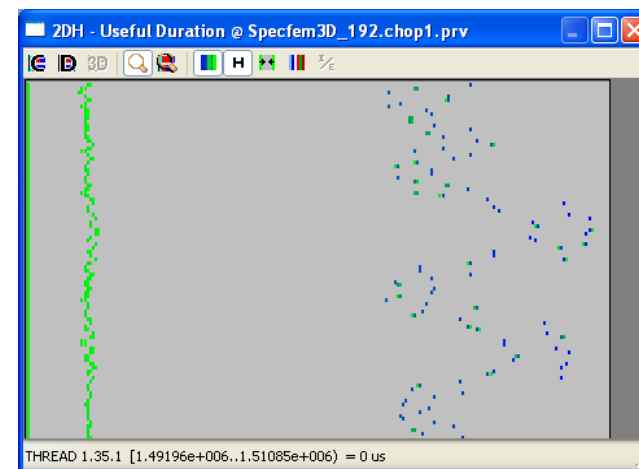  - Each category is assigned a distinct colour

MPI calls profile

MPI calls

- Continuous metrics → gradients
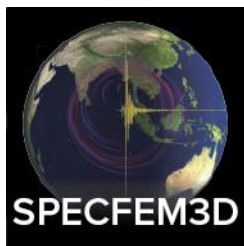  - **Green** to **blue** for **low** to **high** values

Compute-time histogram

Useful duration
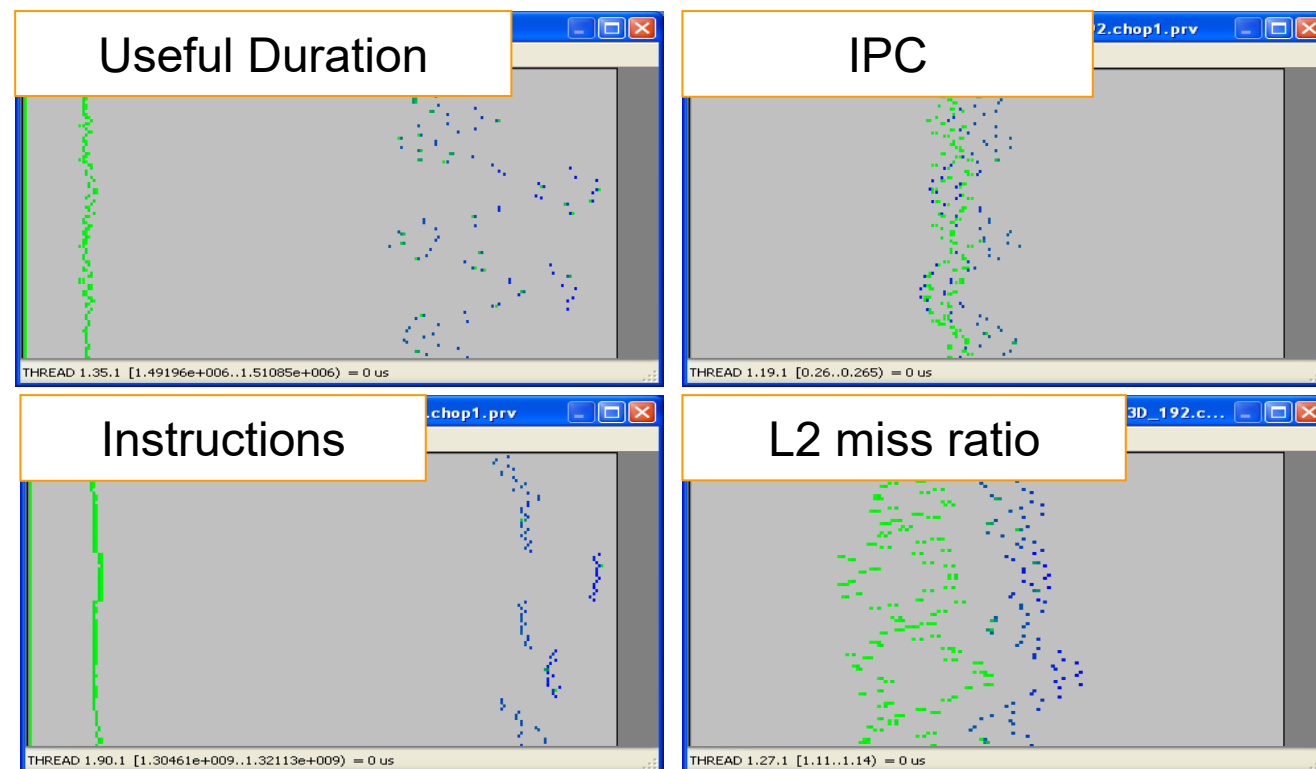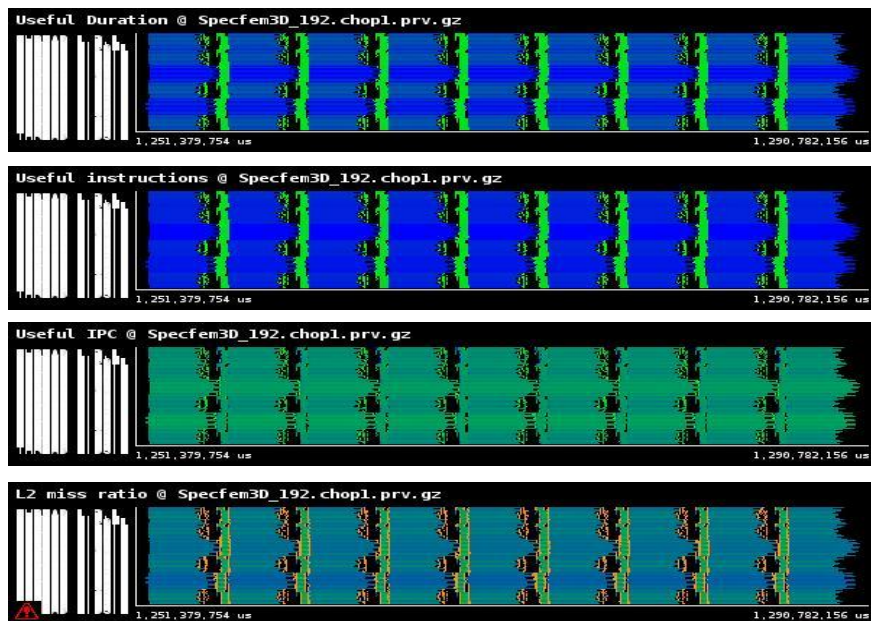
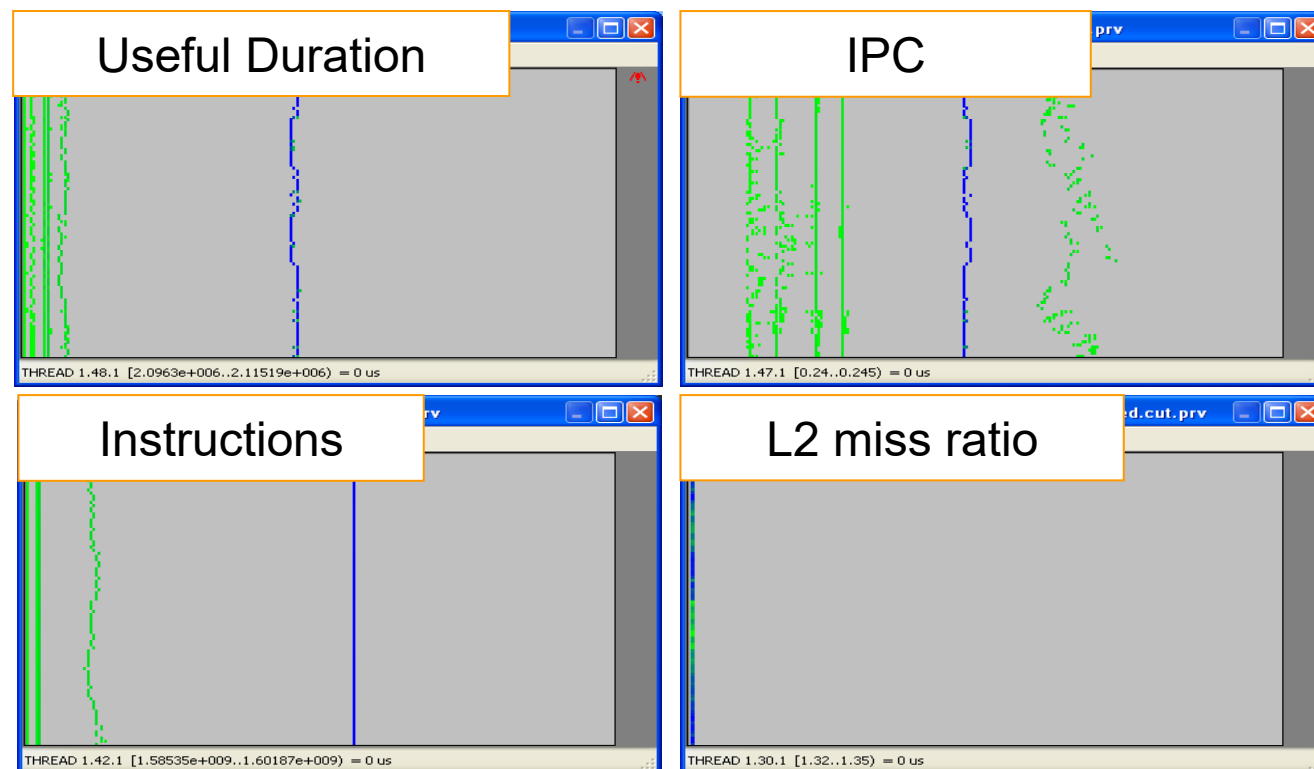# Analyzing variability through timelines and histograms



*Courtesy Dimitri Komatitsch*

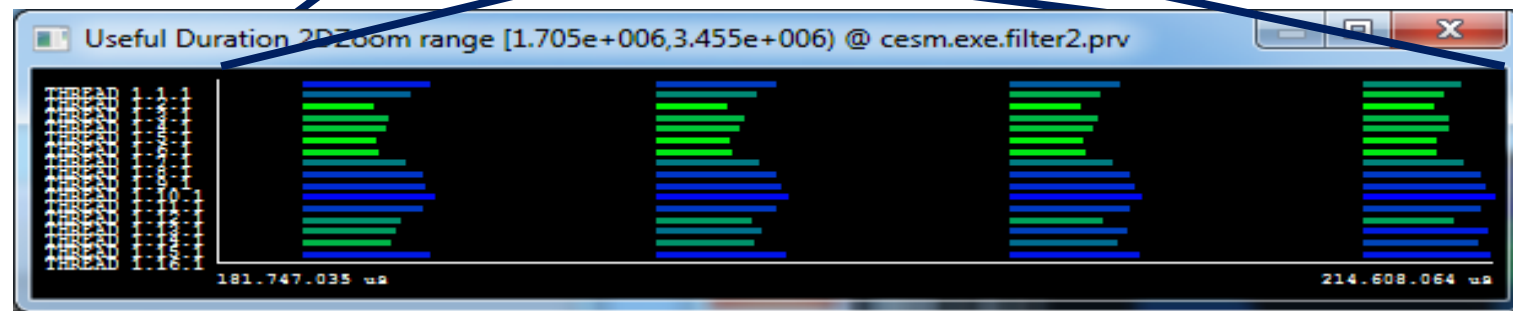# Analyzing variability through histograms and timelines

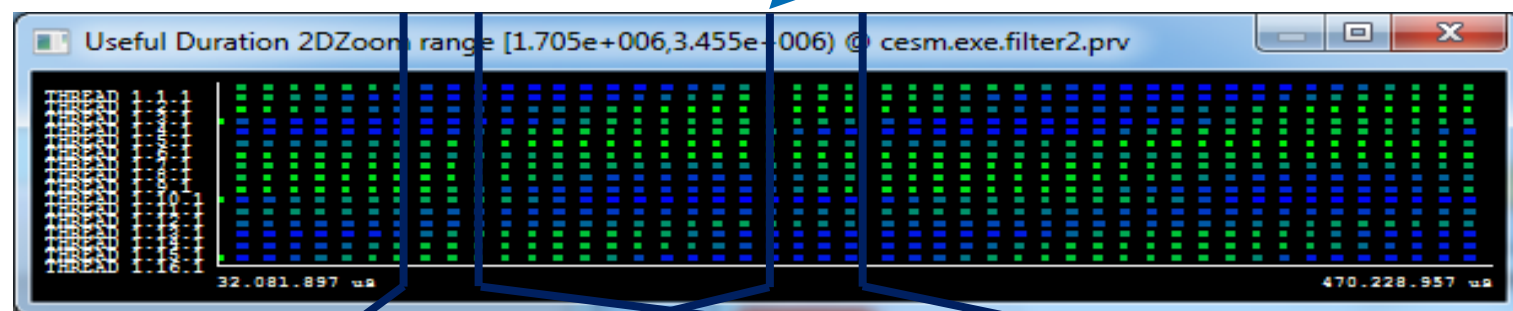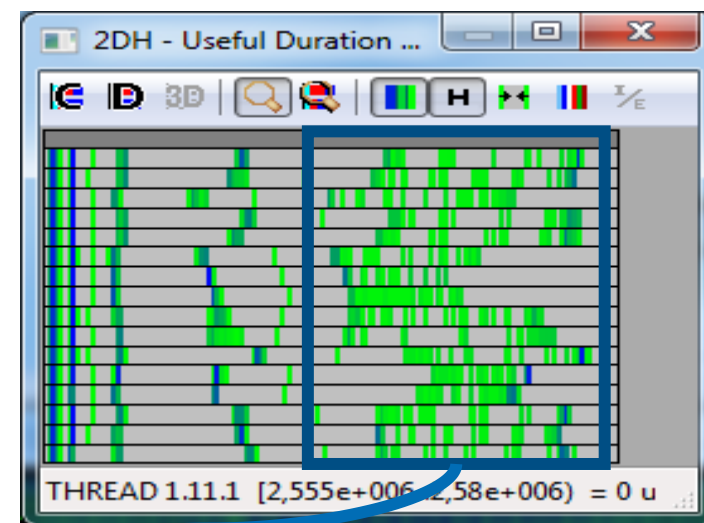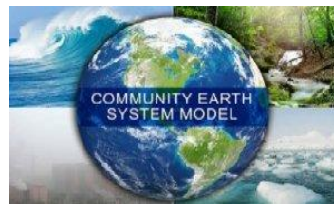- By the way: six months later…
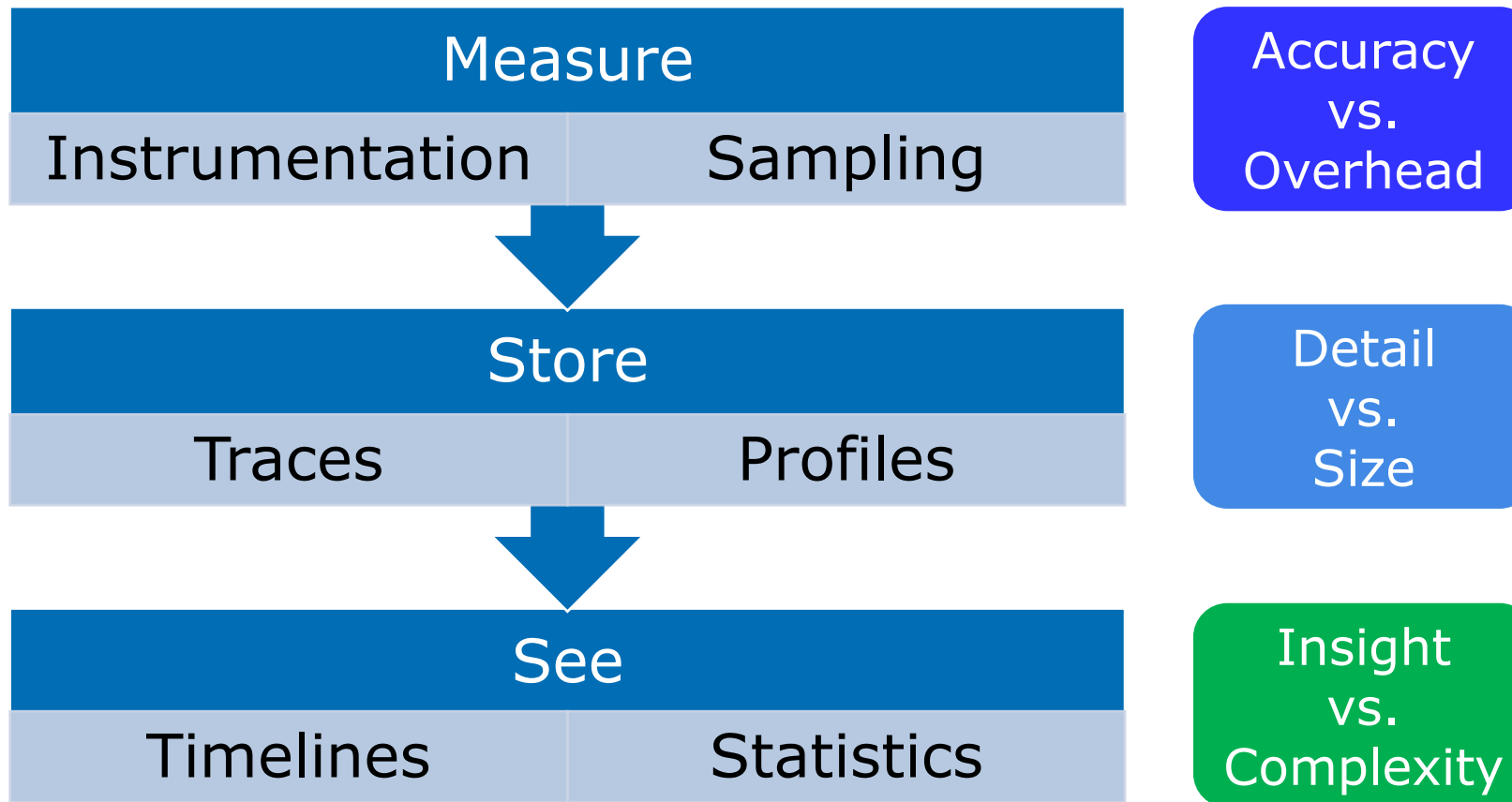


Focus optimization where it matters

# From tables to timelines

- CESM: 16 processes, 2-day simulation

- High variability in useful computation duration

- How is it distributed?

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

# Why traces?

# Performance data trade-offs

| Measure | |
|---|---|
| Instrumentation | Sampling |

**Accuracy vs. Overhead**

| Store | |
|---|---|
| Traces | Profiles |

**Detail vs. Size**

| See | |
|---|---|
| Timelines | Statistics |

**Insight vs. Complexity**

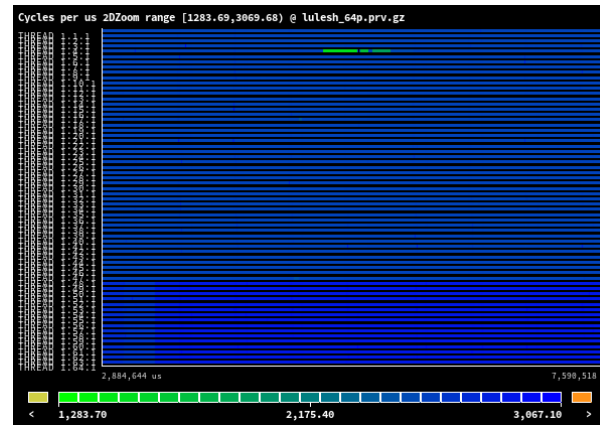**Highly detailed traces can be quite large...**

# Manipulating big traces

- Data processing & summarization

  - **Filtering**
    - Subset of records from the original trace
    - By duration, type, value…

  - **Cutting**
    - All records within a time interval
    - Selected processes only

  - **Software counters**
    - Aggregated metrics as new events
      - MPI call count, HW totals…

  - Remains a Paraver-compatible trace for analysis with the same CFGs (if needed data is kept)
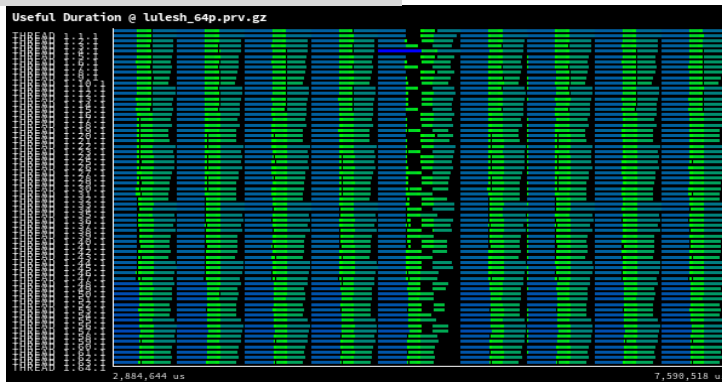
- Automatic analysis → **Performance Analytics**

# The Butterfly Effect...

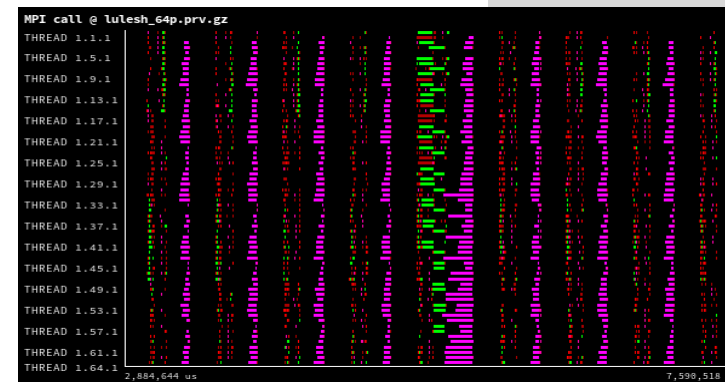▪ A system preemption reduces the cycles assigned to one of the processes for a small interval



**Cycles/us**

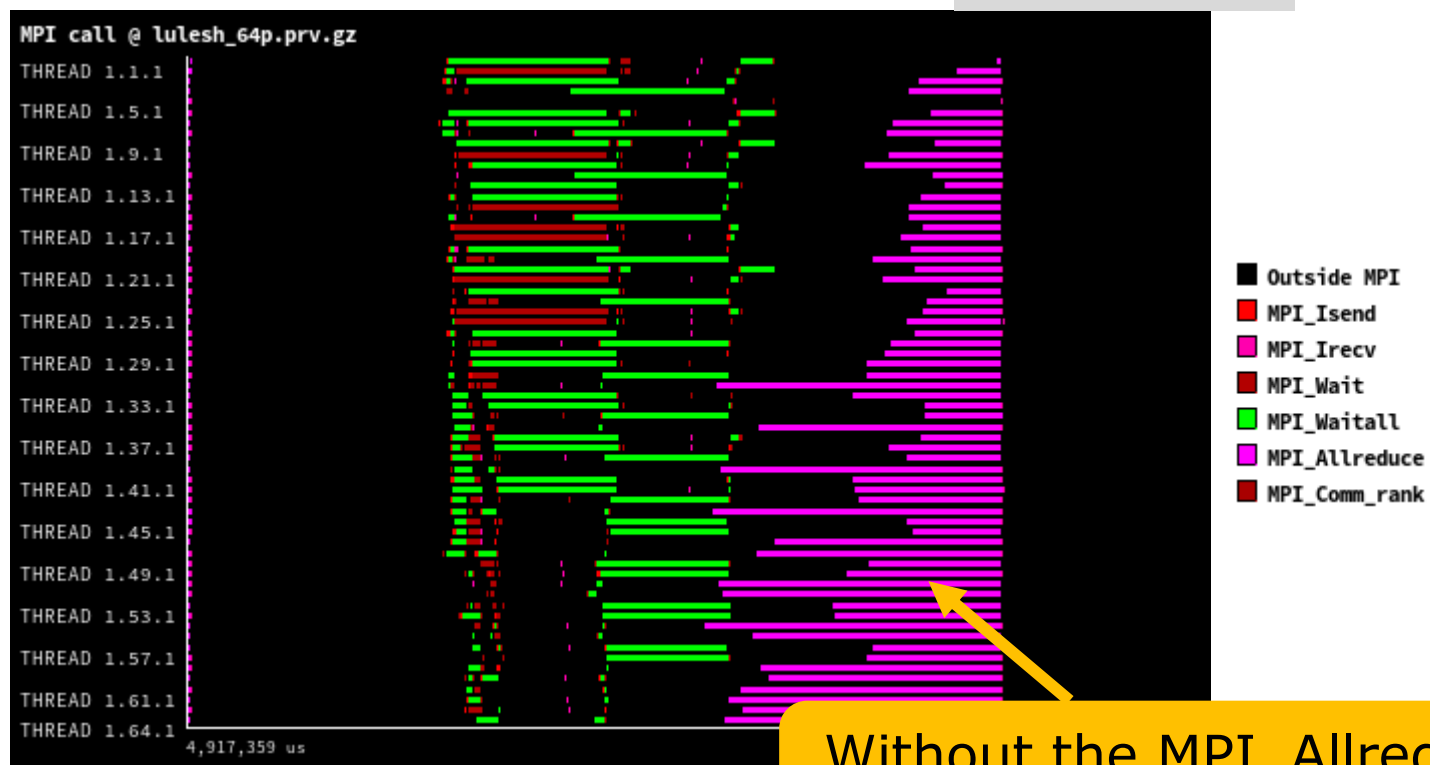**Useful duration**

**MPI calls**

**Affects only one process, but disturbs all processes**
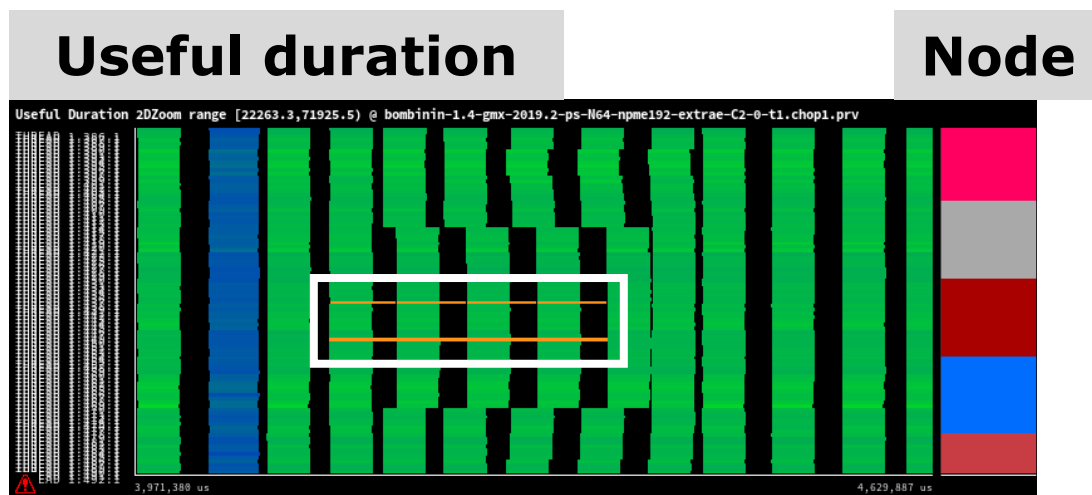
# ... flying through the MPI pattern

- Produces a wave effect that reflects and interferes across all partners
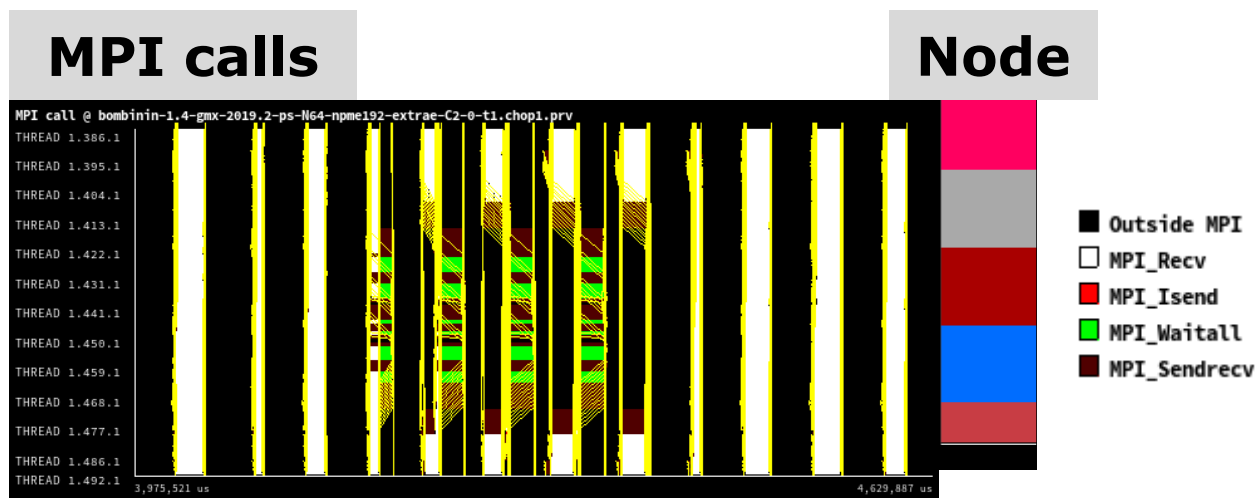
**MPI calls**



Without the MPI_Allreduce it would perturb multiple iterations
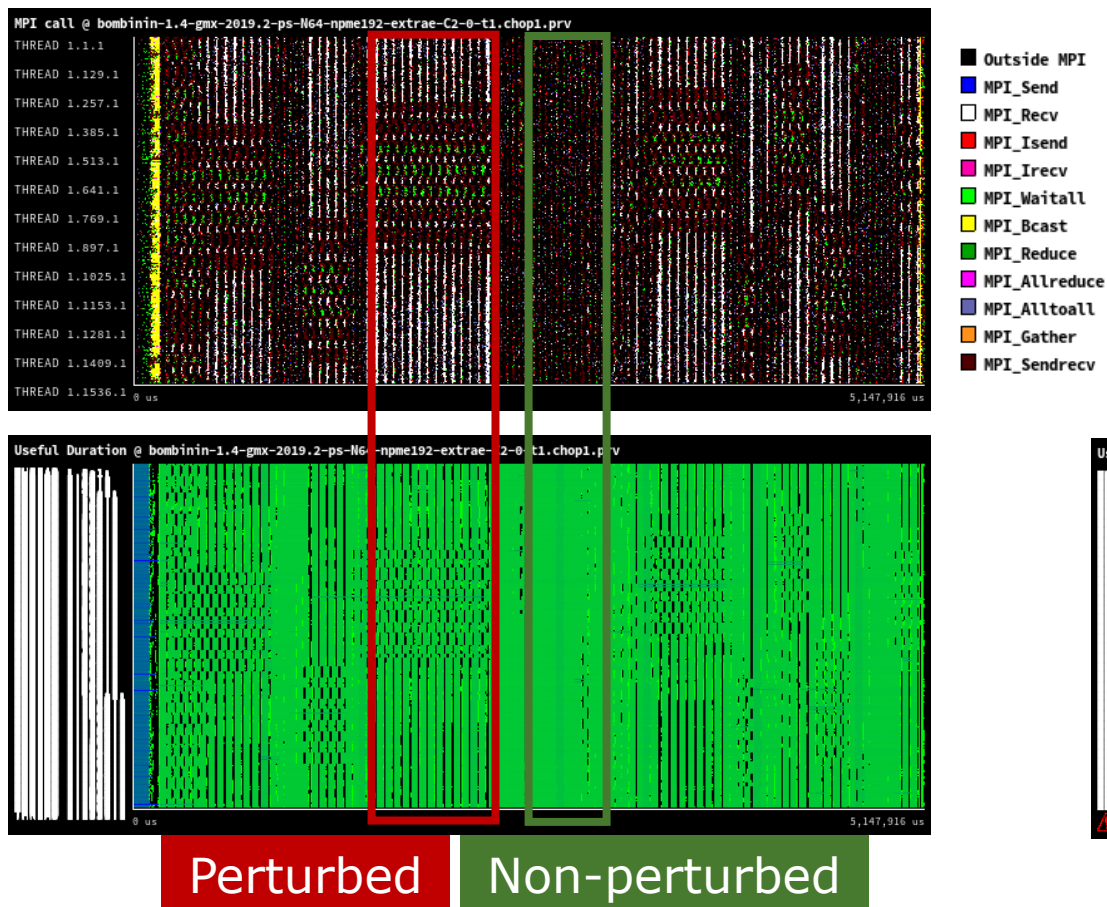
# … flying to other nodes

**Useful duration**

**Node**



- 2 processes perturbed in the same interval
- On the same node

**MPI calls**

**Node**



- ■ Outside MPI
- □ MPI_Recv
- ■ MPI_Isend
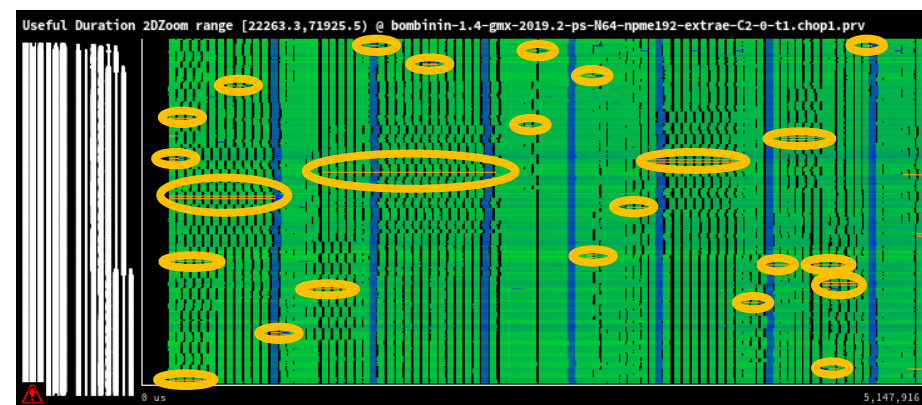- ■ MPI_Waitall
- ■ MPI_Sendrecv

**Affects only one node,
but disturbs many nodes**

# Thousands of butterflies?



**Microscopic effects with large global impact**

- Rescaling the gradient...

Perturbed    Non-perturbed
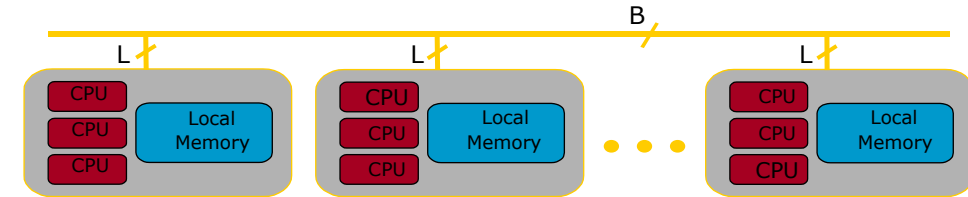
**Pinpoint bursts larger than expected**

# Dimemas

# Dimemas: Coarse-grain, Trace-driven simulation

- **Fast simulation of an abstract interconnect**
  - SMP nodes with local memory for intra-node comms
  - Interconnect defined by L (links) and B (buses)
    - B → Limits concurrent messages (contention)
    - L → Limits per-node traffic (connectivity)
  - Local/remote Latency/Bandwidth

- **Parametric sweeps**
  - On abstract architectures
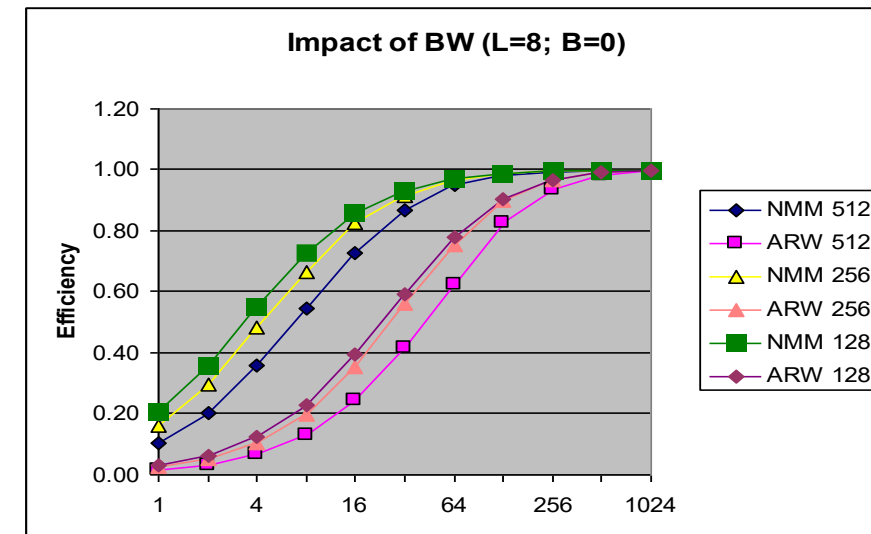  - On application computing regions

- **"What-if…" analysis**
  - Ideal machine (instantaneous network)
  - Would benefit from asynchronous communications?
  - Are all regions equally sensitive to the network?

- **MPI sanity check**
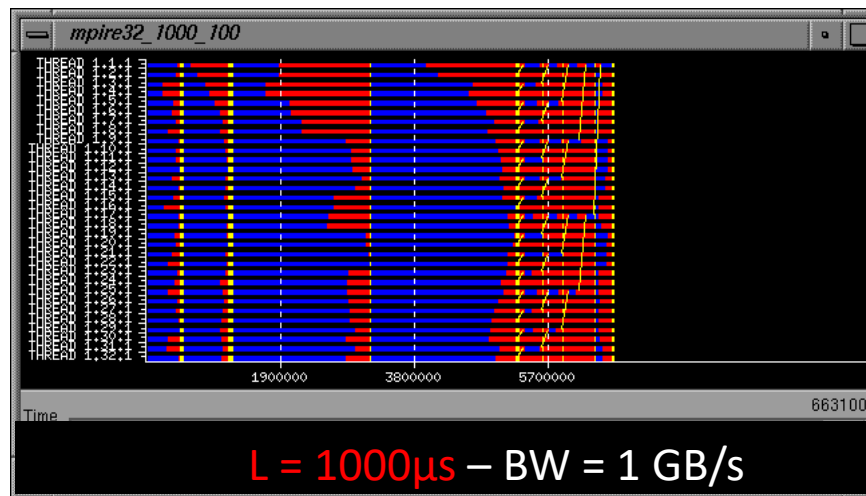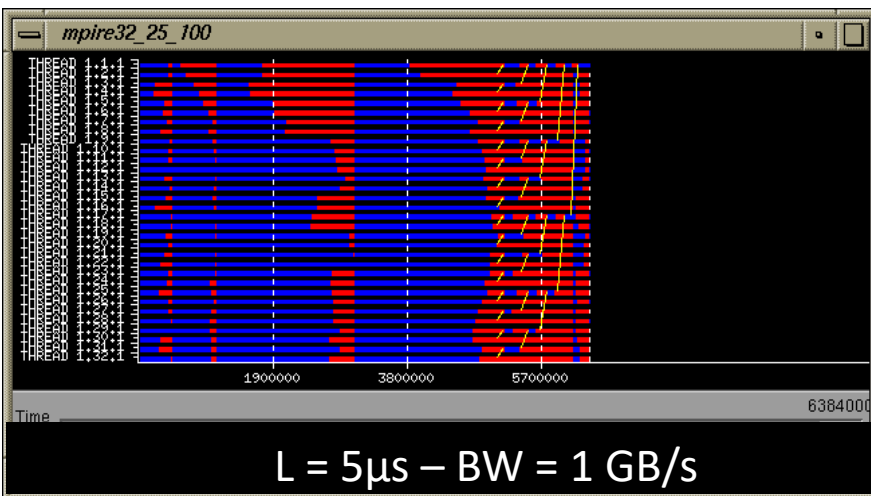  - Nominal modeling

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





Impact of BW (L=8; B=0)

Simulation generates a trace → Detailed feedback

# Network sensitivity

- MPIRE 32 tasks, no contention → Sensitive to BW or Latency?



L = 5µs – BW = 1 GB/s



L = 1000µs – BW = 1 GB/s

Higher latency doesn't hurt

Lower bandwidth kills the P2P phase



L = 5µs – BW = 100MB/s

# What if… we had asynchronous communications?



*Courtesy Dimitri Komatitsch*

Ideal simulation shows not benefit

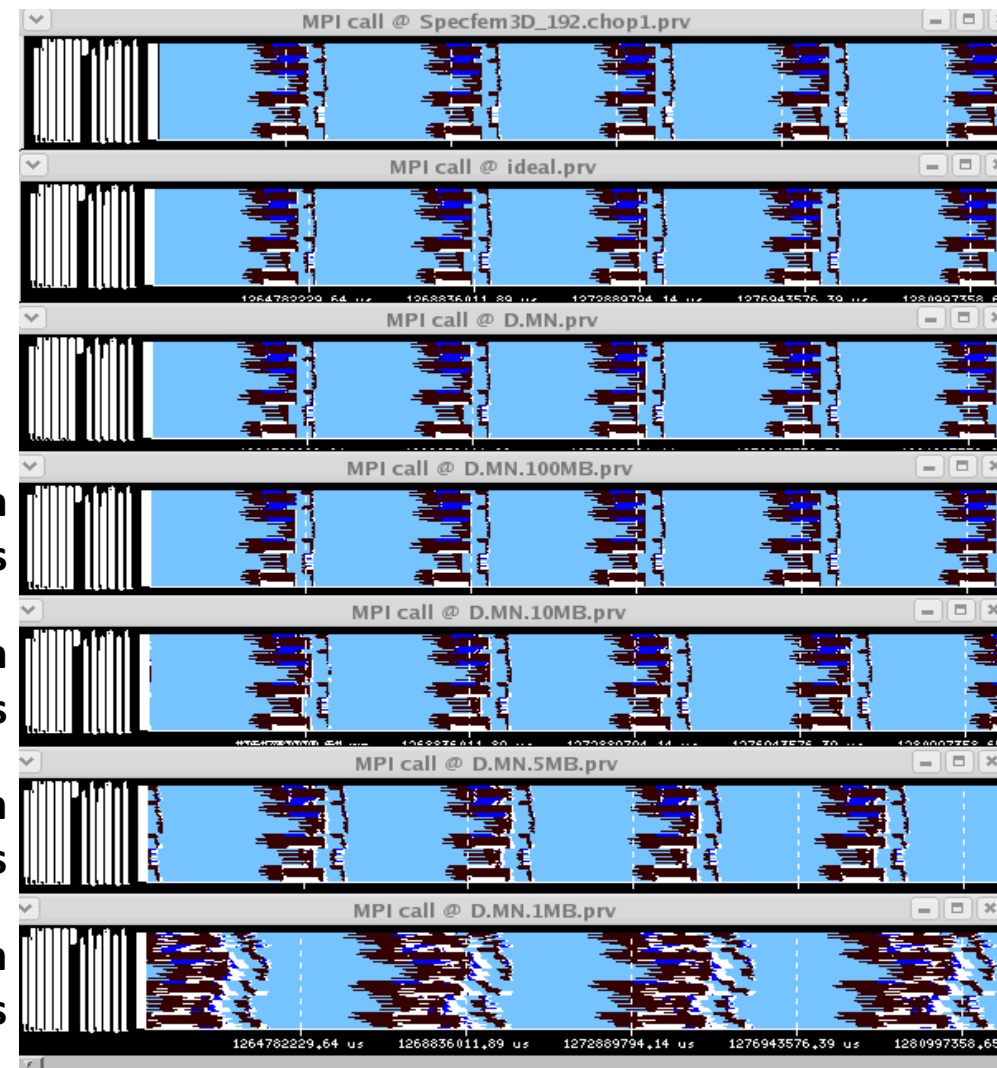Perturbation only appears below 10 MB/s

**Async comms won't help**

Real

Ideal

Prediction MN

Prediction 100MB/s

Prediction 10MB/s

Prediction 5MB/s

Prediction 1MB/s

# The Ideal Machine

- **BW = ∞, L = 0**
  - Data transfer would be instantaneous → MPI time should vanish. Why not?
  - Characterizes intrinsic application behavior
    - Load balance problems?
    - Dependency problems?

Allgather  Allreduce  Alltoall  Sendrecv  Waitall

Real run

GADGET @ Nehalem 256 processes

Ideal network

No data transfer issue, it's dependencies!

💡 Impact on real machines?
If it doesn't improve on the ideal, it won't on a faster real one

# Efficiency Model

# Efficiency Metrics



- Hierarchical model
- Multiplicative metrics
- 1 to 100% scale
- Two kinds of metrics:
  - Efficiency metrics (absolute)
  - Scalability metrics (relative to a base case)
- MPI, Hybrid MPI+OpenMP/CUDA

# Parallel Efficiency model

# Communication Efficiency refinement

**Computing**  **Communication**

Can't blame MPI for this!

MPI_Send     MPI_Recv

MPI_Recv     MPI_Send

$$LB=1$$

LB     *Ser*     *Trf*

$Comm_{ideal}$     $Comm/Ser$

MPI_Send     MPI_Recv

- Splitting Serialization (*Ser*) and Transfer (*Trf*)
  - Simulate with Dimemas an ideal network
    - Instantaneous data transfer ➜ ▇ vanishes (*Trf*=100%)

MPI_Recv     MPI_Send

Parallel Efficiency = LB * *Ser* * *Trf*

Inefficiencies from communication delays, including transfer

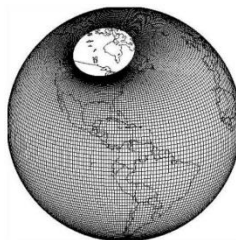Inefficiencies from circular dependences or non-uniform imbalances

Can be measured in Paraver with the real and simulated trace, and directly with BasicAnalysis
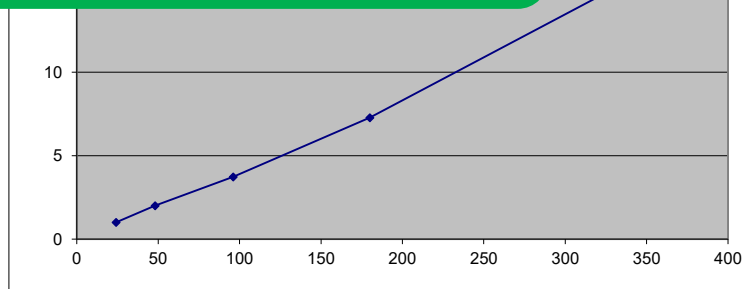
A low value means algorithmic dependency chains, varying load imbalance, or noise

# Why scaling?

- CGPOP ocean modelling
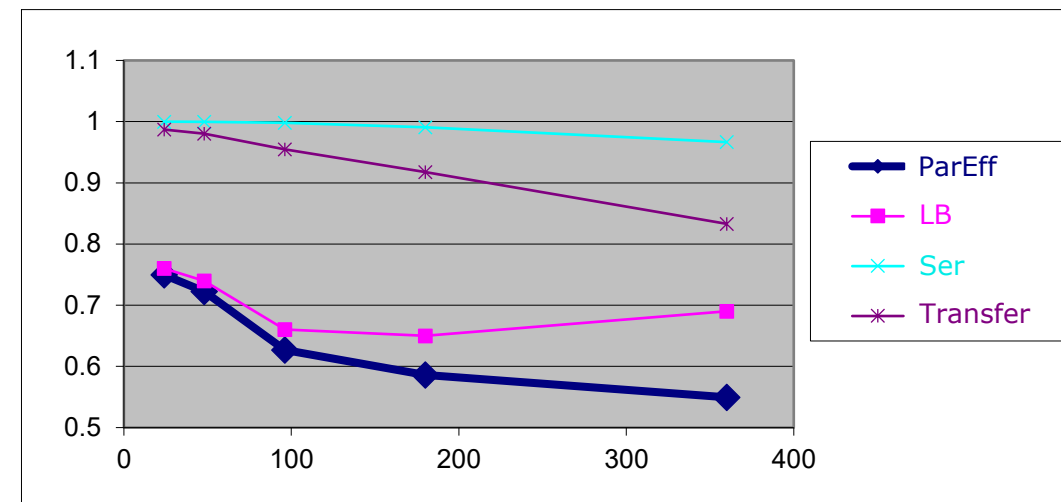  - Intrinsically unbalanced problem

**Or just lucky!?**

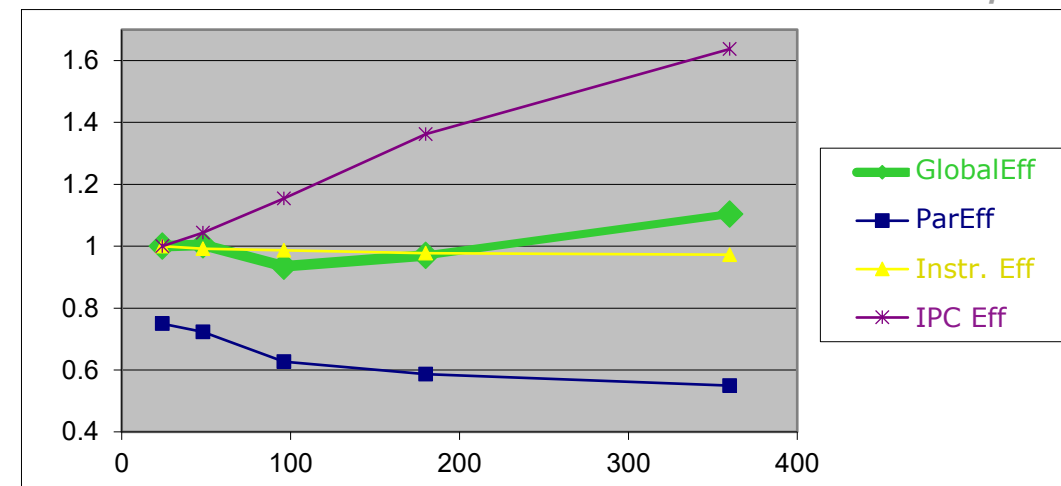**Good speed-up!
Are we happy?**

$$ParEff = LB * Ser * Trf$$



- Transfer Efficiency ↓
- IPC helps… for now!
- **Comms will become a bottleneck**
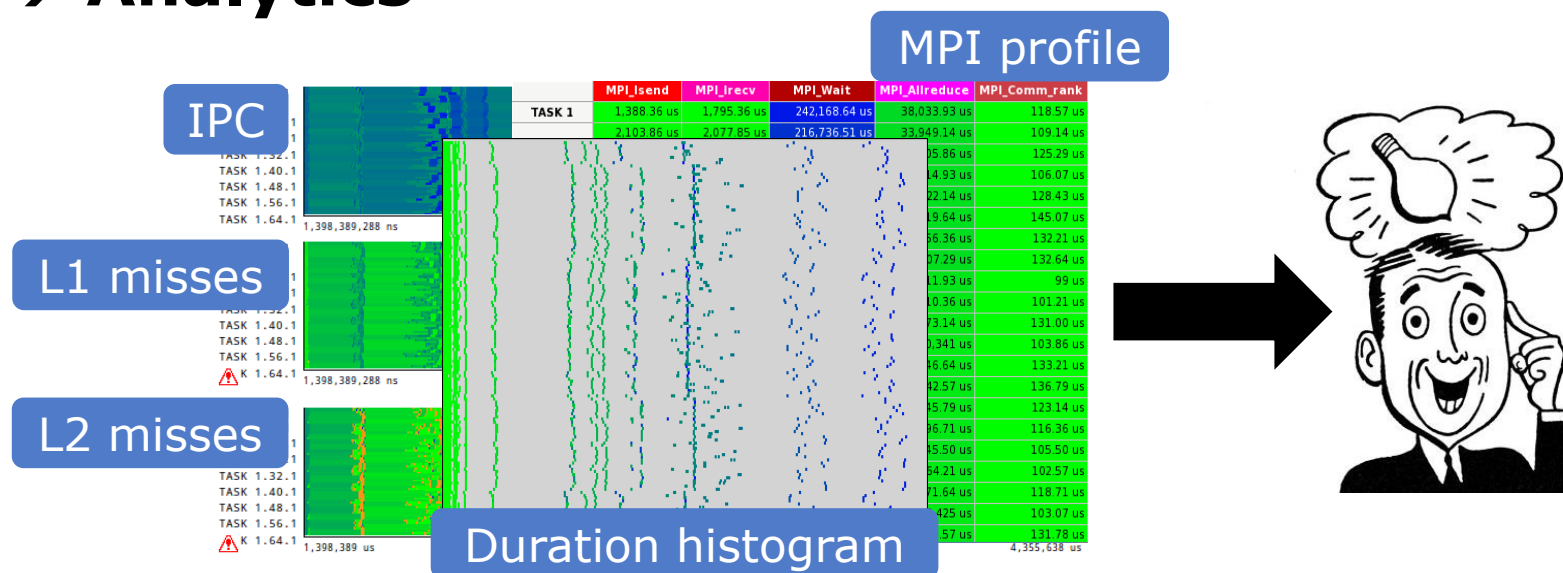
$$GlobalEff = ParEff * Instr * IPC * Freq$$

# Analytics

# Analysis → Analytics



- Dominant practice
  - Lots of data captured
    - But presentation goes from raw data to too general statistics

- Need for performance analytics
  - Leveraging techniques from data science, image processing, signal processing, etc.
  - Towards insight and models

# Basic Analysis

- Automatically compute the efficiency metrics from a Paraver trace (or many for scalability studies)
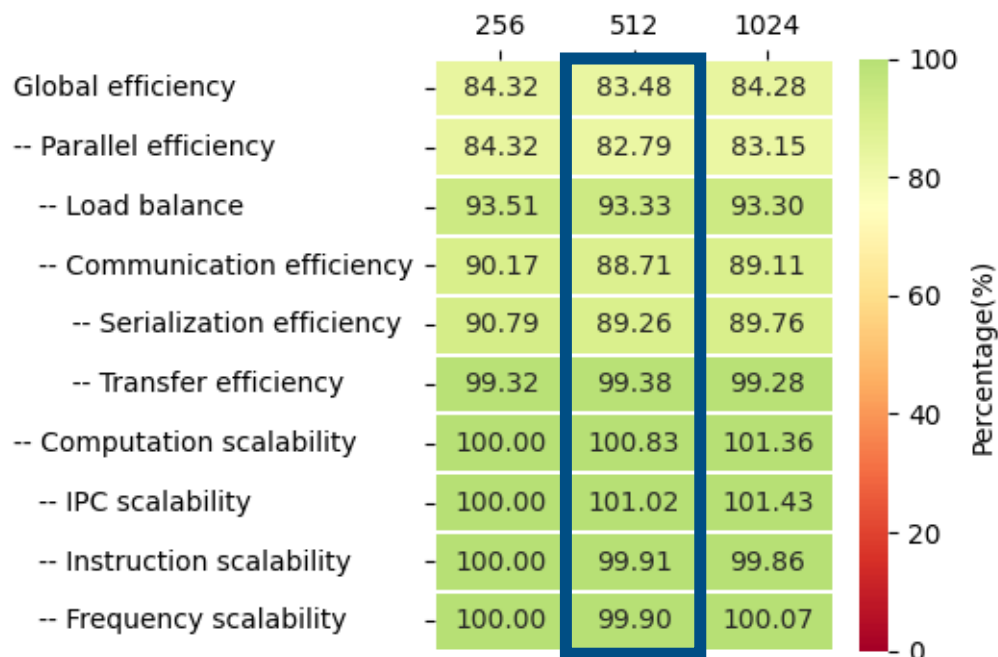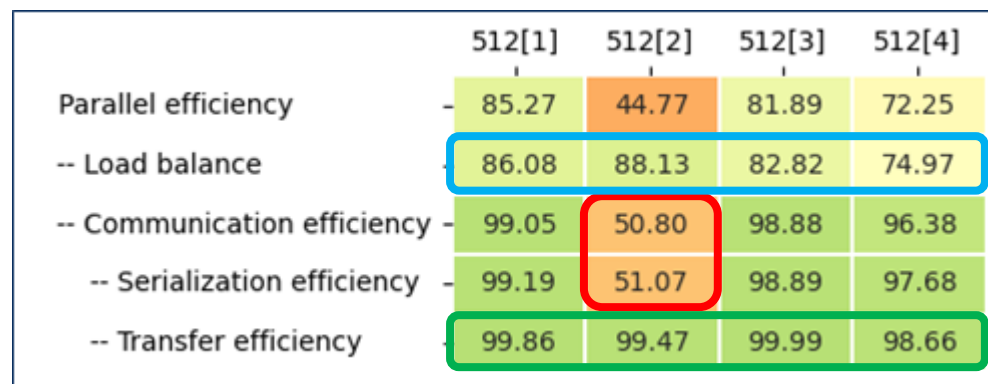- Dig down from global to detailed efficiencies

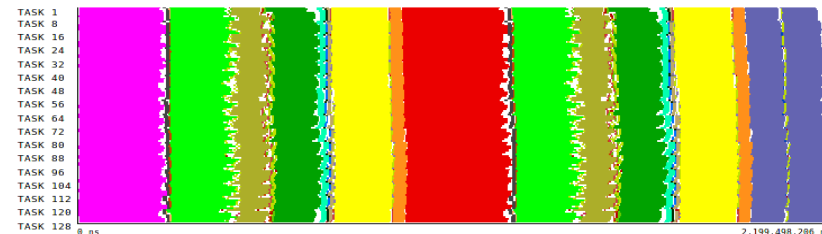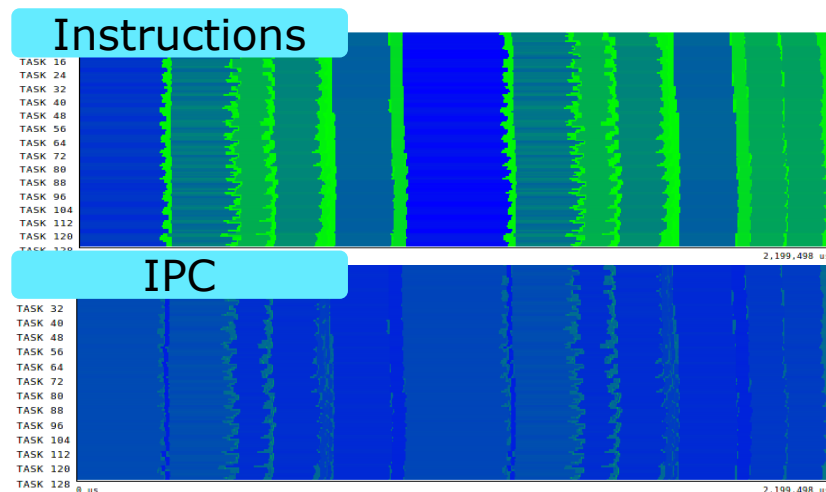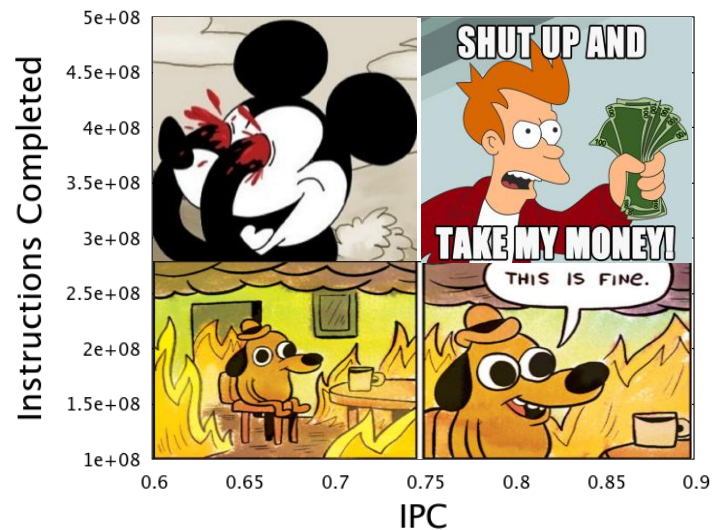**What to look for?**

Low values
Trends
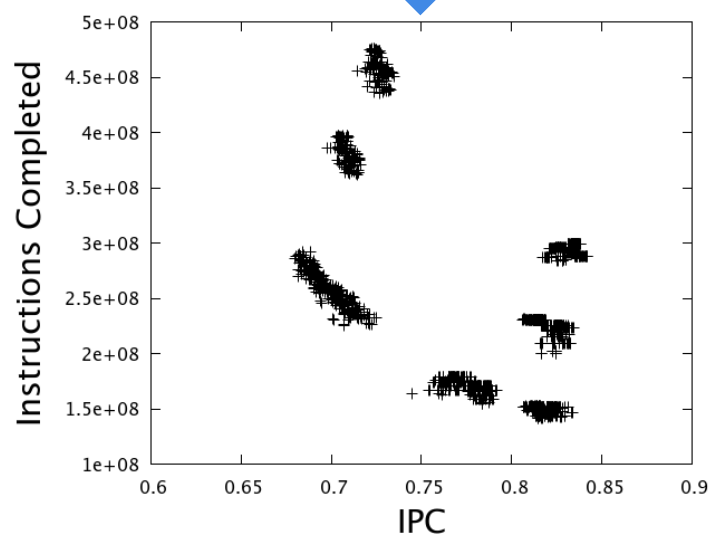High values



Comparing scales
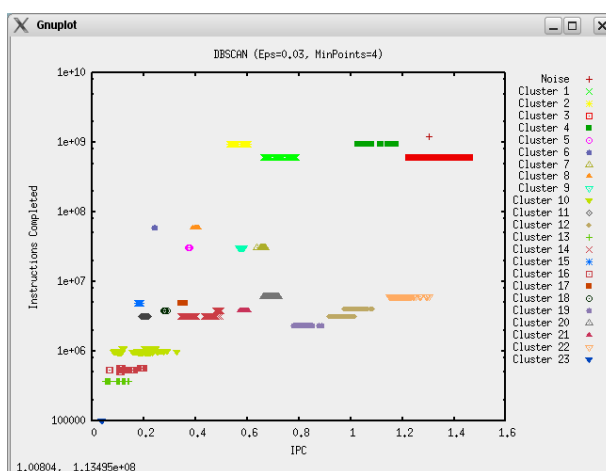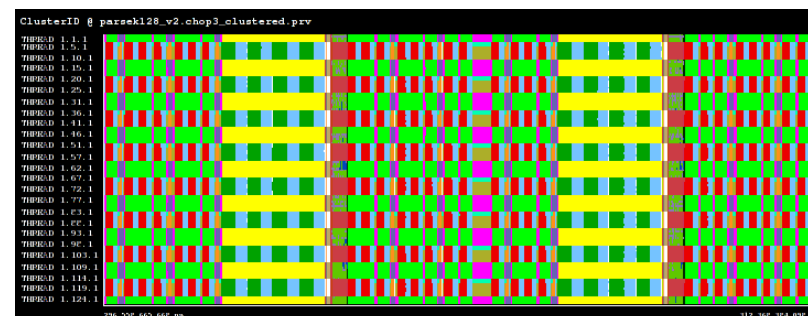


Comparing phases
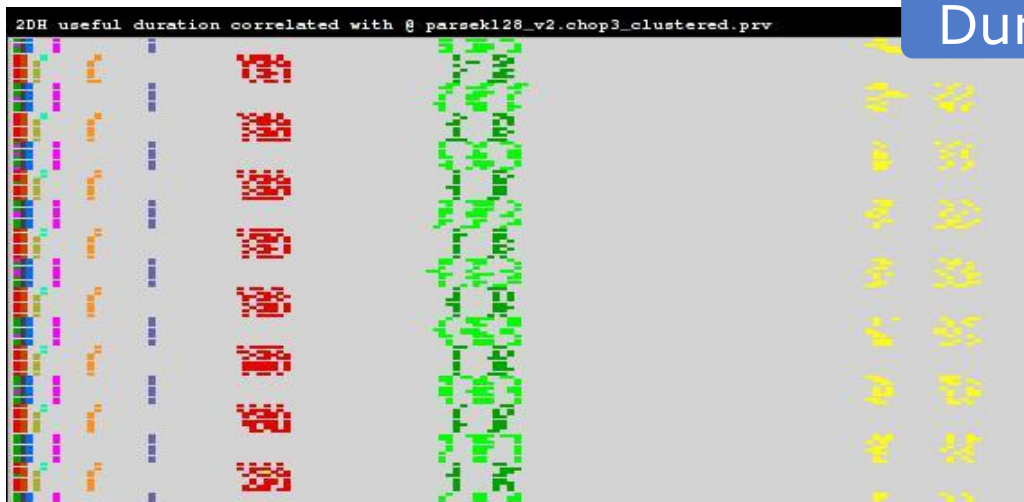
# Clustering to identify structure



Quick insight into program's behavior

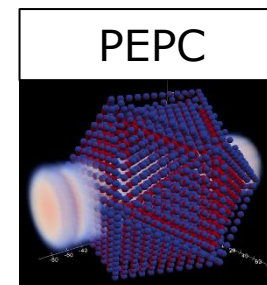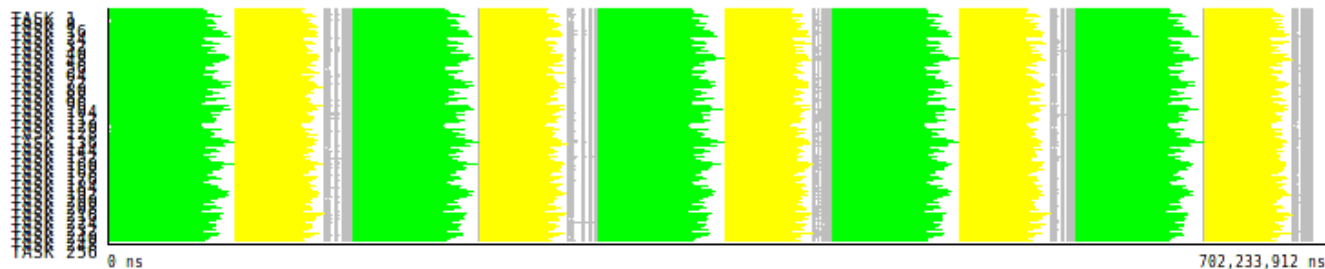# Correlate clusters, histograms & timelines



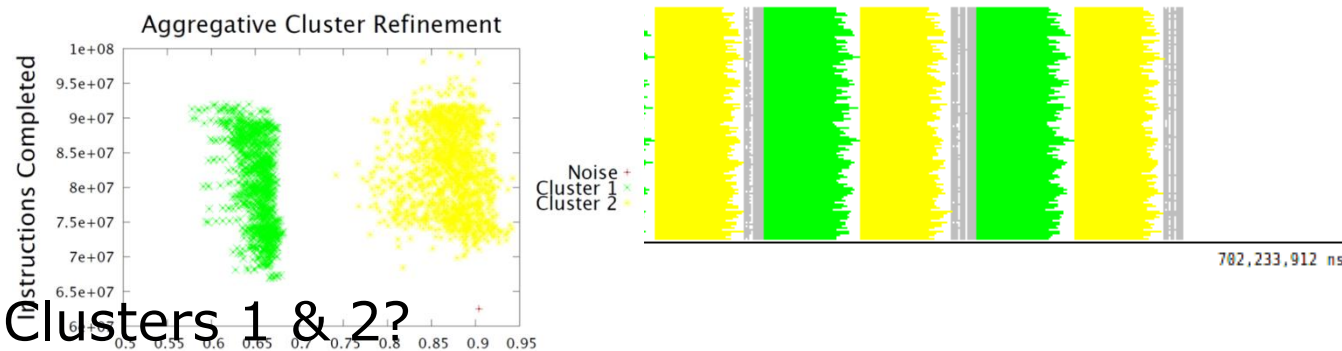Duration vs. Cluster

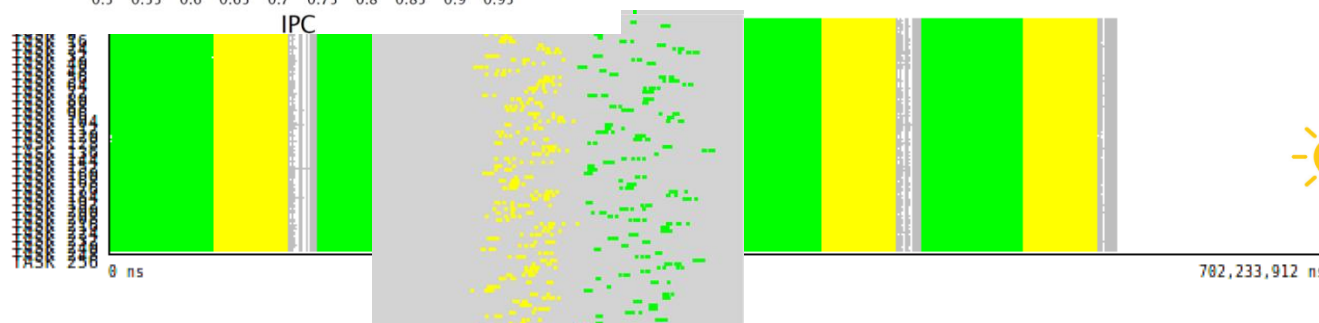Instructions vs. Cluster

# Integrating models and analytics

- What if…



PEPC

… we increase the IPC of Cluster1?



13%
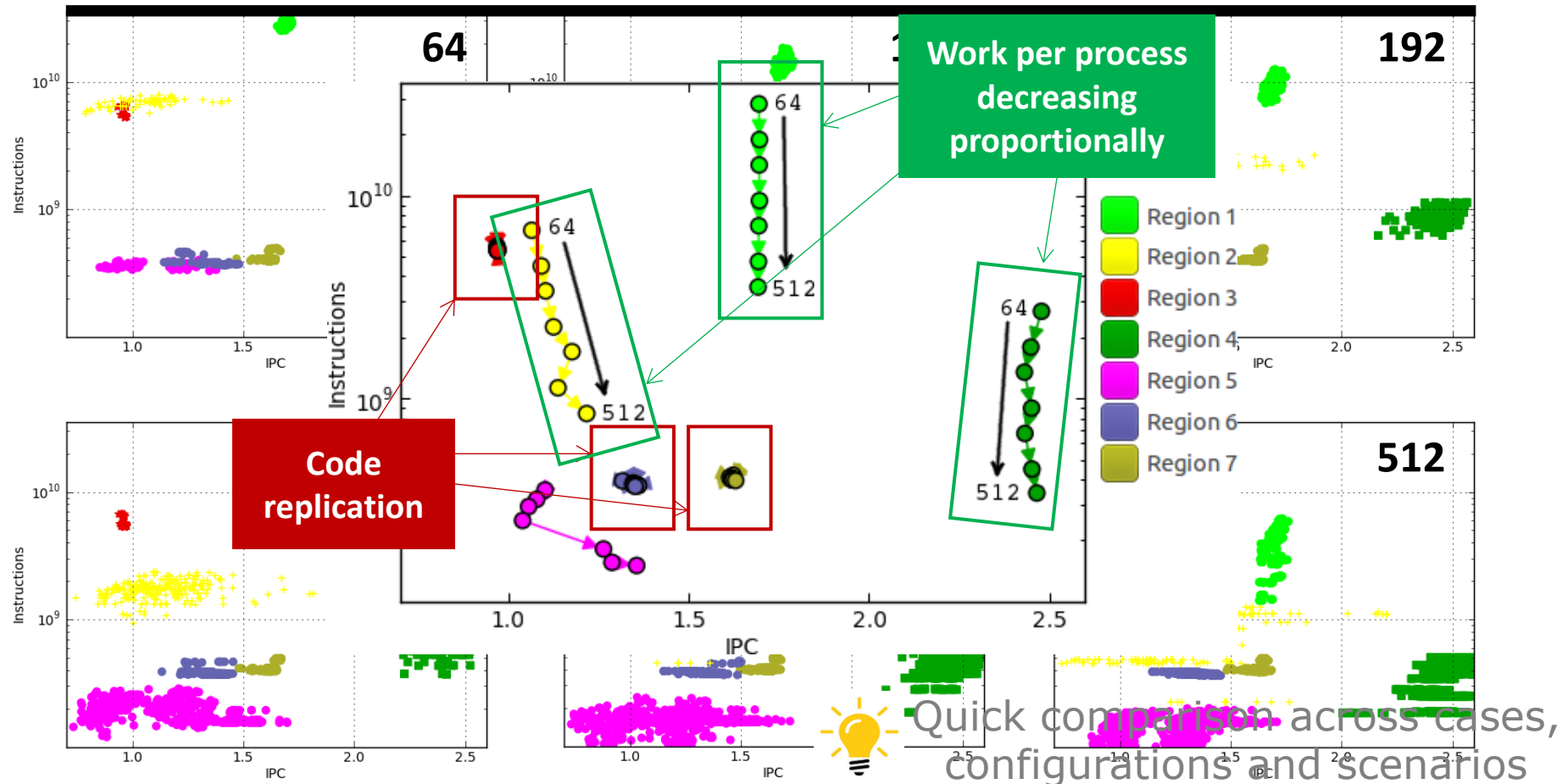
… we balance Clusters 1 & 2?
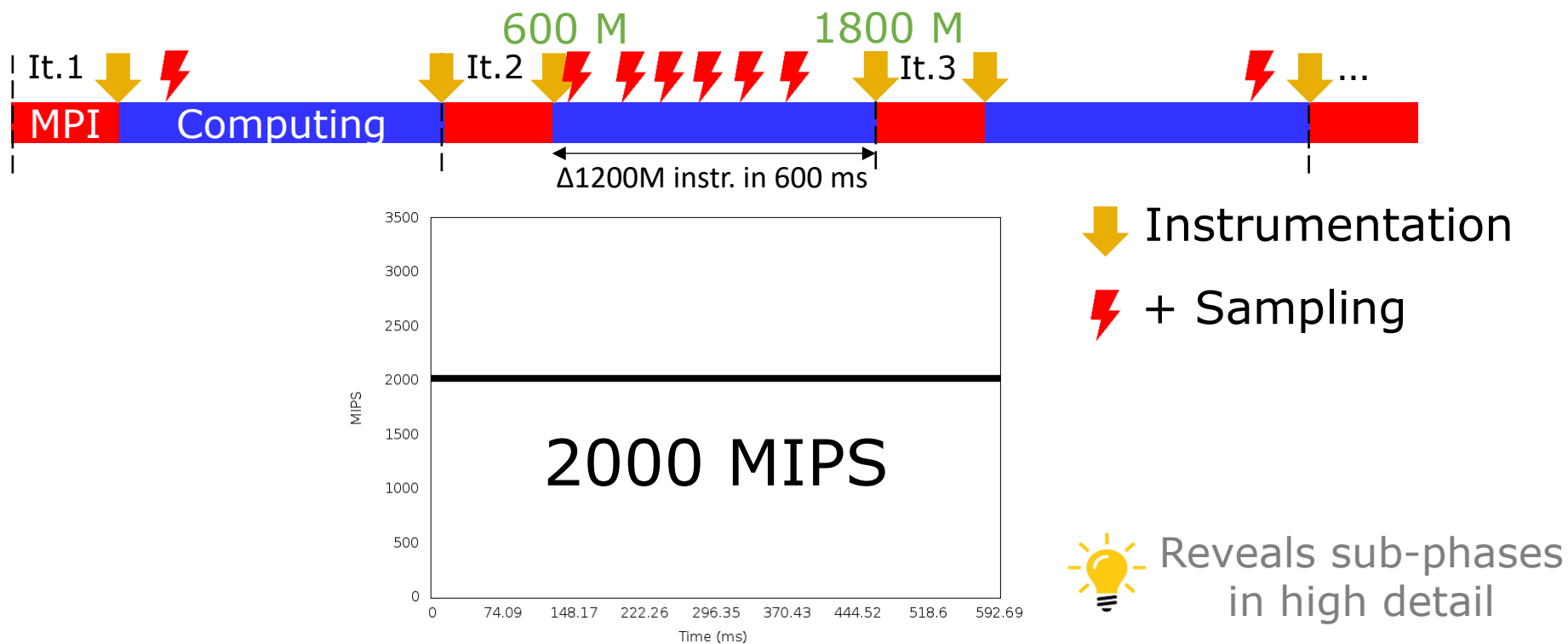


19%

Know where effort pays off.

# Tracking scalability through clustering

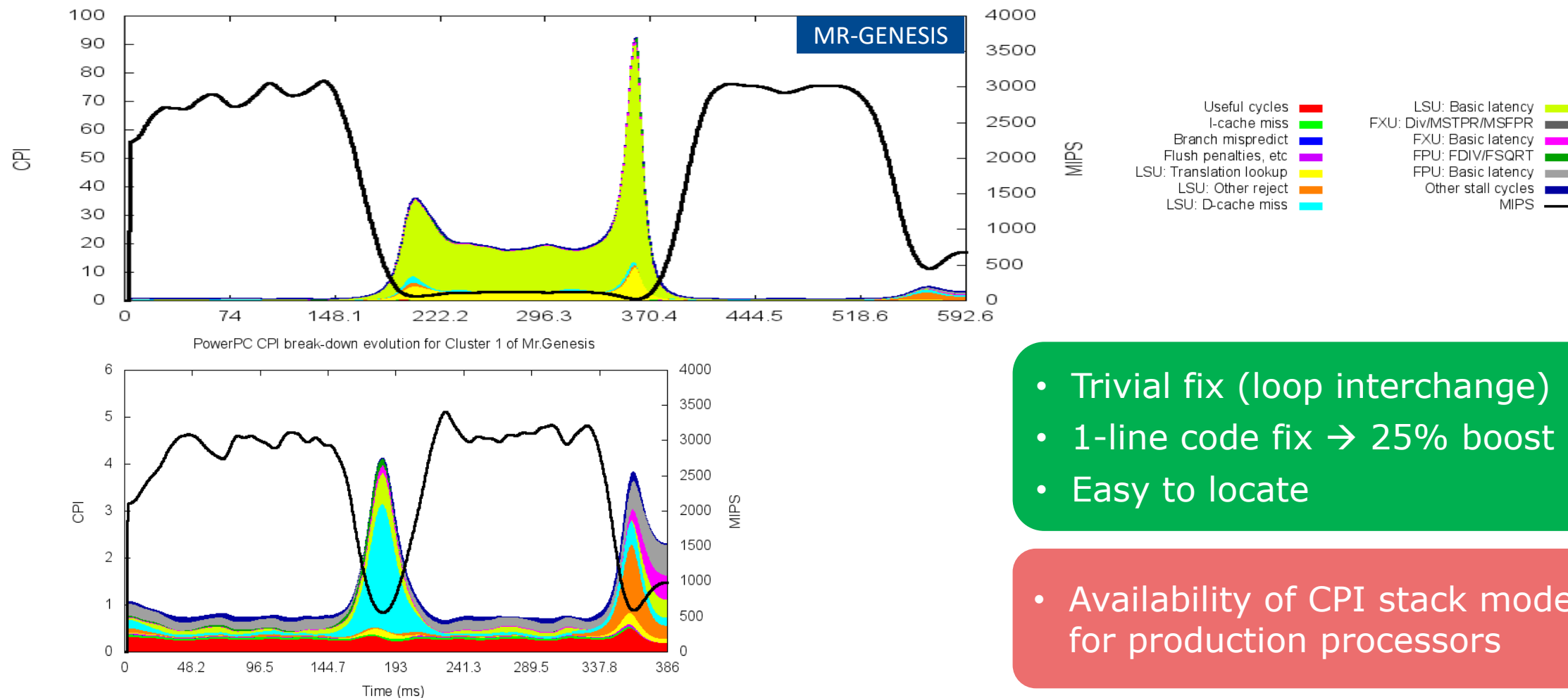▪ Analyze scalability of computing regions across 64 – 512 tasks

# Folding to increase details

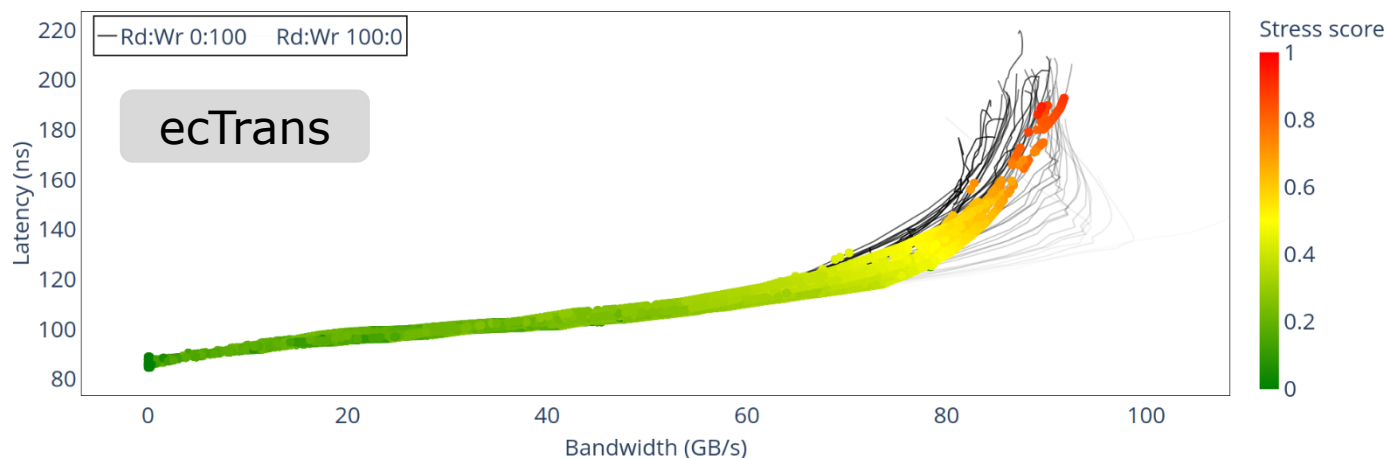▪ What is the performance of a single serial region?



**Instrumentation**

**+ Sampling**

Reveals sub-phases in high detail

"Fold" similar iterations into a single, highly detailed synthetic iteration

# Folding: CPI and HWC stack models



MR-GENESIS

PowerPC CPI break-down evolution for Cluster 1 of Mr.Genesis

Legend:
- Useful cycles
- I-cache miss
- Branch mispredict
- Flush penalties, etc
- LSU: Translation lookup
- LSU: Other reject
- LSU: D-cache miss
- LSU: Basic latency
- FXU: Div/MSTPR/MSFPR
- FXU: Basic latency
- FPU: FDIV/FSQRT
- FPU: Basic latency
- Other stall cycles
- MIPS

- Trivial fix (loop interchange)
- 1-line code fix → 25% boost
- Easy to locate

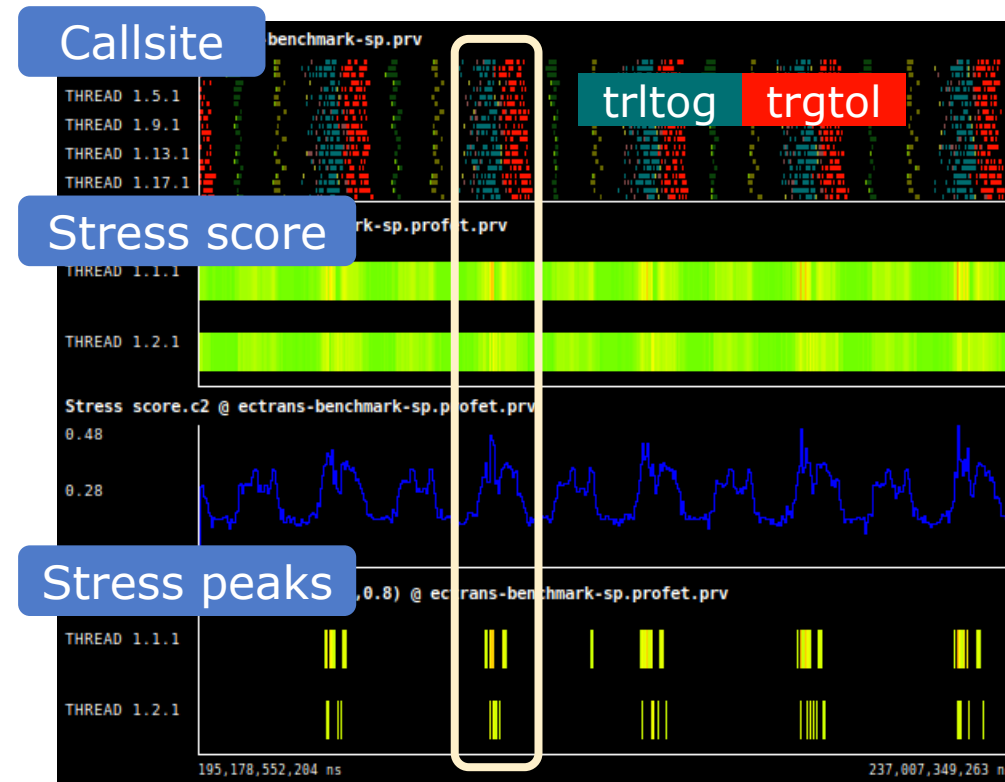- Availability of CPI stack models for production processors

# Understanding memory influence

- Mess: Bandwidth-latency curves describe memory performance from unloaded to fully saturated states
- Integration with Paraver: Easily identify where memory stress is highest
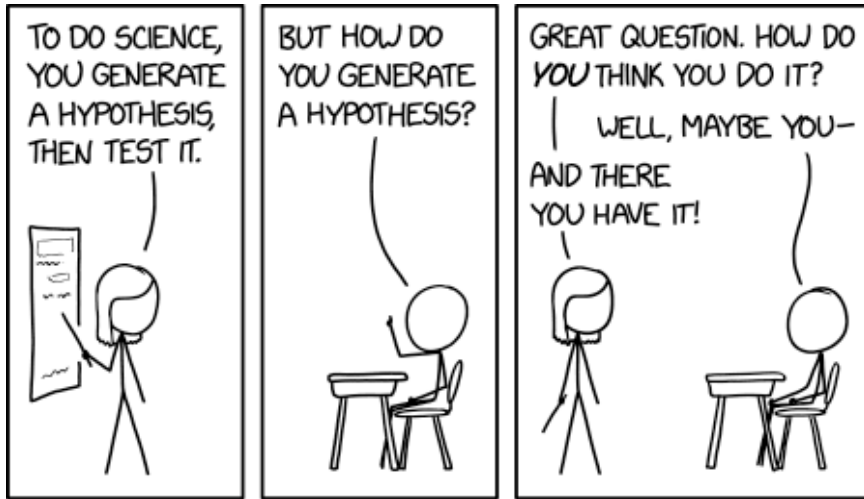


- Prediction of future ones (PROFET)
- Cache-Aware Roofline Model (CARM)

# Methodology

# BSC Performance Tools aim



- Be your copilot in the process of generating & validating hypotheses
  - The tool is the "wheel" and you "drive it"

- Provide quantitative (how much) + qualitative (why/how)
  - Patterns, shapes, structure… beyond raw numbers "for a better understanding"

# First steps of analysis

| Identify Structure | → | Select FoA | → | Efficiency Metrics | → | Detailed Analyis | → | Conclusions & Suggestions |

- **Parallel Efficiency: Parallel resources are mostly doing useful work?**
  - Load Balance → Work (programmer's fault)? Performance (machine's fault)?
  - Serialization → Dependency chains? Sequence of MPI calls? Noise?
  - Transfer → Bandwidth or Latency? Simulations with Dimemas

- **Serial Efficiency: How far from peak performance?**
  - IPC → Cache effects? Correlate with other counters
  - Frequency → Cycles per us? Multi-core sharing? Dynamic scaling? Power?

- **Scalability: Benefit from additional resources?**
  - Code replication → Total instructions?

- **Variability?**

- **Behavioral Structure → Analytics**

Paraver Tutorials:
Introduction to Paraver &
Dimemas methodology

# BSC Tools Website & Contact

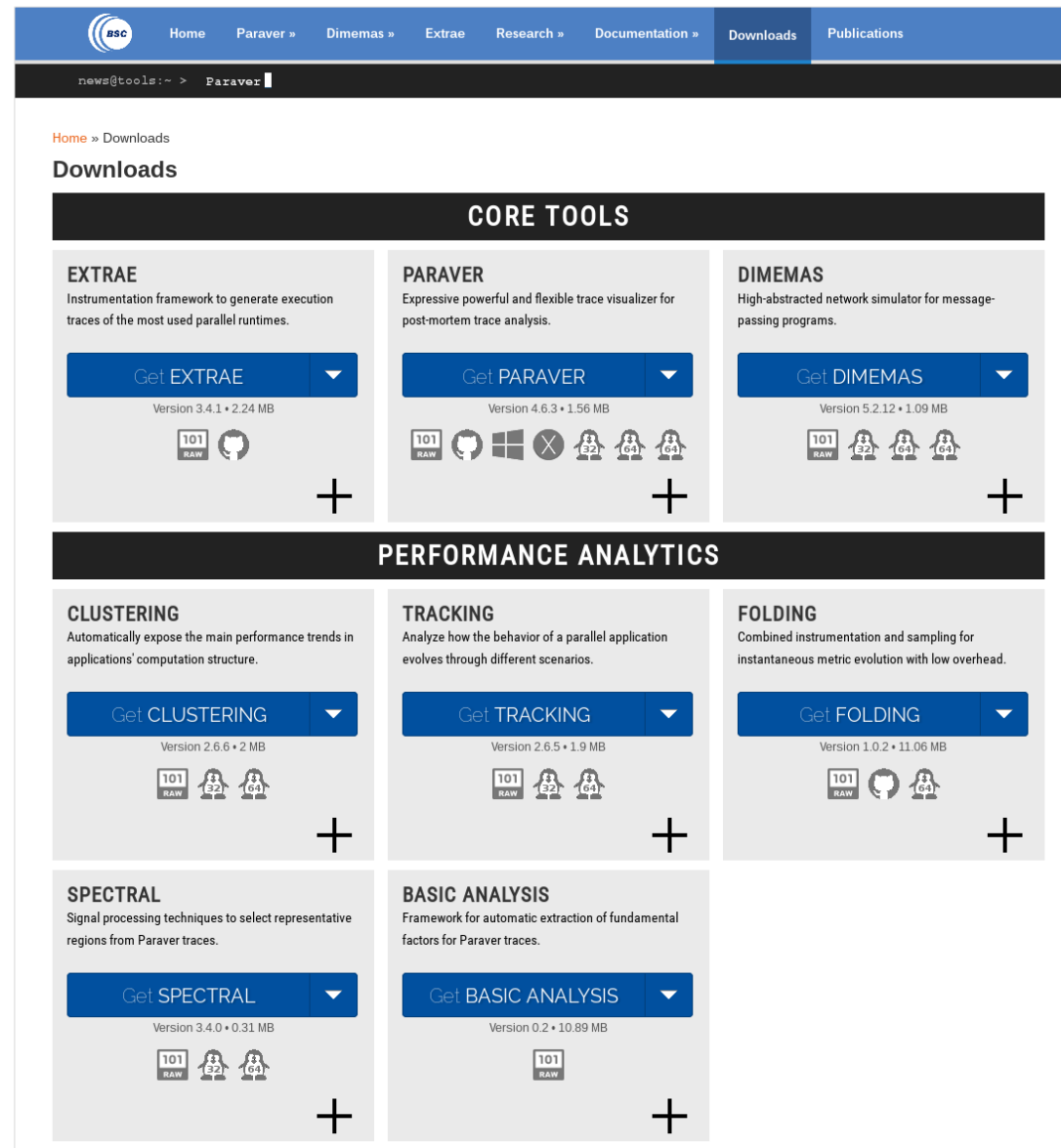- ## https://tools.bsc.es

## Open Source

## Downloads
- Sources & Binaries

## Documentation
- Training guides
- Tutorial slides

## tools@bsc.es

## Quick Start
- Start wxparaver
- Help → Tutorials
- Follow training guides

- 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**

Takeaway

# Performance Optimisation and Productivity 3
## A Centre of Excellence in HPC

- Free performance assessment services for European users

- Request your assessment here

Contact:

🌐    https://www.pop-coe.eu

✉    pop@bsc.es

▶    youtube.com/POPHPC