# Introduction to BSC-Tools

Germán Llort

gllort@bsc.es

05/11/2021

DEEP-SEA Seminar

# BSC Tools

- Since 1991

- Based on traces

- Open source – https://tools.bsc.es

- Focus: Detail, variability, flexibility

- Core tools
  - Extrae – Instrumentation
  - Paraver – Offline trace analysis
  - Dimemas – Message-passing simulator

- Performance Analytics
  - Leveraging techniques from data analytics
  - Towards insight and models

# Extrae: Flexible instrumentation…

- 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 anywhere in your program (Extrae API)

**No need to recompile nor relink!**

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# ... of unmodified binaries...

- Symbol substitution through LD_PRELOAD

```
export LD_PRELOAD=$EXTRAE_HOME/lib/libmpitrace.so
```

- Specific libraries for each runtime and combinations
  - MPI
  - OpenMP
  - OpenMP+MPI
  - ...

**Recommended**

- Dynamic instrumentation
  - Based on Dyninst (developed by U.Wisconsin / U.Maryland)
    - Instrumentation in memory
    - Binary rewriting
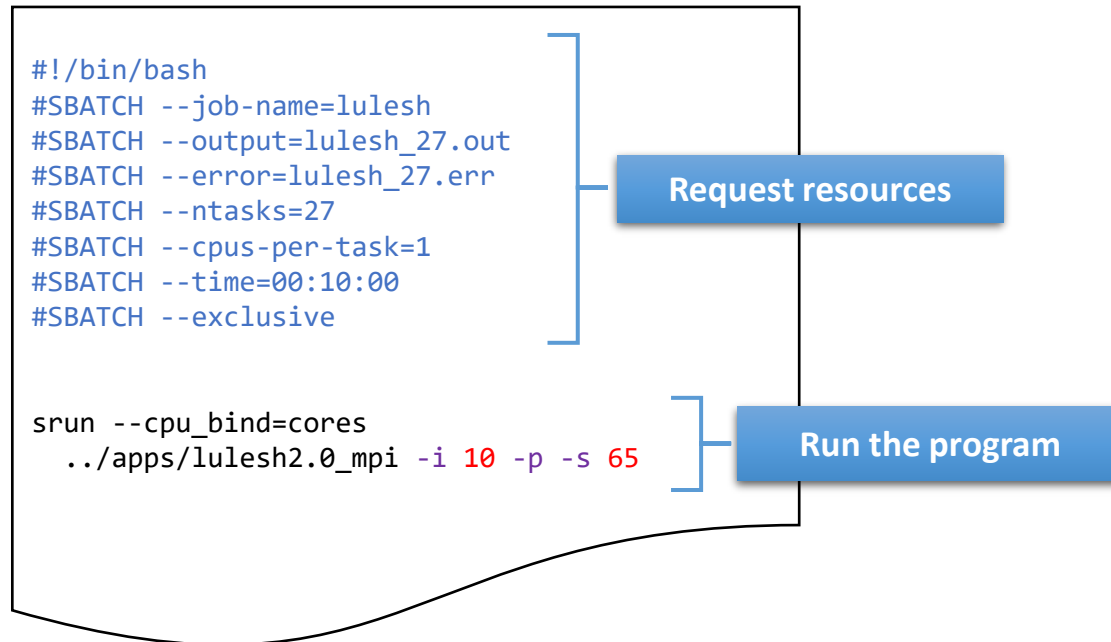
- Static link (i.e., PMPI, Extrae API)

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# … with low overhead

| | MN4 |
|---|---|
| Punctual event | 166 ns |
| Event + PAPI counters | 751 ns |
| Event + 1-level callstack | 2.875 us |
| Event + 6-levels callstack | 6.109 us |

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# Extrae: Quick start guide

- Sample SLURM jobscript

```
#!/bin/bash
#SBATCH --job-name=lulesh
#SBATCH --output=lulesh_27.out
#SBATCH --error=lulesh_27.err
#SBATCH --ntasks=27
#SBATCH --cpus-per-task=1
#SBATCH --time=00:10:00
#SBATCH --exclusive


srun --cpu_bind=cores
    ../apps/lulesh2.0_mpi -i 10 -p -s 65
```

**Request resources**

**Run the program**

# Extrae: Quick start guide (II)

```bash
#!/bin/bash
#SBATCH --job-name=lulesh
#SBATCH --output=lulesh_27.out
#SBATCH --error=lulesh_27.err
#SBATCH --ntasks=27
#SBATCH --cpus-per-task=1
#SBATCH --time=00:10:00
#SBATCH --exclusive


srun --cpu_bind=cores ./trace.sh
   ../apps/lulesh2.0_mpi -i 10 -p -s 65
```

```bash
#!/bin/bash

# Configure Extrae
export EXTRAE_CONFIG_FILE=./extrae.xml
export EXTRAE_HOME=<path-to-installation>

# Load the tracing library (choose RT/C/Fortran)
export LD_PRELOAD=$EXTRAE_HOME/lib/libmpitrace.so

# Run the program
$*
```

**What to trace?**

**Choose tracing library**

libseqtrace
libmpitrace[f]
libomptrace
libpttrace
libcudatrace
libompitrace[f]
libptmpitrace[f]
libcudampitrace[f]
…

# Understanding applications with Paraver

- Timeline views
  - Categorical data – color encoding



MPI calls
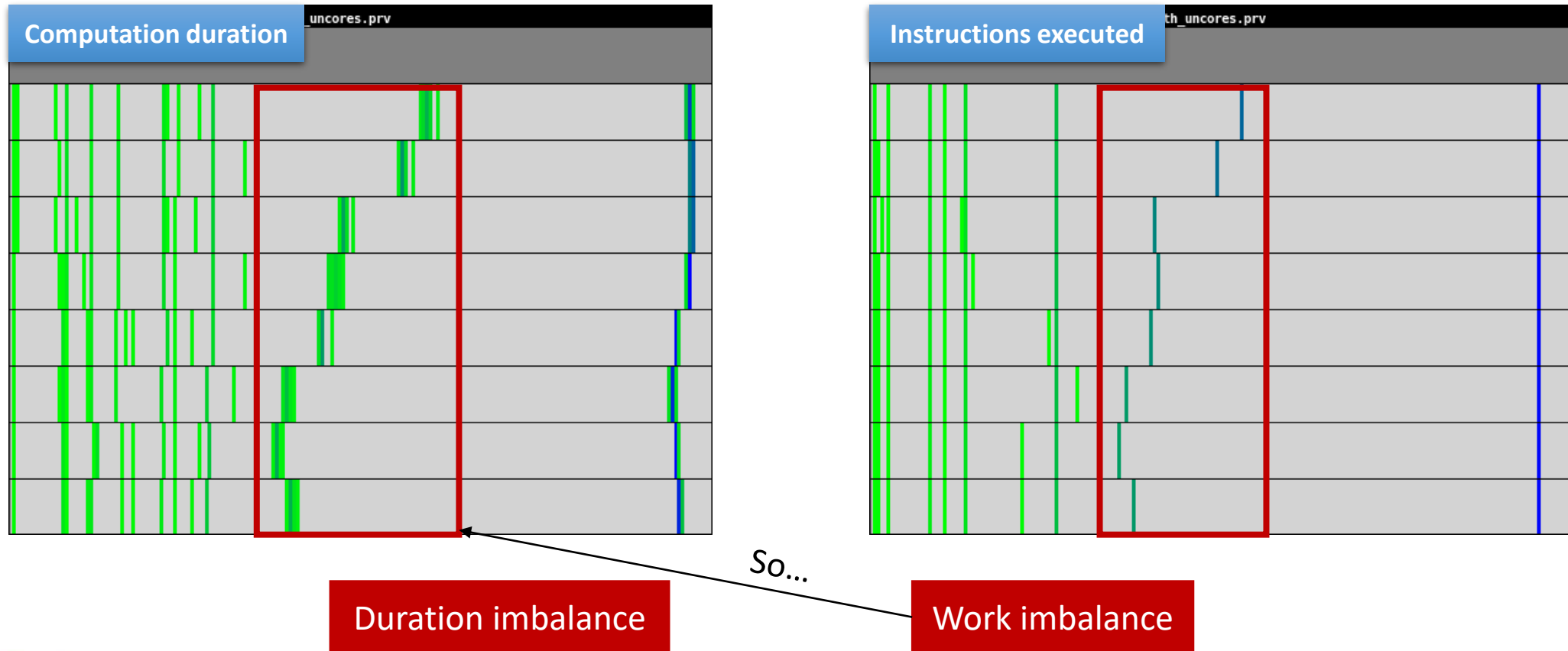
  - Continuous data – gradient (green to blue)



Computation duration
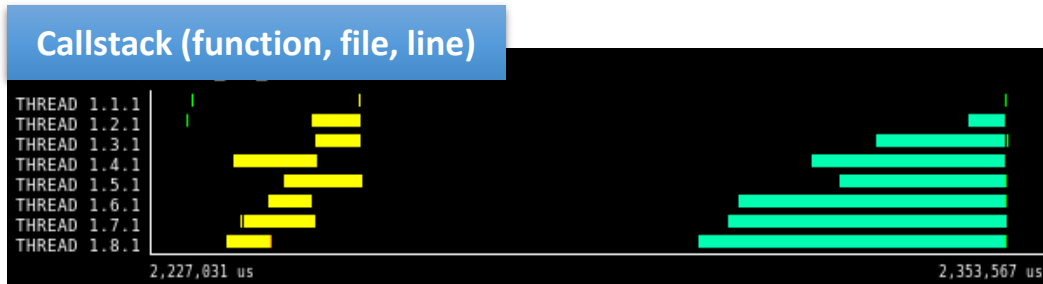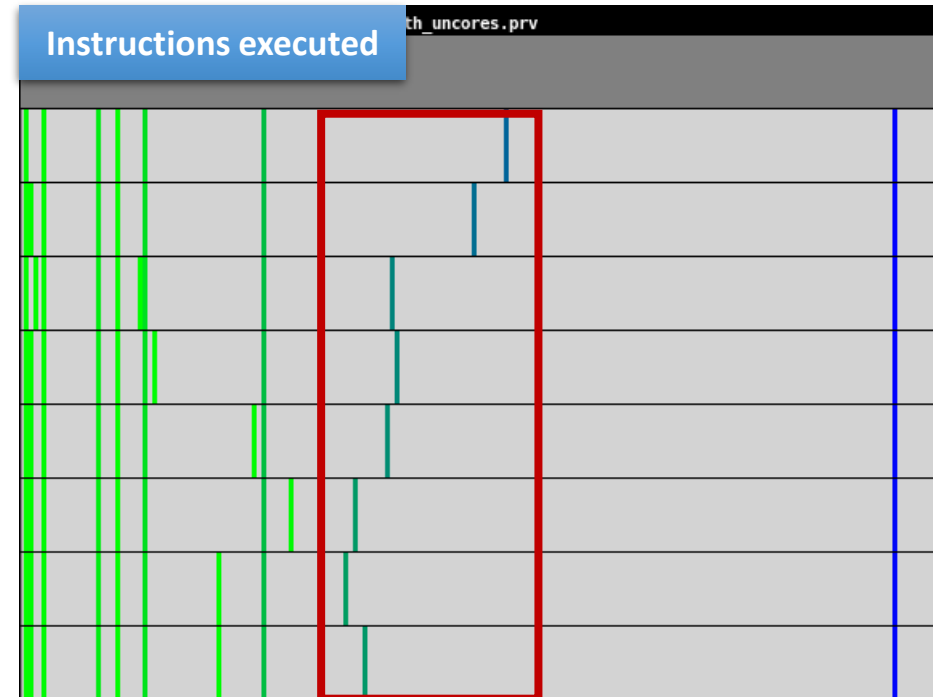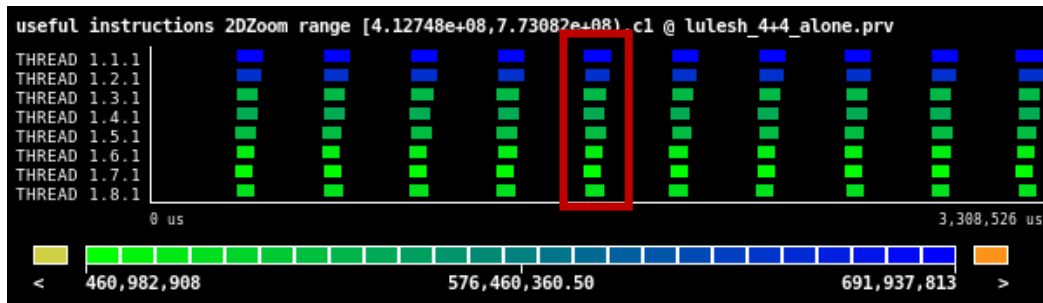
- From timelines to tables



MPI profile



Duration histogram

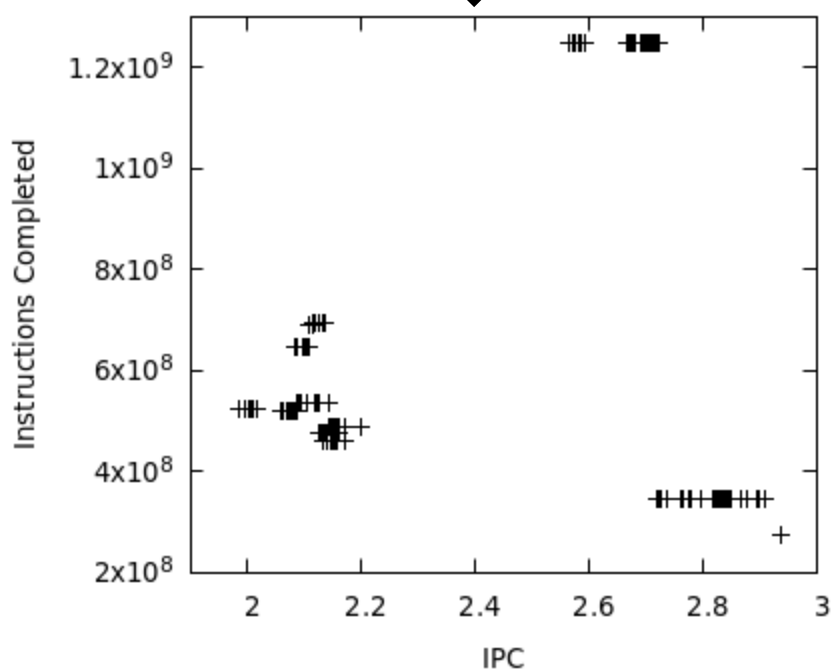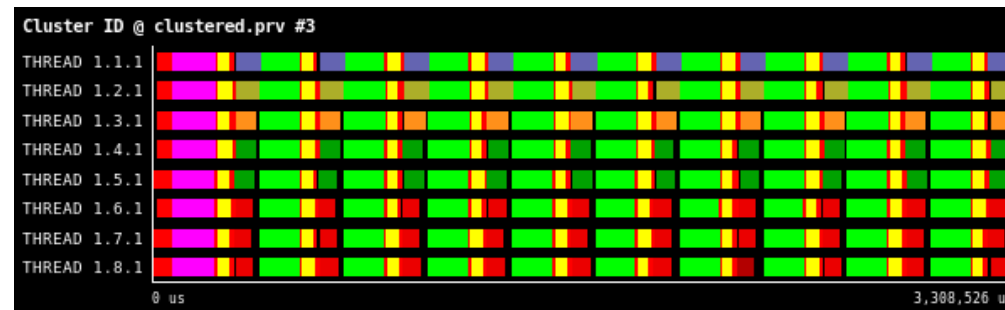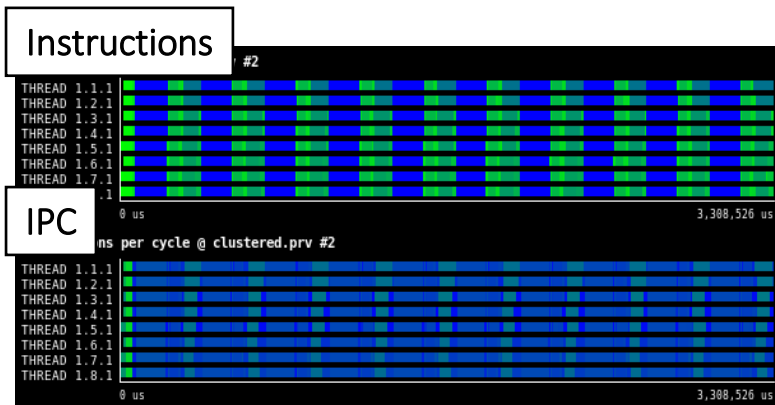# Understanding applications with Paraver (II)

- Correlating multiple views



**Computation duration** — uncores.prv

**Instructions executed** — th_uncores.prv

So...

**Duration imbalance**

**Work imbalance**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Understanding applications with Paraver (III)
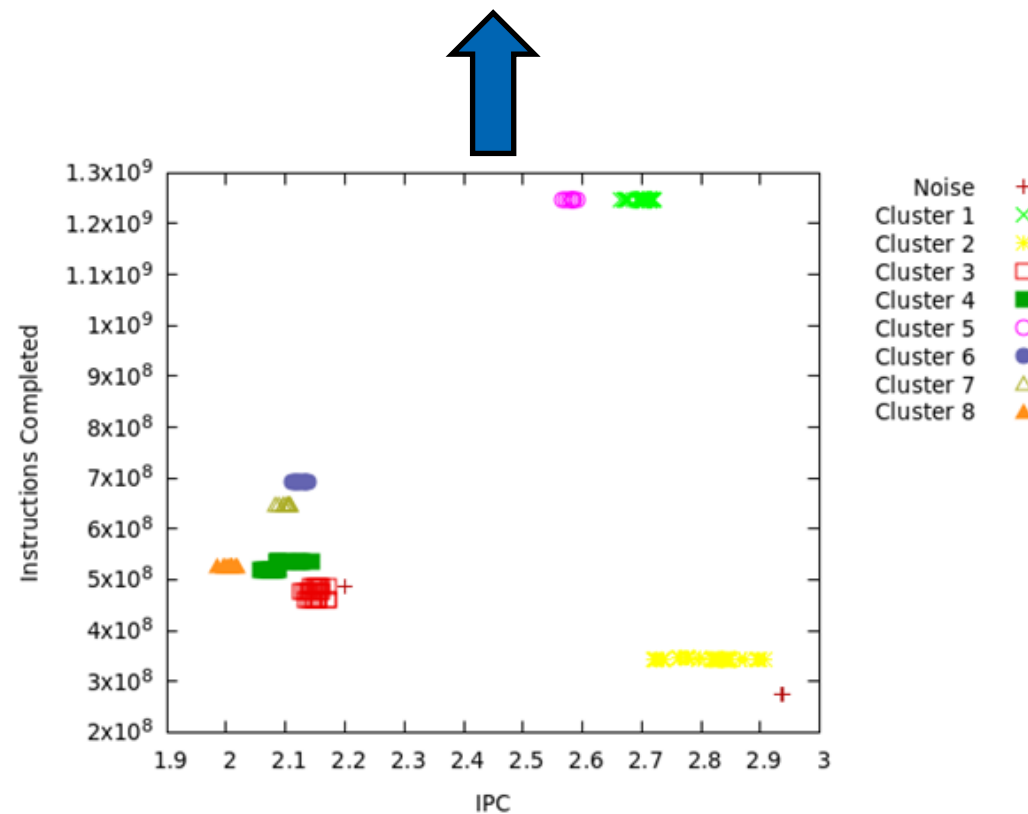
- Going back to the source code



useful instructions 2DZoom range [4.12748e+08,7.73082e+08]_c1 @ lulesh_4+4_alone.prv

**Instructions executed**

**Callstack (function, file, line)**

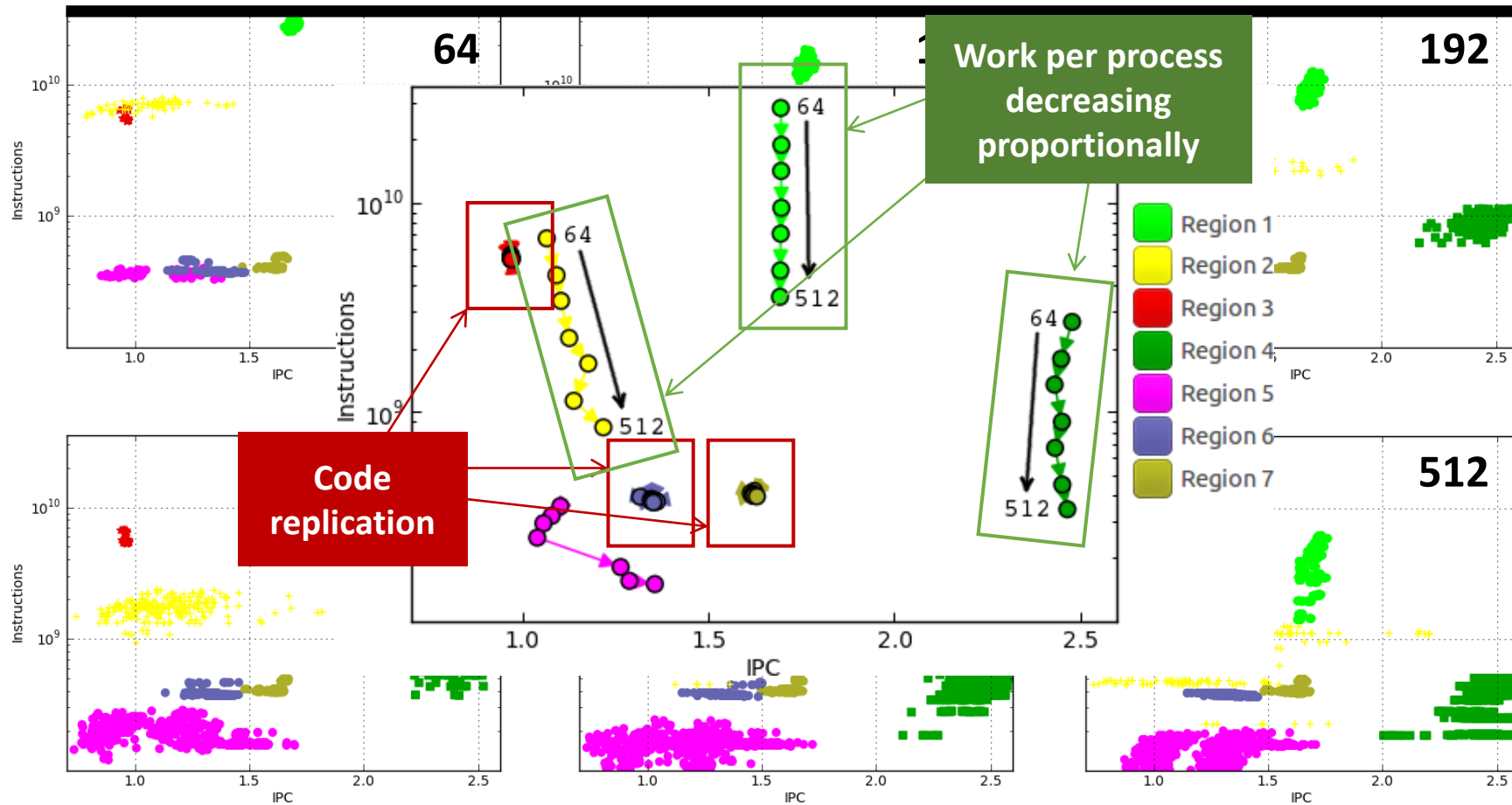☐ CommSend   ☐ TimeIncrement

# Clustering to identify structure

# Tracking to compare experiments

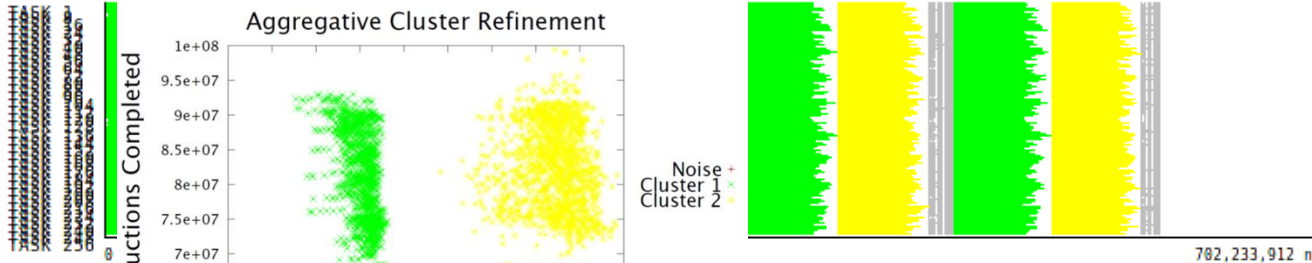- Use case: Study the scalability of the computing regions → From 64 to 512 cores
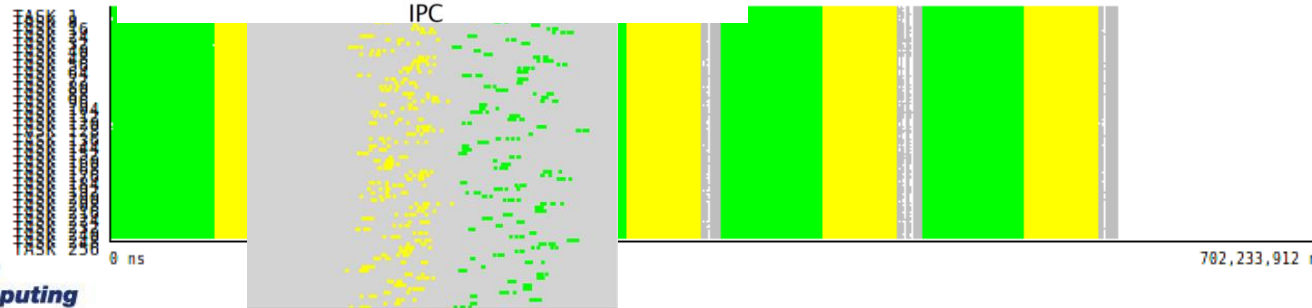
# Dimemas to predict scenarios
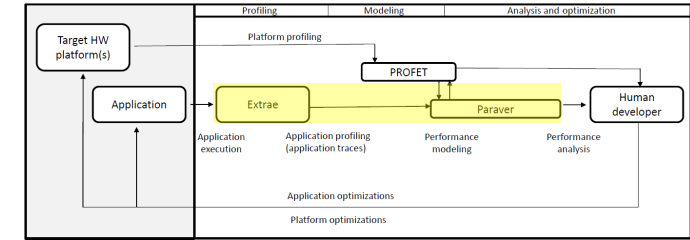
What if ….



… we increase the IPC of Cluster1?



13% gain

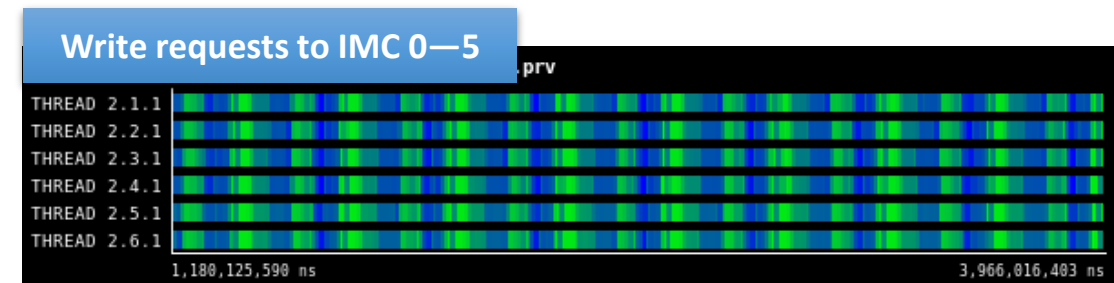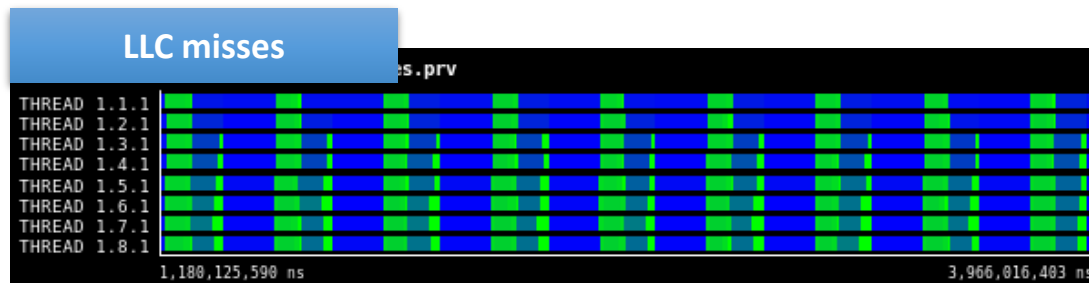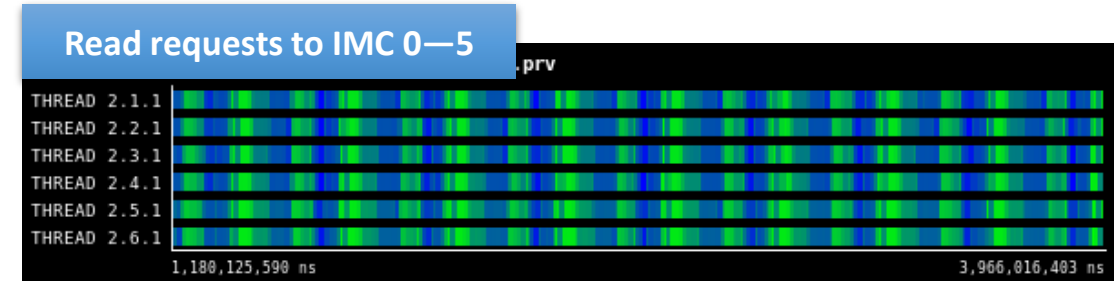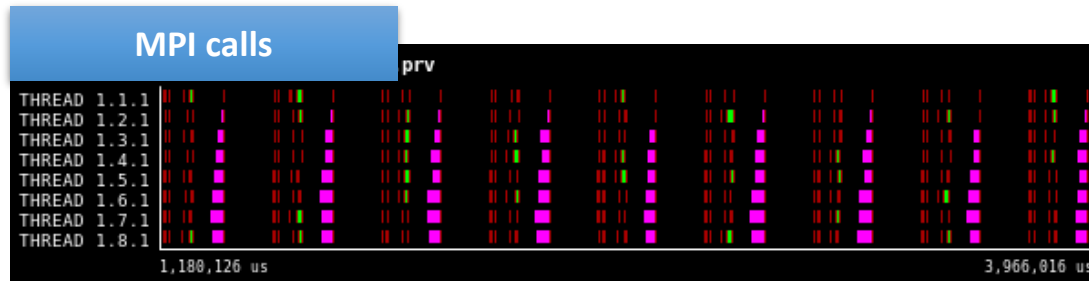… we balance Clusters 1 & 2?



19% gain

Quick estimates on when is an optimization worth the effort
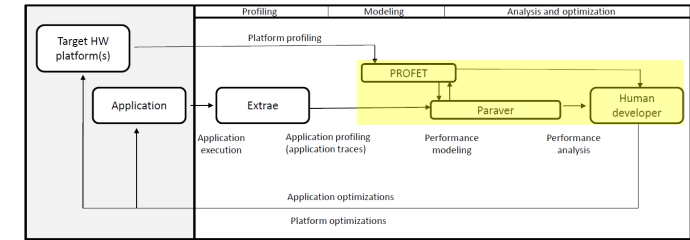
# Tools + PROFET: Current progress



- Extrae writes Paraver traces with required counters for PROFET



**MPI calls**

**Read requests to IMC 0—5**
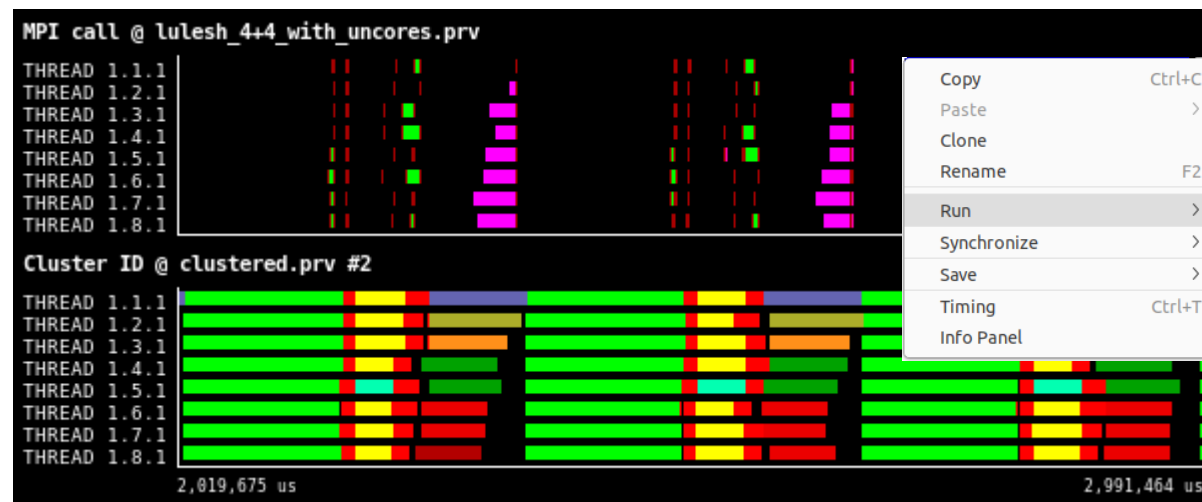
**LLC misses**

**Write requests to IMC 0—5**

Native counters measured
by application threads
at instrumentation points

Uncore counters sampled
periodically by
dedicated threads per IMC
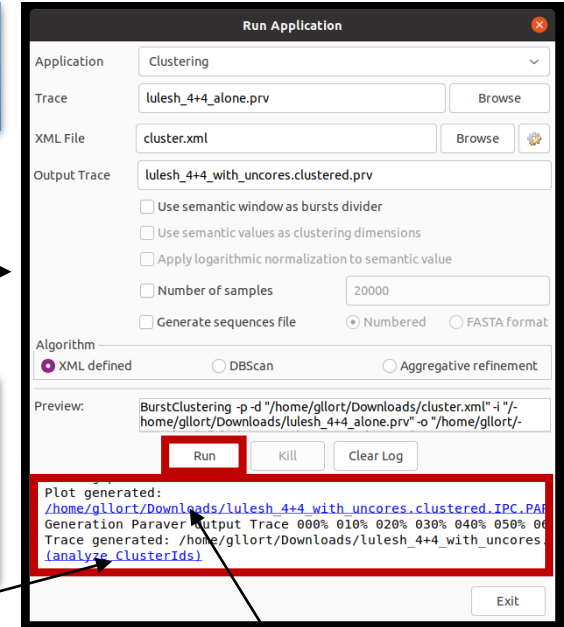
# Integrated environment



- Characterise application's memory usage on selected regions directly from Paraver context menu
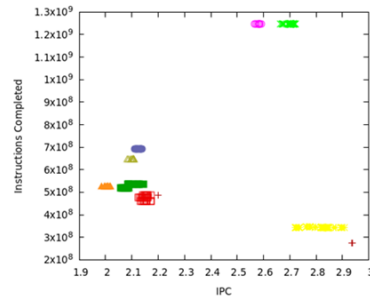


**1. Invoke tool with presets**

**2. Hit 'Run' → Tool output is reported with links**

**3. Clicking on a link automatically brings up new Paraver windows**
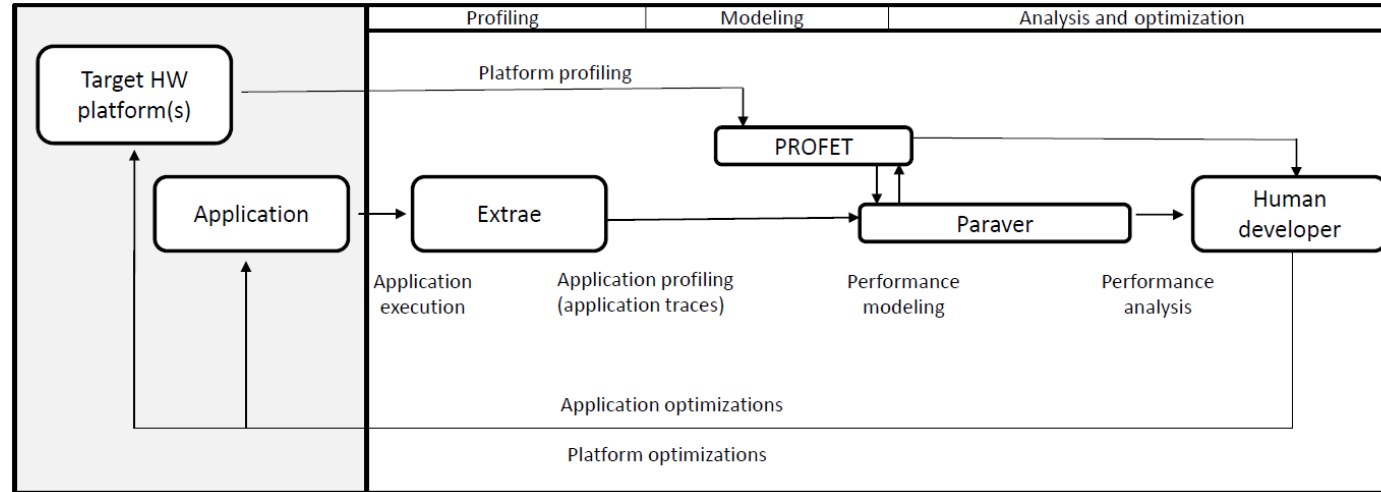
**4. Or additional reports (e.g. GNUplot)**

- Same approach for integrating PROFET
  - In-trace feedback
    - Memory usage (latency, bandwidth)
  - Additional model outputs

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# App developers: Try our cycle!



- It's easy to use!
- Just add LD_PRELOAD to your jobscript to get a trace
- Directly invoke PROFET from Paraver context menu → Run
- PROFET outputs will jump back to Paraver
- We are happy to assist you!
  - Just need the jobscript and the binary to give you useful insight!

https://tools.bsc.es