



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

BSC TOOLS: INSTRUMENTATION & ANALYSIS

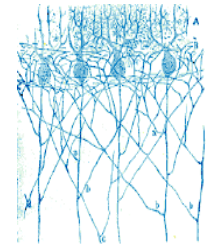
Juan Gonzalez

Why tools?

Measurements as science enablers

Vital for app. development at Exascale

- Flight instrumentation
- Work in the right direction



Important for “lawyers”

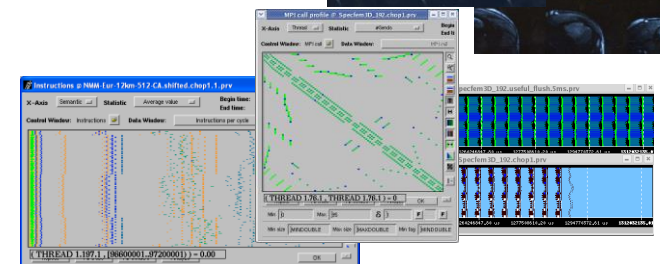
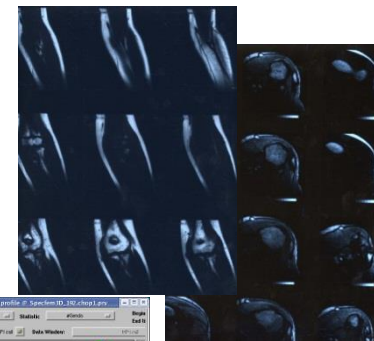
- Who to blame?

Vital for system architects

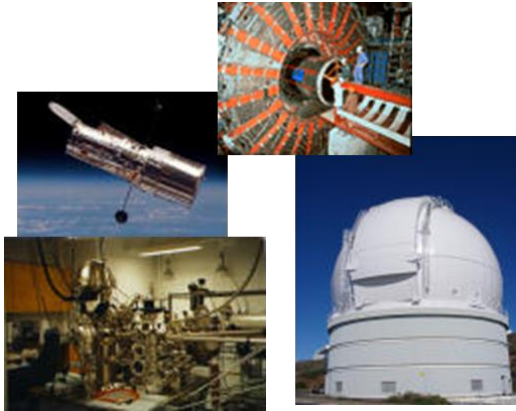
- Understand our system → Increase productivity

Performance analyst

- Expert understanding displays



Science



Physics

- One “expensive” experiment
- Lots of data to analyse
- To understand a nature system

Computer science

- `printf()`
- Timers
- Lots of speculation
 - We see $\int_a^b f(t)dt$
 - We talk about $f(t)$

Parallel computing

- Just another system
- We should use similar practices and techniques

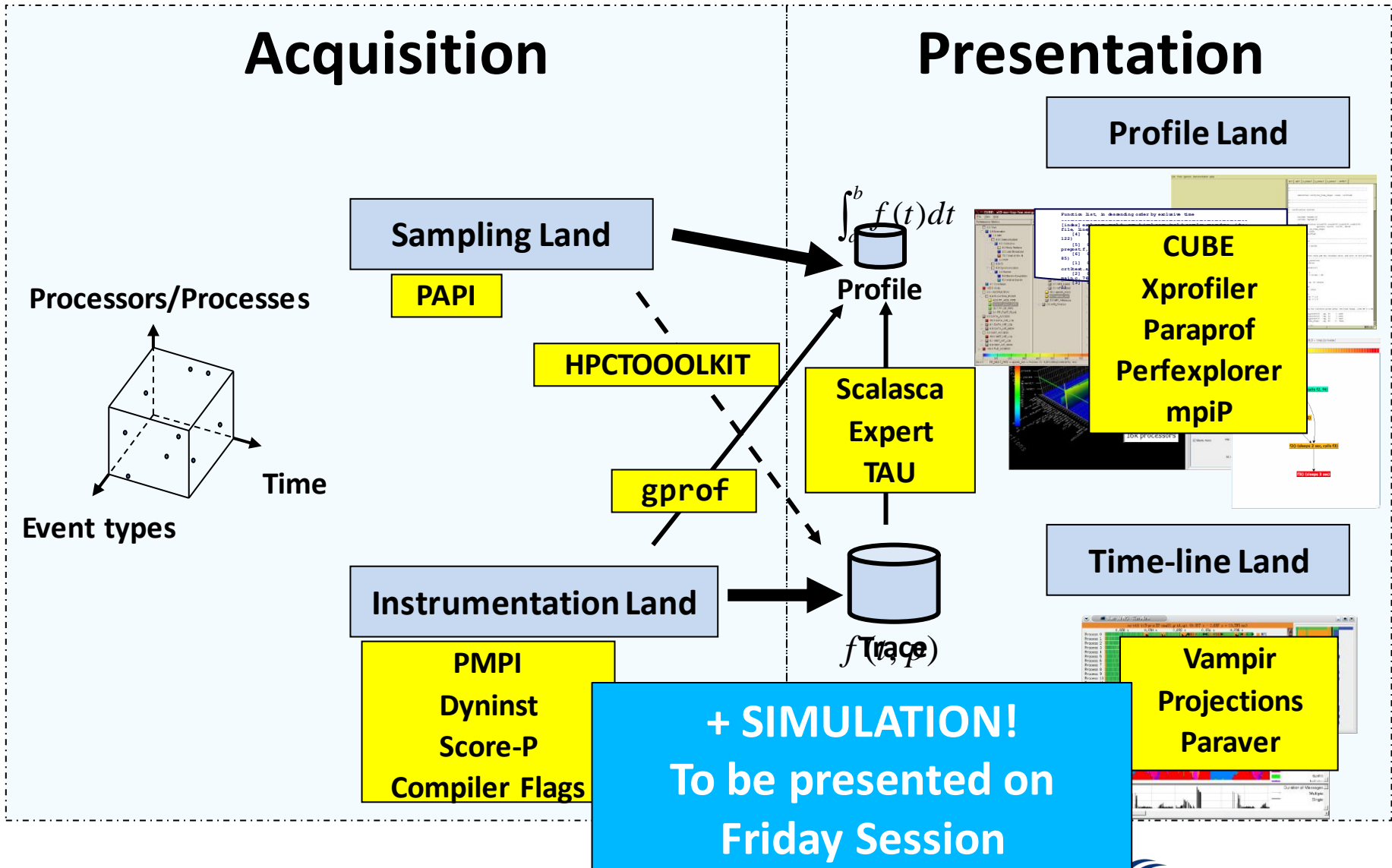
Performance Analysis Tools

- Ecosystem
- Data acquisition
- Data presentation

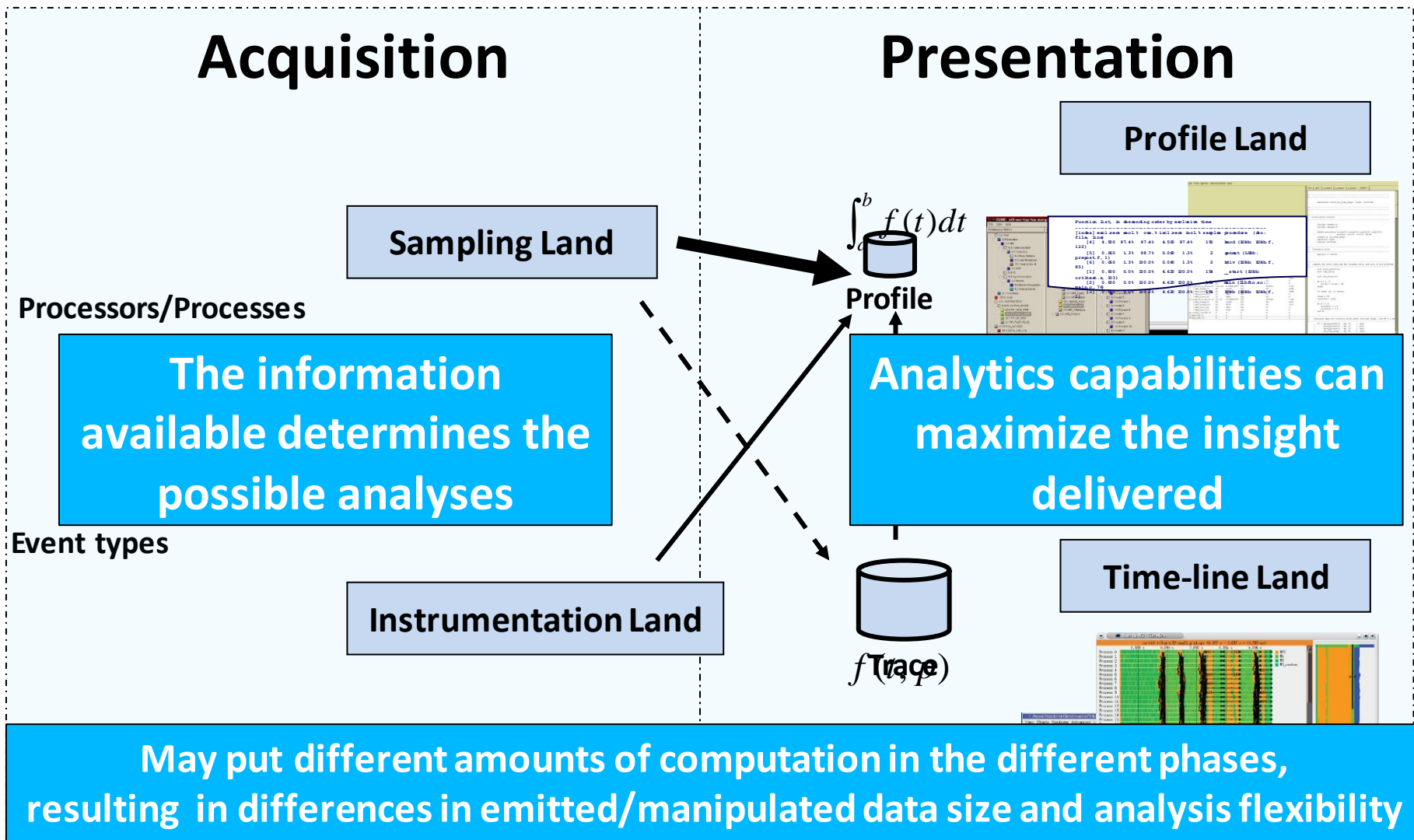
BSC Tools

- Extrae
- Paraver
- Performance Analytics

Performance analysis universe...



Performance analysis universe...



Performance tools

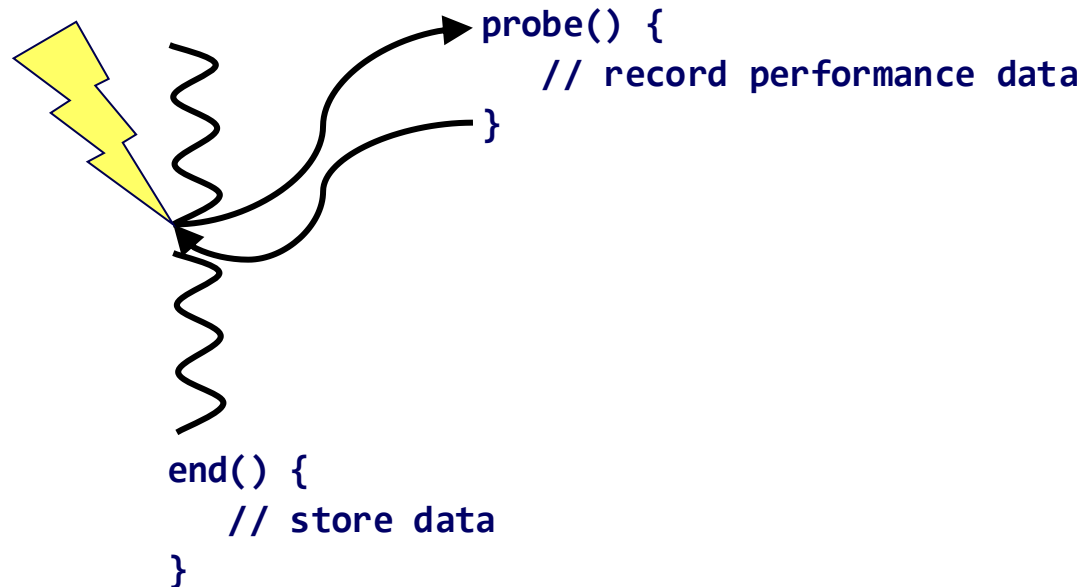
OBJECTIVE:

- Identify performance problems to help optimizing apps.

PHASES:

- Data acquisition
- (Data emission)
 - Compaction
 - Summarization
- Data presentation

Insert probes into the running program



Issues

- **Control flow: *When a probe is called?***
- ***Which is the information available?***

Sampling

- Punctual records of app. activity
- Interrupt driven: correlated (or not) with app. activity
 - POSIX clocks, PAPI overflows
- Need to project samples to total application behaviour

Instrumentation

- Probes driven by relevant application events
- Statically: link with an instrumentation library
 - PMPI / OMPTI / OMPSs / Charm++
- Dynamic: rewrite the binary
 - DynInst / DPCL

Data acquisition: Information available

Control flow

- Program counter (PC) register / Call stack

Control flow arguments

- e.g. MPI call parameters through PMPI

Communications / Synchronizations

- Messages transmitted / Tasks Dependencies

Hardware Counters

- PAPI / PMAPI

Operating System

- rusage

Runtime internals

- PERUSE (MPI statistics)

... and timing
information too!! 😊

Data acquisition: perturbation

⌘ Probe effect in the application execution

⌘ Granularity: number of probes

- Application
- Sampling frequency / what to instrument

⌘ Overhead: cost of each probe

- Control flow
- Time measurement
- Inline processing
- Storage to buffer (or disk!)

⌋ What metrics?

- Time / Counts
- Absolute / Relative

⌋ How?

- Textual
- Graphical

⌋ Which type?

- Profile: accumulated statistics per app. abstraction
- Timeline: instantaneous value of metrics vs. time vs. process

Keep in mind this objective

- **Maximize the flow of information you require...**
 - Qualitatively: colours, shapes
 - Quantitatively: numbers
- **... by a proper balance of the approaches**

Analyse your requirements!

⌋ Which of my subroutines consumes the most time?

- Do I need time measurements?
- Do I need MPI measurements?
- Do I need a timeline?

⌋ Where is the computation performing inefficiently?

- Do I need MPI measurements?
- Do I need hardware counters?
- Do I need a timeline?

⌋ I want to simulate the MPI message exchange

- Do I need time measurements?
- Do I need MPI measurements?
- Do I need a timeline?

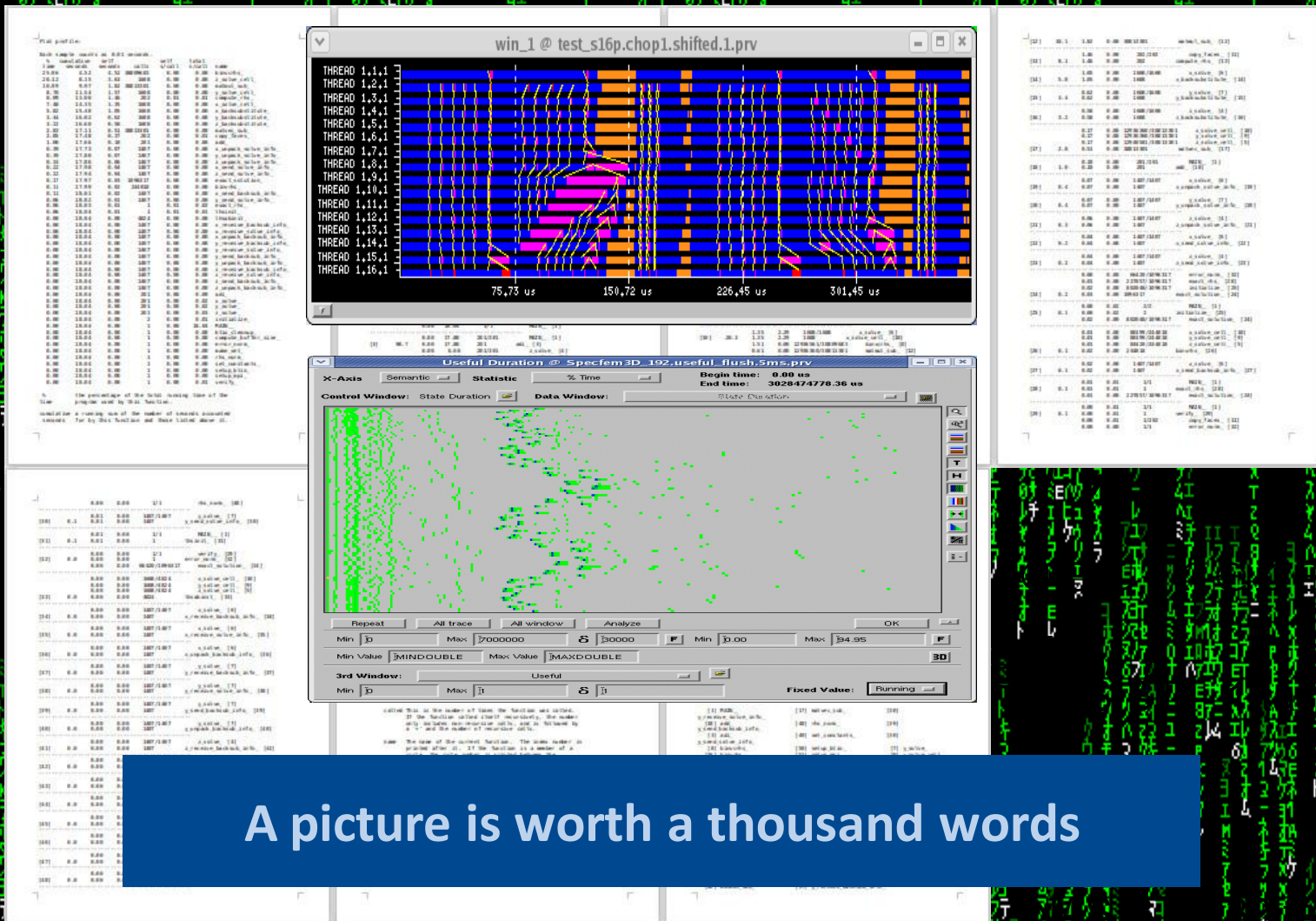


**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

BSC TOOLS

The ways of debugging & performance analysis



A picture is worth a thousand words

Since 1991

Based on traces

Flexibility and detail

Core Tools

Trace generation - Extrae

Trace analyzer - Paraver

Message passing simulator - Dimemas

Open-source

Do not speculate about your code performance

LOOK AT IT

A “different” view...

⌋ Behavioural structure vs. syntactic structure

- Algorithmic and performance
- In space and time

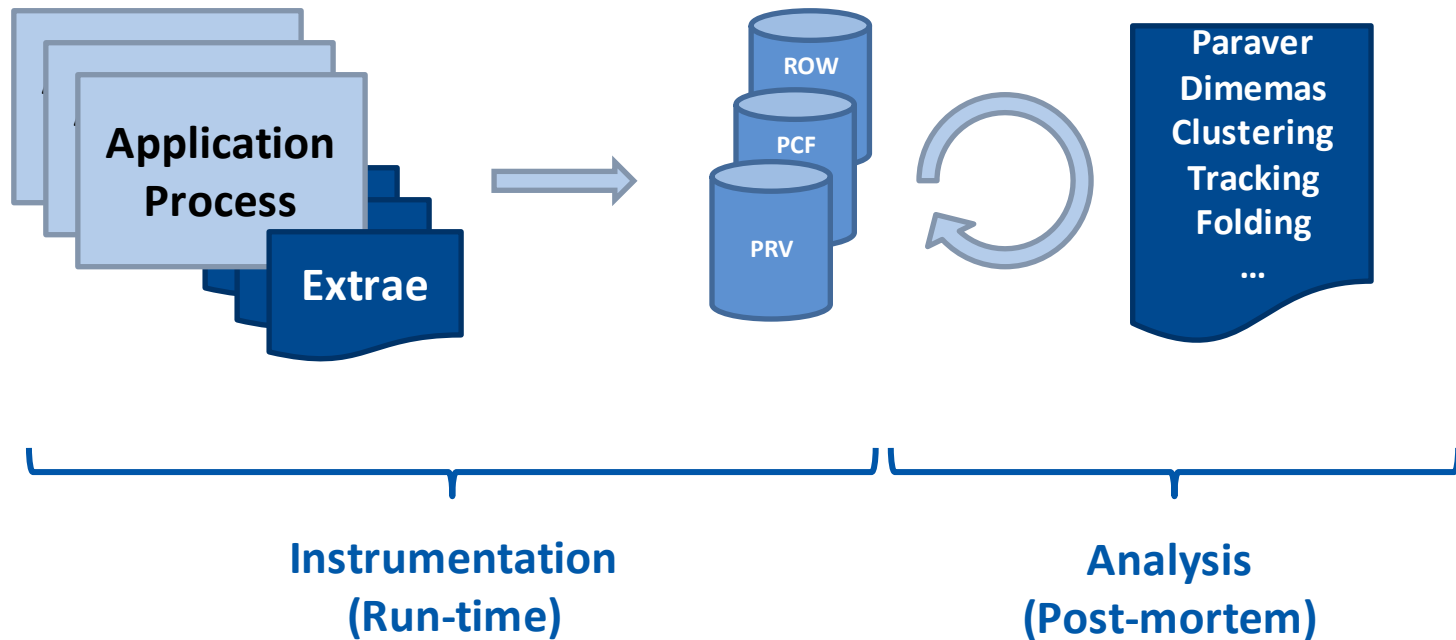
⌋ Variability

- Multimodal distributions
- Variability + synchronization → critical non linear effects

⌋ Flexibility

- Let the analyst navigate through capture data
- Gain insight minimizing a single app. execution

Basic workflow



The (Paraver) trace

⌘ Sequence of time stamped records (.prv)

- **Punctual events**
 - Something happened: when (time) & where (object/entity e.g. thread)
 - One record (type:value) per specific information
 - About the event
 - About the interval from previous event
- **Relations between objects (... communications)**
 - Source and destination
 - Attributes (... tag, value)

⌘ Symbolic information (.pcf)

- **Strings to name types and events**

⌘ System information (.row)

- **Strings to name the trace entities**



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

EXTRAE

Extrae features

Parallel programming models

- MPI, OpenMP, pthreads, OmpSs, CUDA, OpenCL, Intel MIC...

Performance Counters

- Using PAPI and PMAPI interfaces

Link to source code

- Callstack at MPI routines
- OpenMP outlined routines and their containers
- Selected user functions

Periodic samples

User events (Extrae API)

How does Extrae work?

Symbol substitution through LD_PRELOAD

- For production binaries
- Specific library for runtime/s
- No knowledge about the app. required

Programming model	Library
Serial	libseqtrace
Pure MPI	libmpitrace[f] ¹
Pure OpenMP	libomptrace
Pure Pthreads	libpttrace
CUDA	libcudatrace
MPI + OpenMP	libompitrace[f] ¹
MPI + Pthreads	libptmpitrace[f] ¹
Mpi + CUDA	libcudampitrace[f] ¹

¹ for Fortran codes

Dynamic instrumentation

- For production binaries
- Specify the functions to be instrumented

Based on Dyninst
U.Wisconsin/U.Maryland

Other possibilities

- Static linking (e.g. PMPI @ BG/Q)
- OmpSs instrumentation calls
- Extrae API

How to use Extrae?

⌋ Adapt your submission script

- Point the Extrae library require in the LD_PRELOAD
- Point to the XML instrumentation control file

⌋ Specify the data to be captured

- Editing the XML instrumentation control file

⌋ Run and get the trace

Extrae 2.3.4 User's Guide available in

<http://www.bsc.es/computer-sciences/performance-tools/documentation>

Default control files and further examples within installation in
\$EXTRAE_HOME/share/example

application.job

```
#!/bin/bash
...
# @ total_tasks = 4
# @ cpus_per_task = 1
# @ tasks_per_node = 4
...

srun ./myMPIbin/my_MPI_binary
```

trace.sh

```
#!/bin/sh

export EXTRAE_HOME=...
export EXTRAE_CONFIG_FILE=extrae.xml
export LD_PRELOAD=${EXTRAE_HOME}/lib/libmpitrace.so

# Run the desired program
$*
```

⌋ If the source code is available you can use it to...

⌋ ... scalability

- `Extrac_shutdown()`
- `Extrac_restart()`

⌋ ... inject your own events

- `Extrac_[n]event(type[s], value[s])`
- `Extrac_[n]eventandcounters(...)`

Approaches towards scalability

Manual selection

- **Tracing On-Off**
 - Internal / External
- **Trace/Buffer maximum sizes**
- **Information gathered**
 - Call path depth
- ***Bursts Mode* + Software Counters**

Automatic reduction

- **Spectral analysis**
- **Computation structure detection**

Emitting “relevant” data (i.e. *Burst Mode*)

⌋ Important information → Details

⌋ Not that important → Software counters

⌋ What is important?

- First order approach: Computation !!!
- MPI: a gas. Fills whatever space you give it. Very often not the major cause of problems

⌋ Major computation bursts (i.e. $> X$ ms)

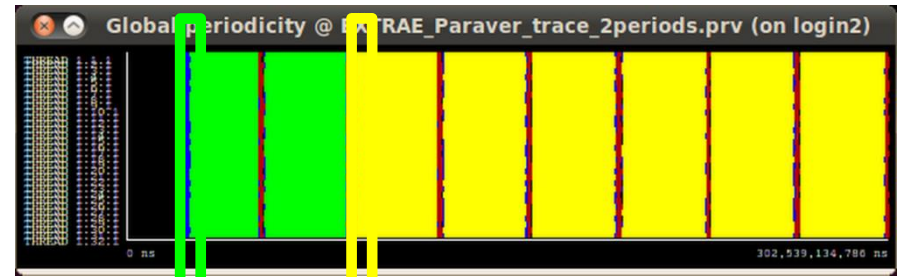
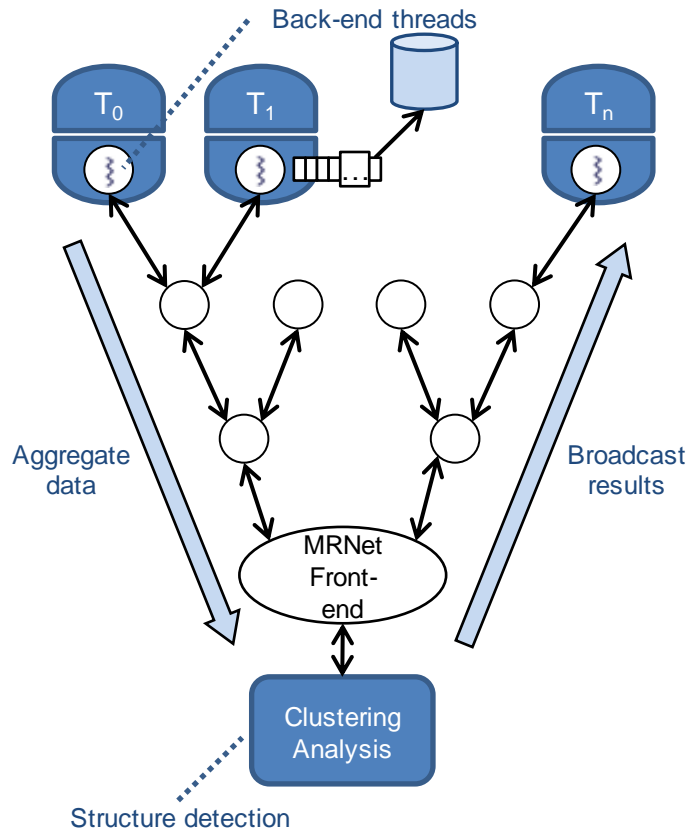
- Entry and exit timestamps and hardware counters

⌋ Communication phases.

- Software counters:
 - # MPI calls, aggregated bytes, %time in MPI, ...

J. Labarta et. al., *Scalability of tracing and visualization tools*, PARCO 2005

Scalability: online automatic interval selection



Detail trace only for small interval

G. Llort et al., *Scalable tracing with dynamic levels of detail*, ICPADS 2011



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

PARAVER

What is Paraver

⌘ A browser ...

⌘ ...to manipulate (visualize, filter, cut, combine, ...) ...

⌘ ... sequences of time-stamped events ...

⌘ ... with a multispectral philosophy ...

⌘ ... and a mathematical foundation ...

⌘ ... that happens to be mainly used for performance analysis

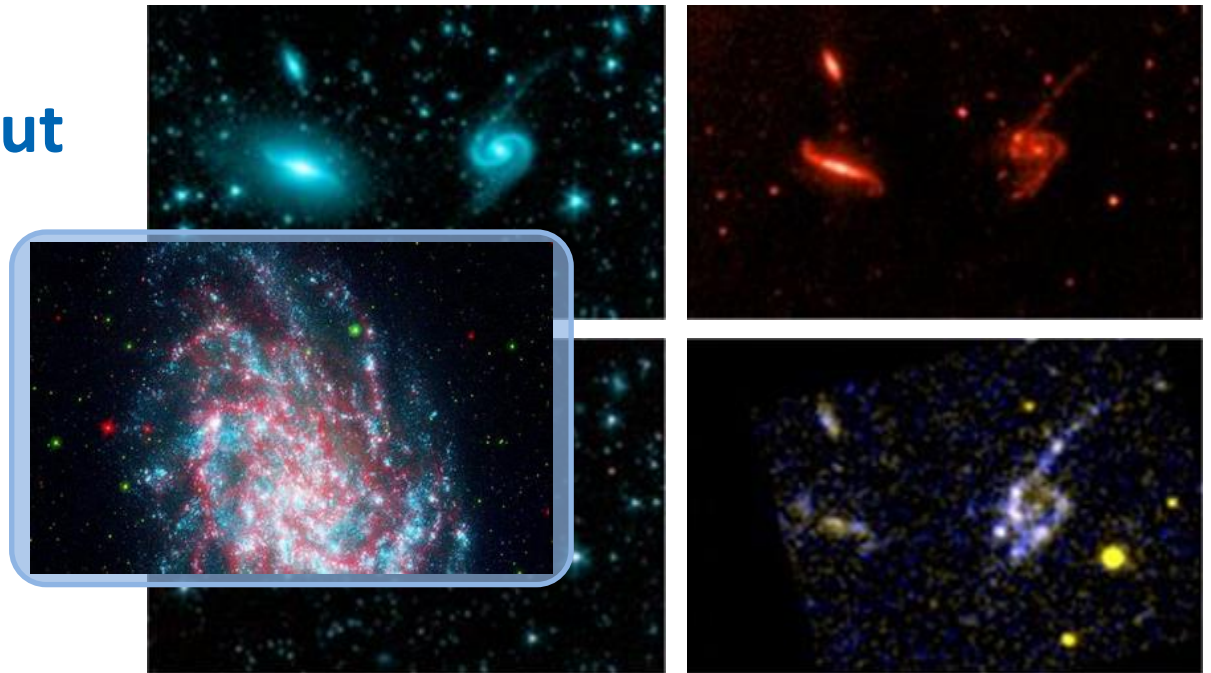
Multispectral imaging

⌋ Different looks at one reality

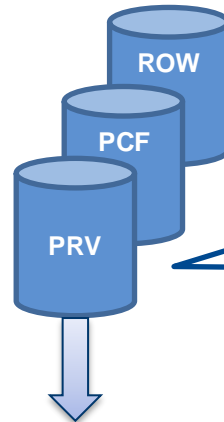
- Different light sources using filters

⌋ Highlight different aspects

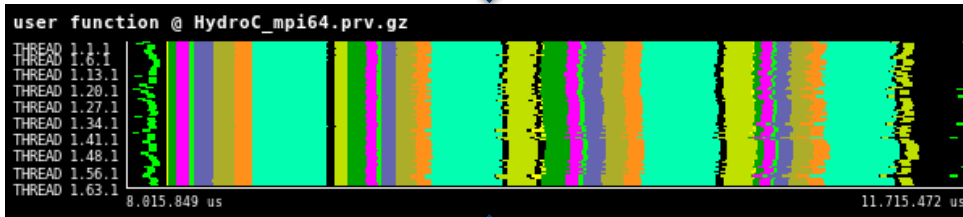
⌋ Zooming In & Out



Paraver displays



The trace
MPI calls, OpenMP regions, user functions, peer-to-peer & collective communications, performance counters, samples...



Timelines



2D / 3D Tables (statistics)

Timelines: description

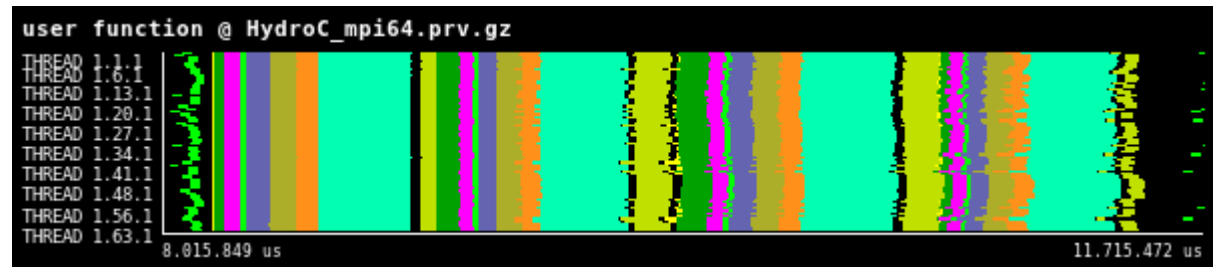
Objects

Process dimension

- Thread (default)
- Process
- Application
- Workload

Resource dimension

- CPU
- Node
- System



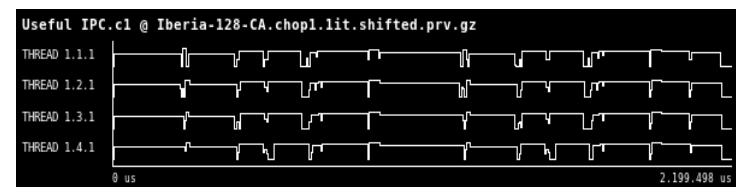
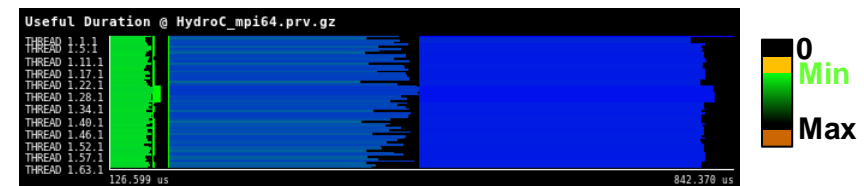
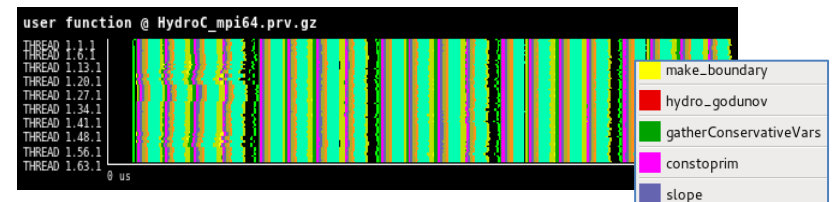
Time

Timelines: Semantics

⌋ Each window computes a function of time per object

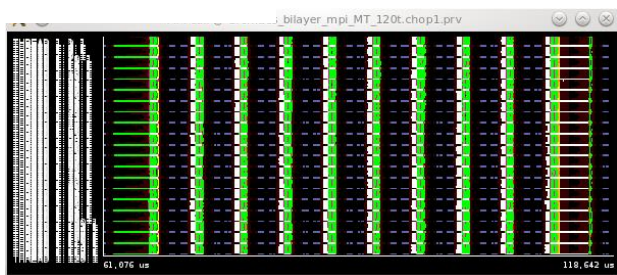
⌋ Two types of functions

- **Categorical**
 - State, User function, MPI call...
 - Color encoding
 - 1 color per value
- **Numerical**
 - IPC, instructions, cache misses, computation duration...
 - Gradient encoding
 - **Black** (or background) for zero
 - From **light green to dark blue**
 - Limits in **yellow** and **orange**
 - Function line encoding

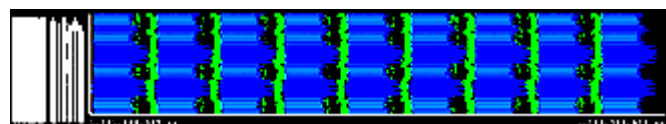
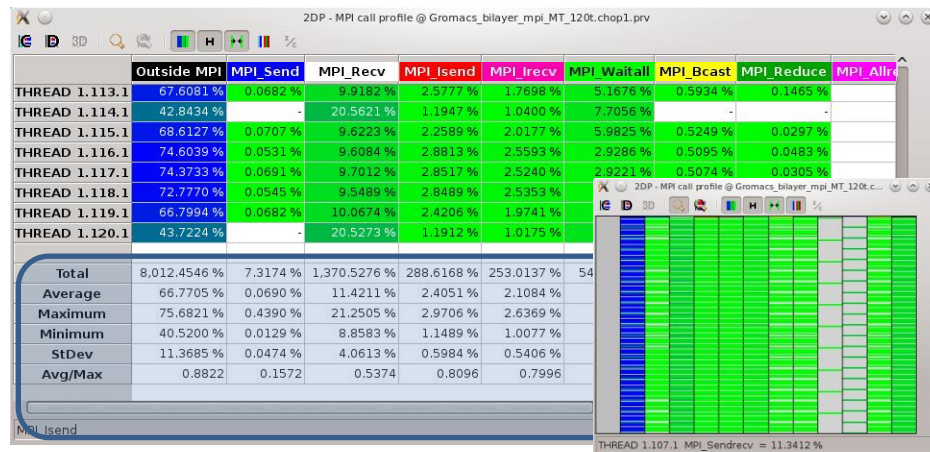


From timelines to tables

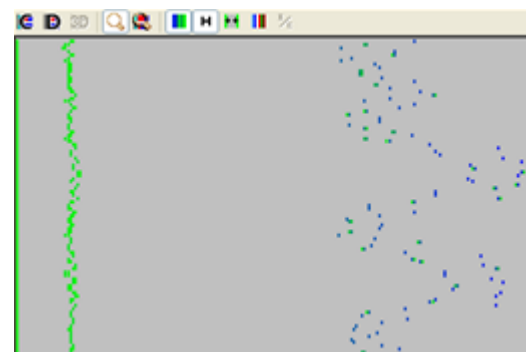
MPI calls



MPI calls profile

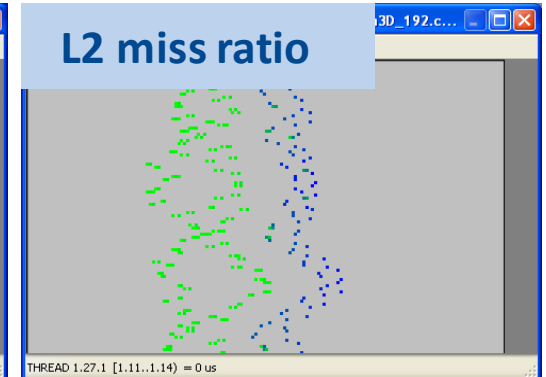
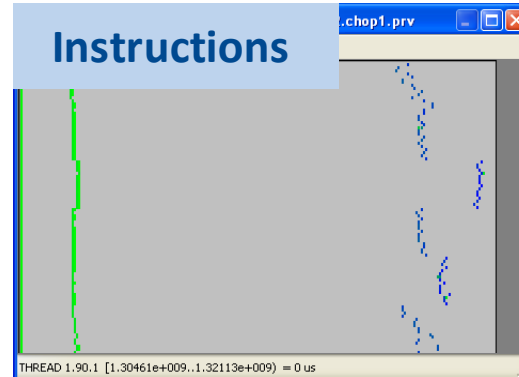
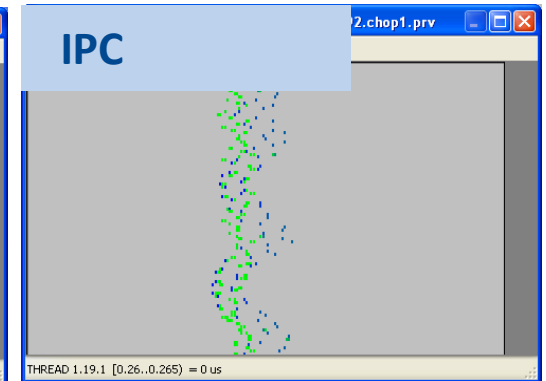
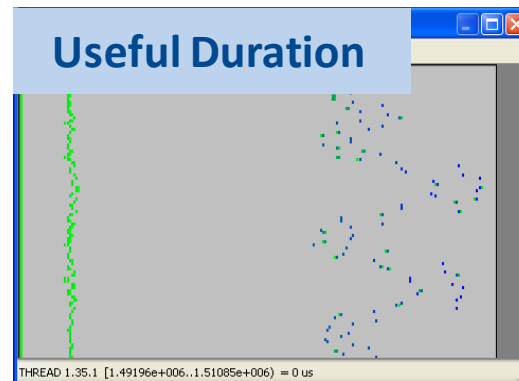
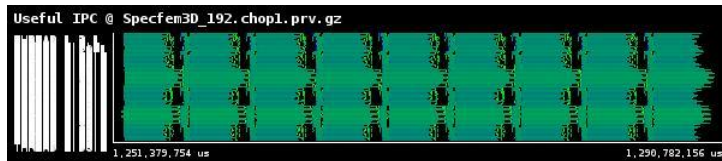
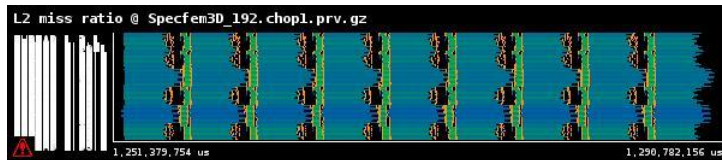
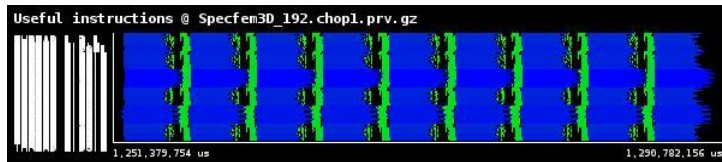
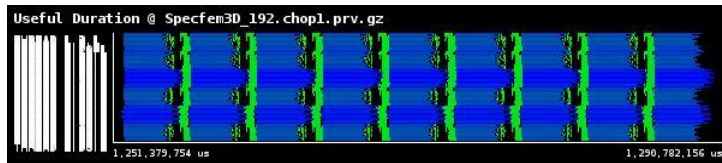


Computation duration



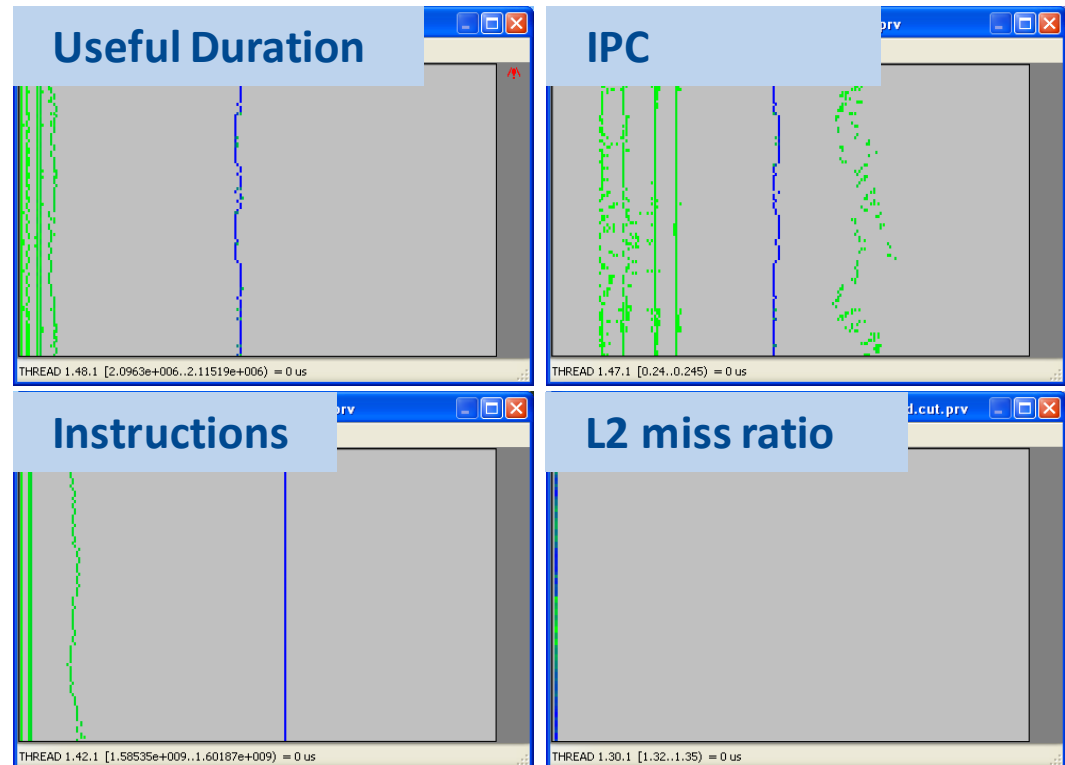
Computation duration histogram

Analyzing variability through histograms and timelines



Analyzing variability through histograms and timelines

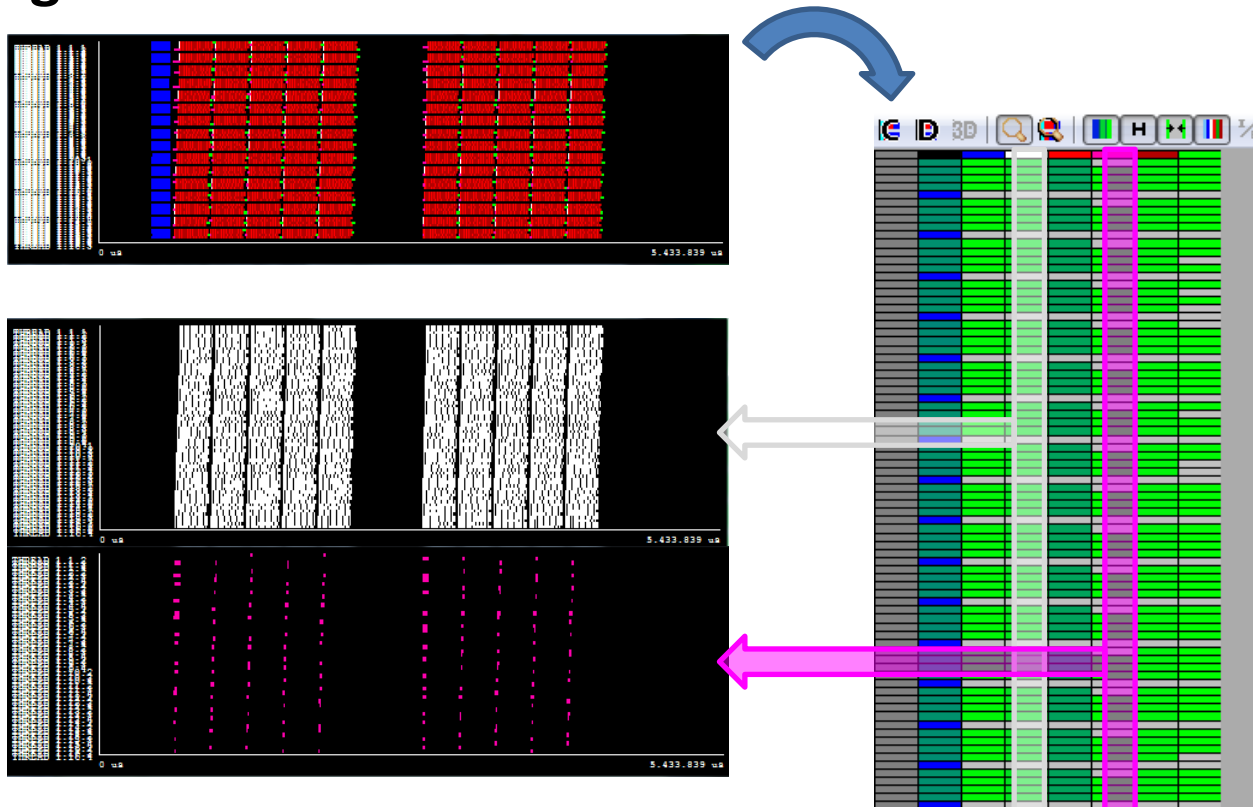
⌋ By the way: six months later...



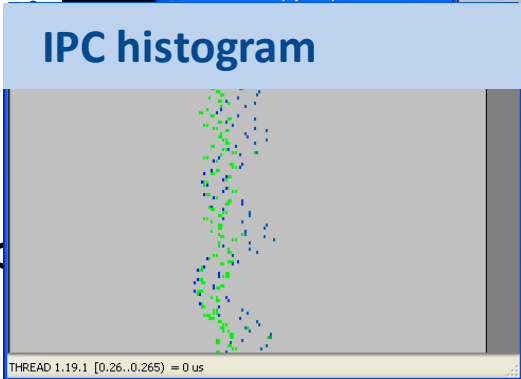
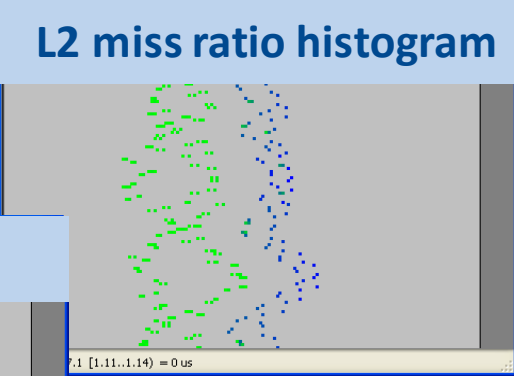
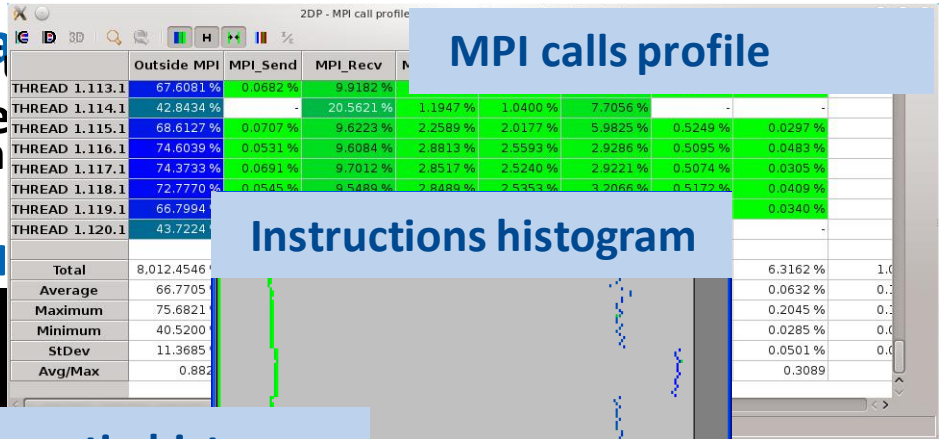
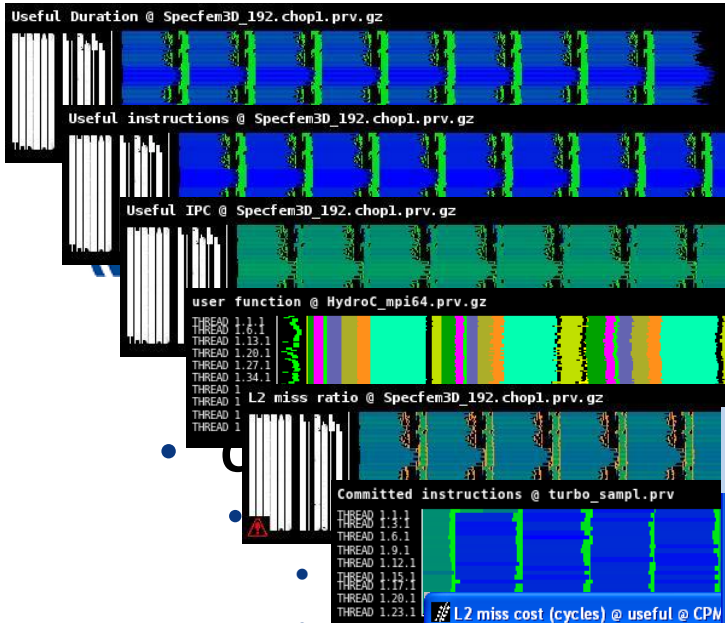
Tables: back to timelines

⌋ Where in the timeline do certain values appear?

- e.g. which is the time distribution of a given routine?
- e.g. when does a routine occur in the timeline?



Configuration files



Cycles wasted per L2 miss
(cache misses, PS, MIPS, IPC...)

Comm. bandwidth

- M
- Op
- ... and many more:

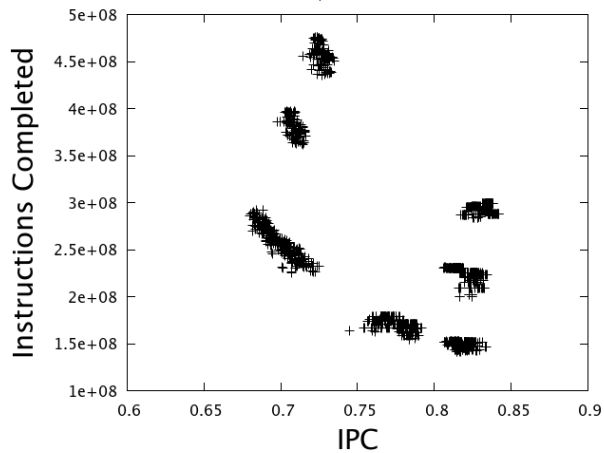
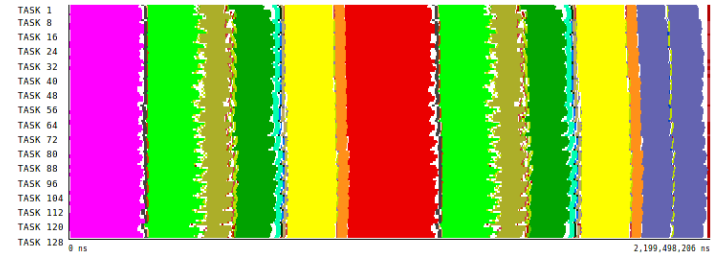
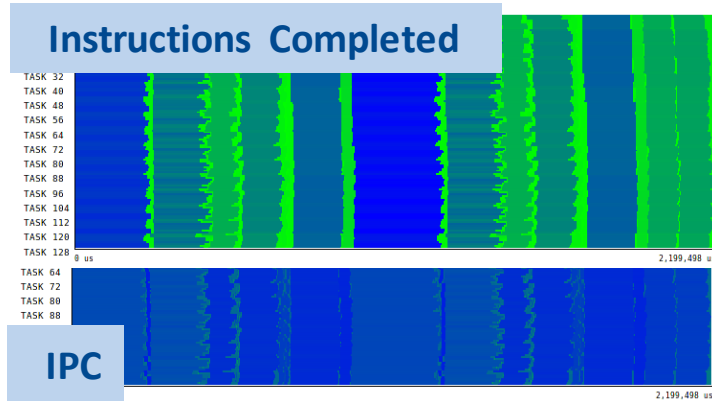


**Barcelona
Supercomputing
Center**

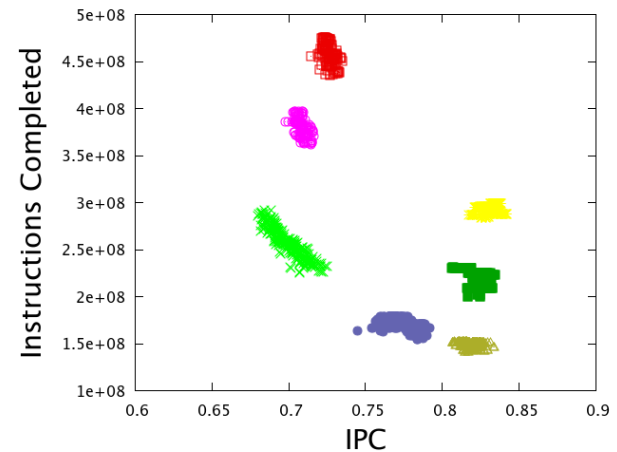
Centro Nacional de Supercomputación

PERFORMANCE ANALYTICS

Clustering to identify structure



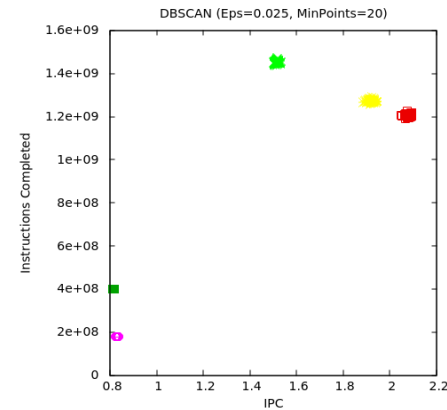
DBSCAN



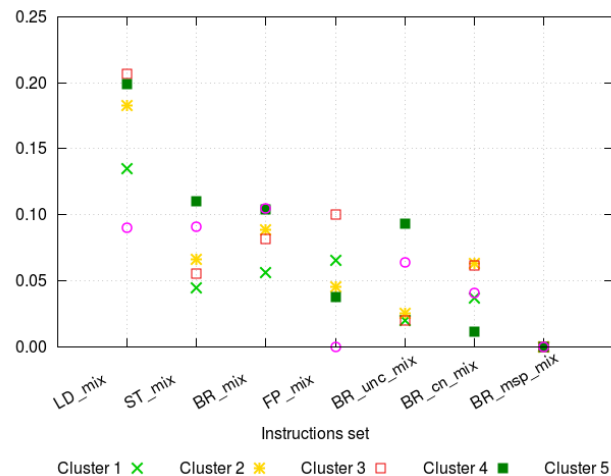
J. Gonzalez et. al., *Automatic detection of parallel applications computation phases*, IPDPS'09

Projecting hardware counters based on clustering

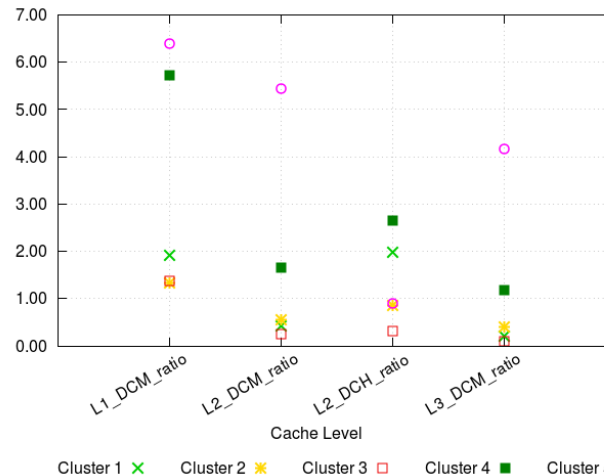
Full per region HWC characterization from a single run



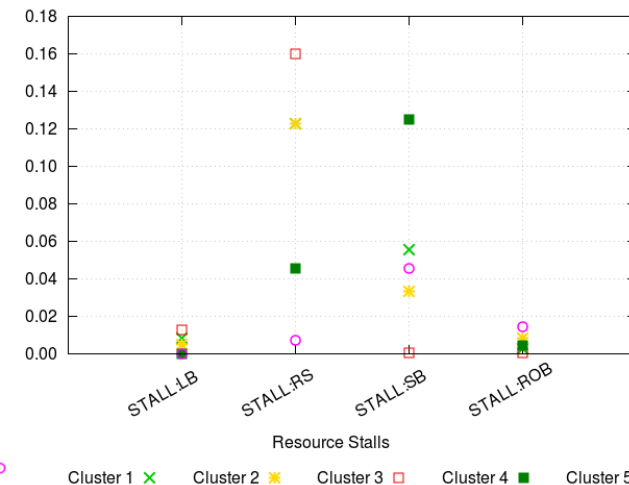
Instruction mix



Miss ratios

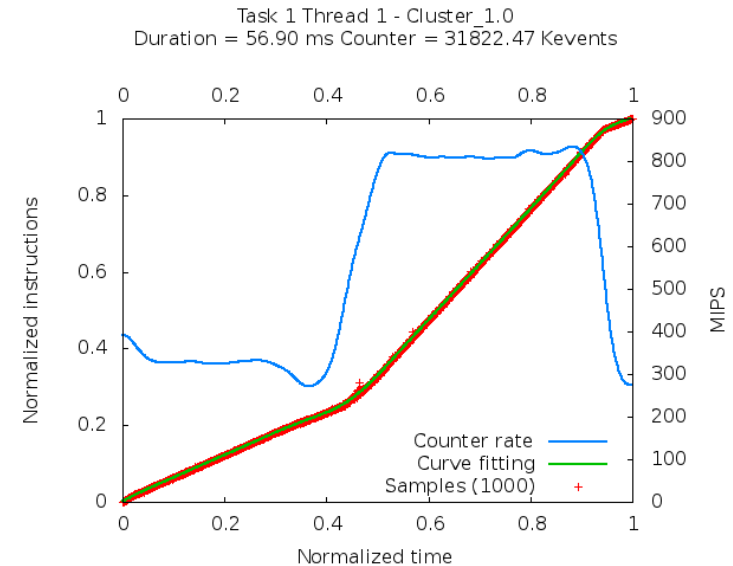
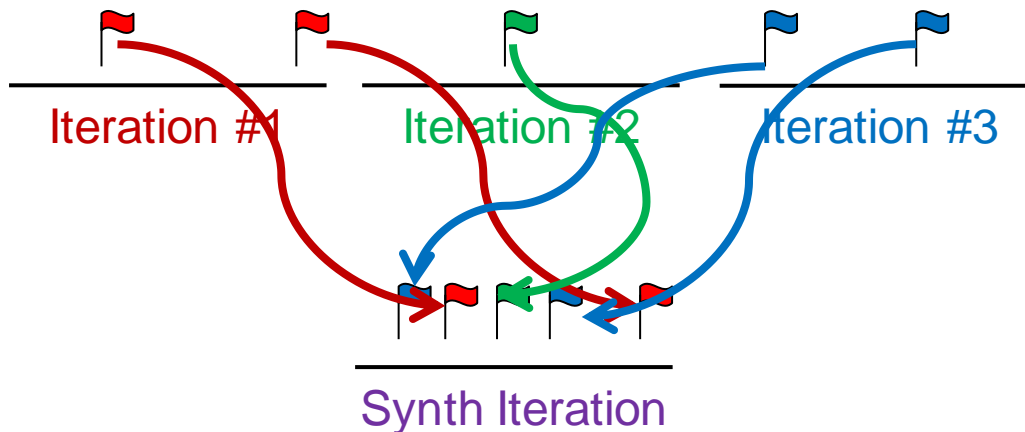


Stalls



Folding mechanism

- Detailed time evolution of metrics with low-frequency sampling
- Combining sampling and instrumentation
- Taking advantage of repetitiveness

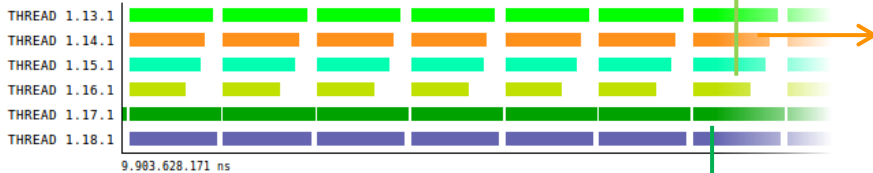


Unveiling Internal Evolution of Parallel Application Computation Phases (ICPP 2011)

Combined clustering + folding

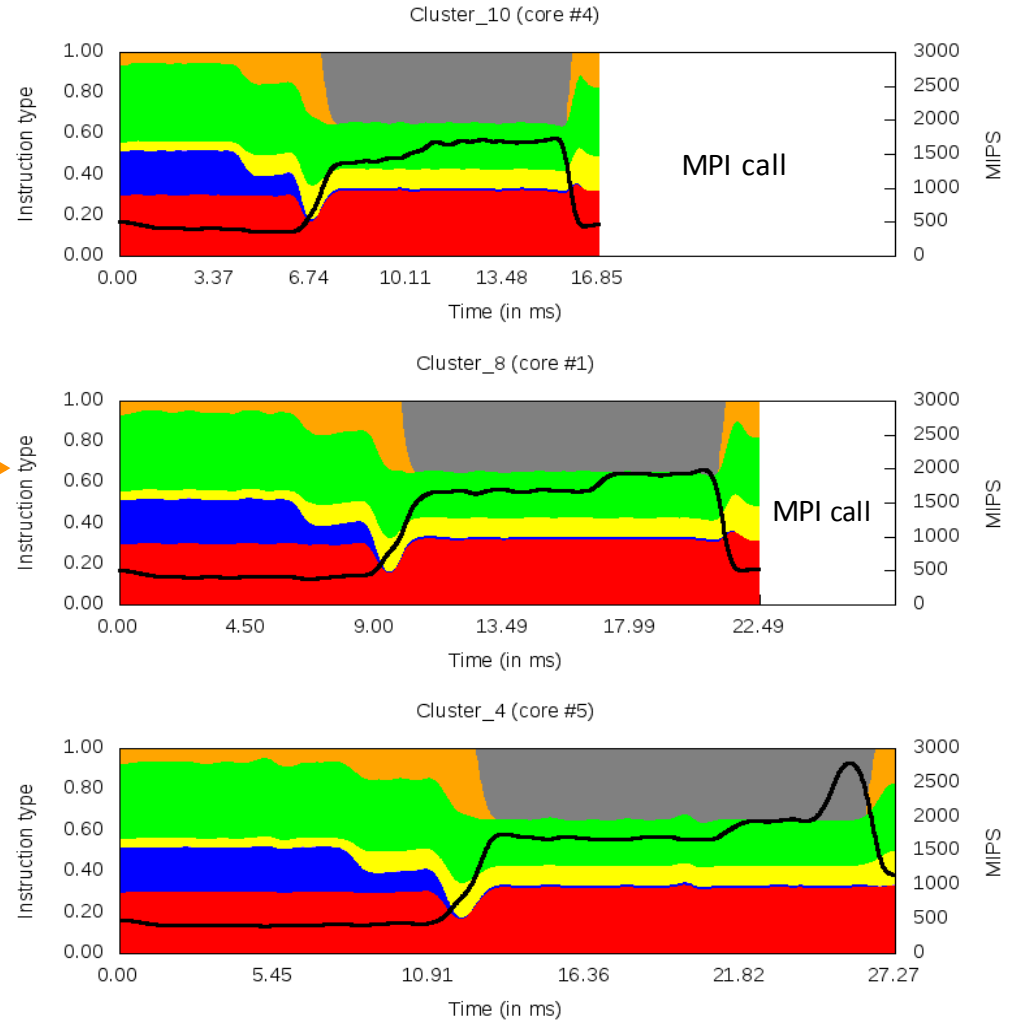
- ⌋ Instantaneous values
- ⌋ All metrics
- ⌋ From a single run
- ⌋ “No” overhead

ClusterID @ cgpop.Linux_icc.180x120.chop2.clustered.prv



CGPOP -1D

Instruction mix model for the unbalanced CGPOP on different cores of the same hexacore chip





**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

CONCLUSIONS

Measure!

⌘ Measuring your application performance is crucial

⌘ Different acquisition and presentation techniques

- Sampling / Instrumentation
- Profiles / Time-lines
- Provide different granularity

⌘ Make the appropriate questions and select the appropriate tools

⌋ Extrae: tracing mechanism

- Supports the most commonly used parallel runtimes
- Is able to capture multiple types of information

⌋ Paraver: trace visualizer & analyser

- Time-lines & Tables
- Flexible & Highly configurable (CFGs!)

⌋ Performance Analytics

- Maximum insight from a single execution
- Detailed and understandable



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

TOOLS DEMO + HANDS-ON



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

EXTRA!

TUNING THE EXTRAE XML FILE

Trace control xml

extrae.xml

```
<?xml version='1.0'?>
```

```
<trace enabled="yes" home="/gpfs/apps/CEPBAT00LS/64.hwc">
```

```
<mpi enabled="yes">  
  <counters enabled="yes" />  
</mpi>
```

```
<openmp enabled="no">  
  <locks enabled="no" />  
  <counters enabled="yes" />  
</openmp>
```

```
<callers enabled="yes">  
  <mpi enabled="yes">1-3</mpi>  
</callers>
```

```
<user-functions enabled="no" list="/home/bsc41/bsc41273/user-functions.dat">  
  <max-depth enabled="no">3</max-depth>  
  <counters enabled="yes" />  
</user-functions>
```

...

Activate/not MPI/OpenMP tracing and features

Add call stack info (number of levels) at tracing points

Add instrumentation at specified user functions
Requires Dyninst based mpitrace

Trace control xml

extrae.xml (cont)

Specification of counters emitted to trace

When to rotate between groups

Groups

```
...  
  
<counters enabled="yes">  
  <cpu enabled="yes" starting-set-distribution="1">  
    <set enabled="yes" domain="all" changeat-globalops="5">  
      PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_L1_DCM  
    </set>  
    <set enabled="yes" domain="user" changeat-globalops="5">  
      PAPI_TOT_INS,PAPI_FP_INS,PAPI_TOT_CYC  
    </set>  
  </cpu>  
  
  <network enabled="yes" />  
  
  <resource-usage enabled="yes" />  
</counters>  
  
...
```

Interconnection network counters
Just at end of trace because of large
acquisition overhead

OS info (context switches,...)

Trace control xml

extrae.xml (cont)

```
...  
<storage enabled="no">  
  <trace-prefix enabled="yes">TRACE</trace-prefix>  
  <size enabled="no">5</size>  
  <temporal-directory enabled="yes" make-dir="no">/scratch</temporal-directory>  
  <final-directory enabled="yes" make-dir="no">/gpfs/scratch/bsc41/bsc41273</final-  
directory>  
  <gather-mpits enabled="no" />  
</storage>  
  
<buffer enabled="yes">  
  <size enabled="yes">150000</size>  
  <circular enabled="no" />  
</buffer>  
  
<trace-control enabled="yes">  
  <file enabled="no" frequency="5m">/gpfs/scratch/bsc41/bsc41273/control</file>  
  <global-ops enabled="no"></global-ops>  
</trace-control>  
...
```

Control of emitted trace ...

... name, tmp and final dir ...

... max (MB) per process size (stop tracing when reached)

Size of in core buffer (#events)

External activation of tracing (creation of file will start tracing)

Trace control xml

mpitrace.xml (cont)

```
...  
  
<others enabled="yes">  
  <minimum-time enabled="no">10m</minimum-time>  
  <terminate-on-signal enabled="no">USR2</terminate-on-signal>  
</others>  
  
<bursts enabled="no">  
  <threshold enabled="yes">500u</threshold>  
  <counters enabled="yes" />  
  <mpi-statistics enabled="yes" />  
</bursts>  
  
<cell enabled="no">  
  <spu-file-size enabled="yes">5</spu-file-size>  
  <spu-buffer-size enabled="yes">64</spu-buffer-size>  
  <spu-dma-channel enabled="yes">2</spu-dma-channel>  
</cell>  
  
</trace>
```

Stop tracing after elapsed time ...

... or when signal received

... emit only computation bursts of a minimal duration ...

... plus summarized MPI events

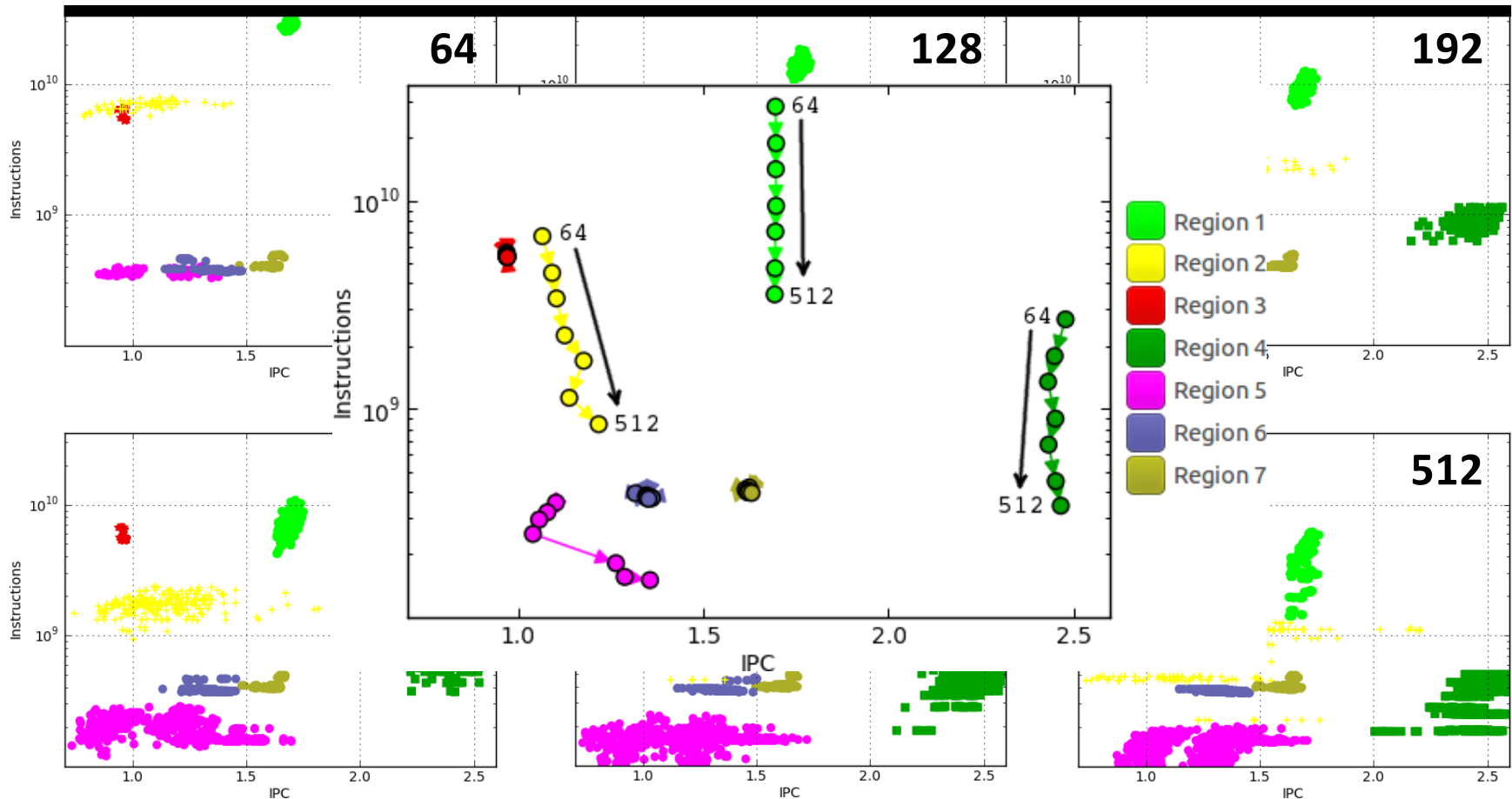
Cell specifics

MORE ON ANALYTICS

Tracking structural evolution

Frame sequence: clustered scatterplot as core counts increases

OpenMX
Strong scaling



Sampling frequency

⌋ Trade-off:

- Too low → no detail
- Too high → too much overhead

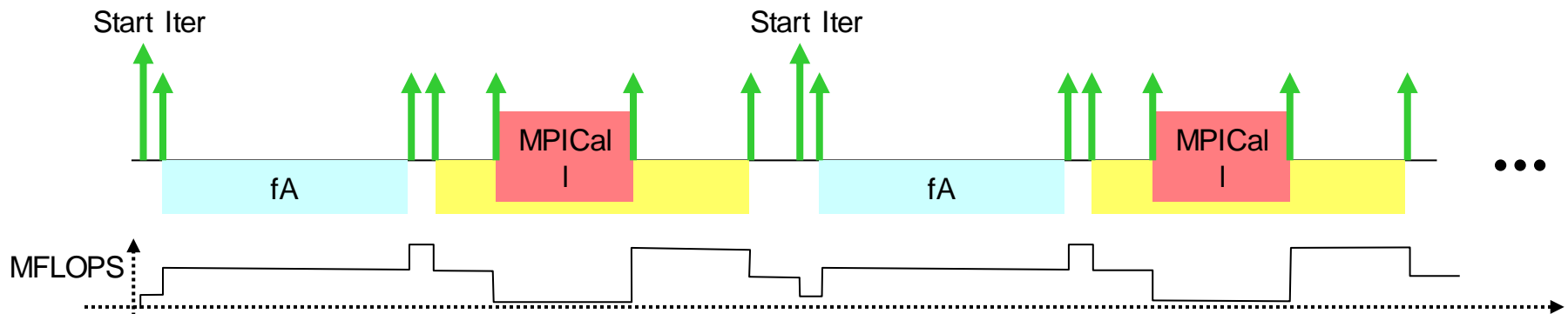
⌋ Challenge: Can we get

- lot of detail, very fine grain information :
 - i.e. “instantaneous” performance metric rates
- With very little overhead:
 - ie. sampling a few times per second

Instrumentation

Events correlated to specific program activity

- Start/exit iterations, functions, loops,...



Different intervals:

- May be very large, may be very short
- Variable precision

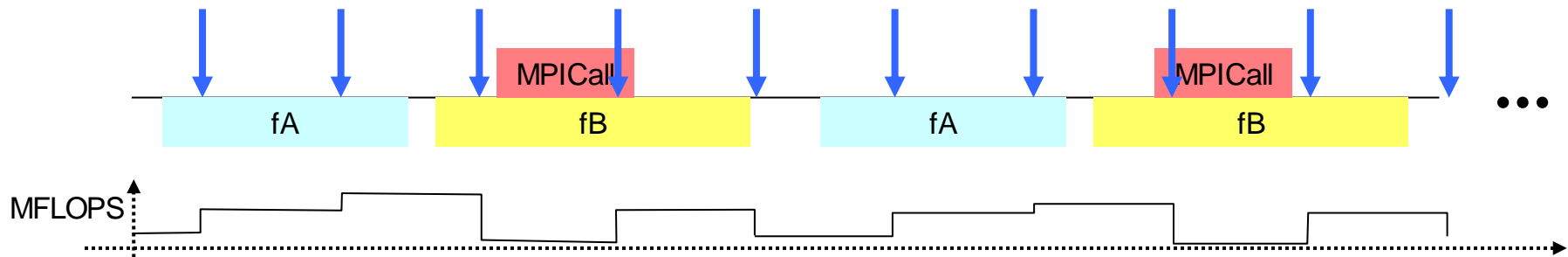
Captured data:: Hardware counters, call arguments, call path,....

Accurate statistics: profiles, ...

Sampling

Events uncorrelated to program activity (at least not specific)

- Time (or counter) overflow



Controlled granularity:

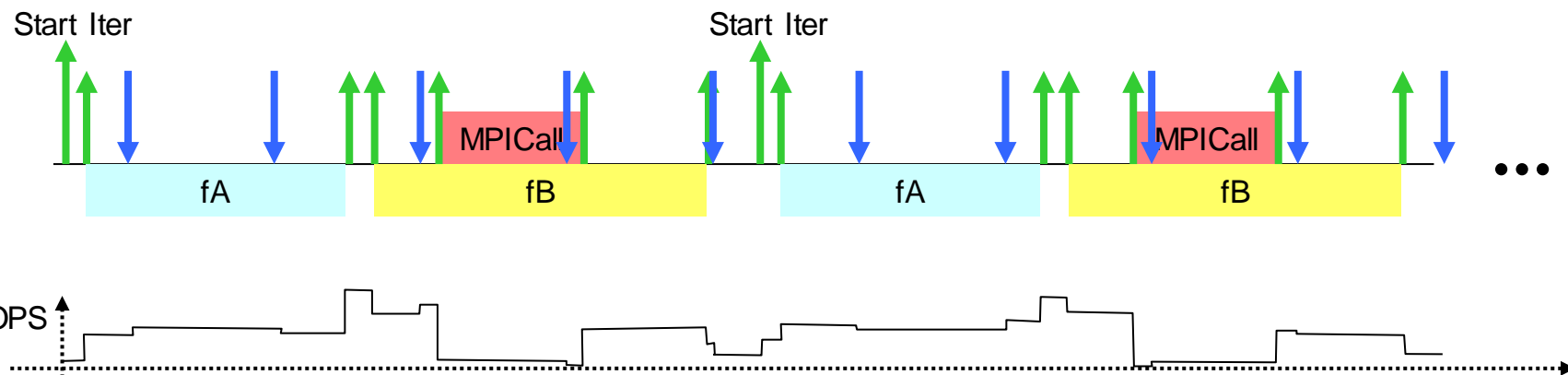
- Sufficiently large to minimize overhead
- Guaranteed acquisition interval/precision

Statistical projection

- %time (or metric) = $f(\%counts)$
- Assuming no correlation, sufficiently large #samples

Instrumentation + sampling

Both



Guaranteed interval

Captured data:

- Hardware counters (since previous probe)
- call path
- Call arguments in some probes