**Barcelona
Supercomputing
Center**
*Centro Nacional de Supercomputación*

# BSC Tools update
## Using clustering and Folding

Judit Gimenez (judit@bsc.es)

CScADS – Snowbird, June 2012

# Outline

**《 Computation Structure detection**

– Short intro

– Aggregative Refinement

– Tracking program evolution

– Scaling clustering algorithm

**《 Instantaneous performance metric**

– Clustering + Folding

**Barcelona**
**Supercomputing**
**Center**
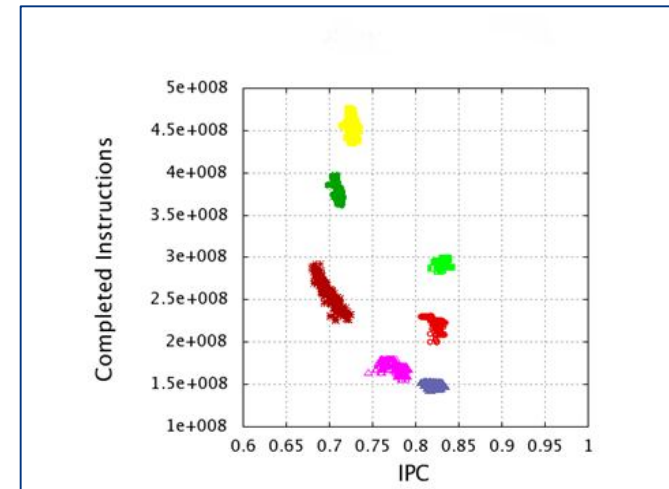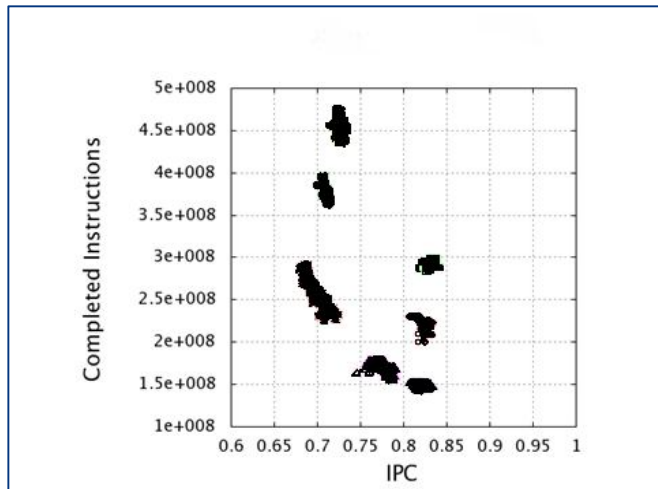Centro Nacional de Supercomputación

# Clustering

- 《 **Identification of computation structure**
  - CPU burst = region between consecutive runtime calls
    - Described with performance hardware counters
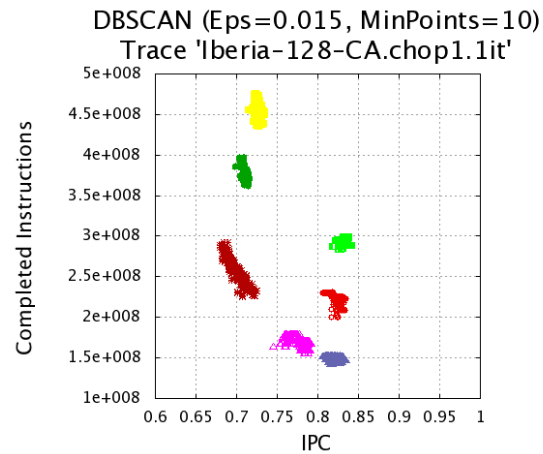    - Associated with call stack data
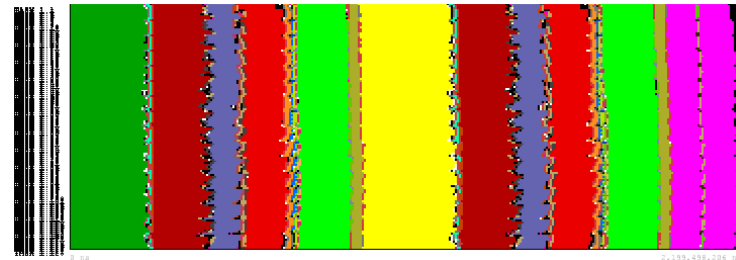- 《 **Using DBSCAN density-cluster algorithm**
  - Data not necessarily Gaussian

# Outputs

## Scatter Plot of Clustering Metrics



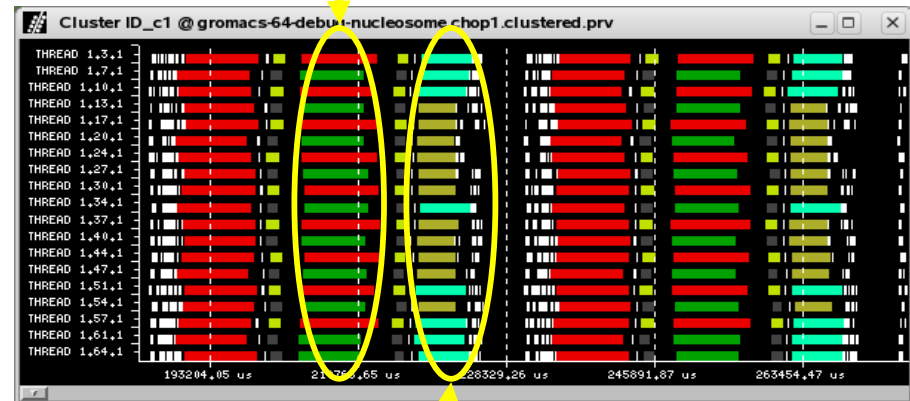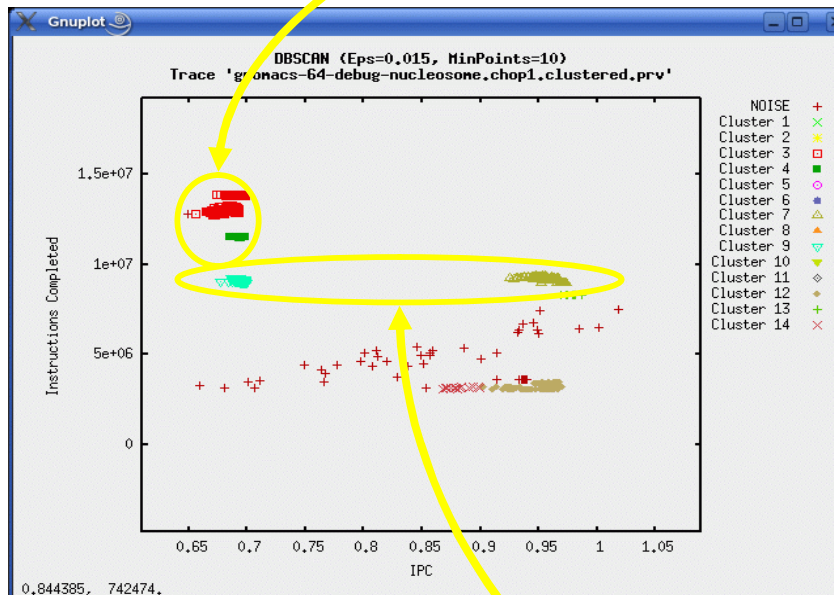## Clusters Distribution Along Time



## Cluster Statistics

| CLUSTER | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| % TIME | 36.29 | 29.52 | 10.13 | 9.68 | 3.73 | 1.71 |
| AVG. BURST DUR. (MS) | 220.46 | 177.70 | 60.81 | 29.09 | 38.71 | 44.83 |
| IPC | 0.53 | 0.50 | 0.62 | 0.77 | 0.66 | 0.59 |
| MIPS | 1210.07 | 1164.36 | 1403.19 | 1743.32 | 1499.47 | 1338.24 |
| L1M/KINSTR | 22.72 | 32.63 | 12.65 | 8.39 | 16.12 | 6.86 |
| L2M/KINSTR | 0.59 | 1.23 | 1.08 | 0.61 | 1.23 | 1.73 |
| MEM.BW (MB/S) | 90.77 | 182.65 | 193.32 | 136.33 | 236.15 | 295.71 |

## Code Linking

| CLUSTER | CODE SECTION |
|---|---|
| 1 | solve_nmm.f:[2037 - 2310] |
| 2 | solve_nmm.f:[1478 - 1782] |
|   | solve_nmm.f:[2030 - 1782] |
| 3 | solve_nmm.f:[1241 - 1345] |
| 4 | solve_nmm.f:[2771 - 2865] |
|   | solve_nmm.f:[2388 - 2489] |
| 5 | solve_nmm.f:[1478 - 1569] |
| 6 | solve_nmm.f:[1607 - 1633] |

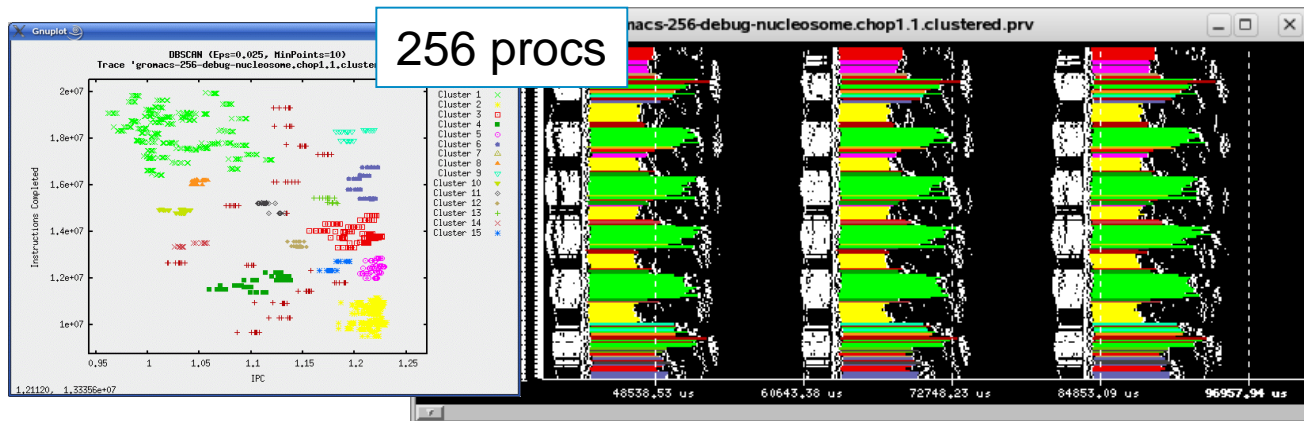# Using clusters to understand apps behavior (GROMACS)
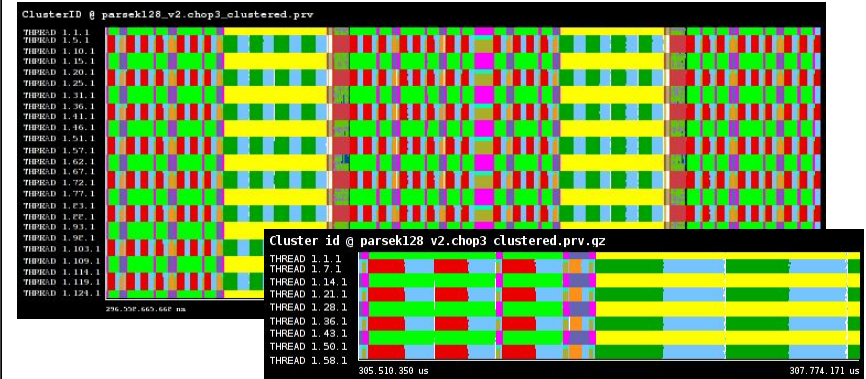
# Using clusters to understand apps behavior (GROMACS)
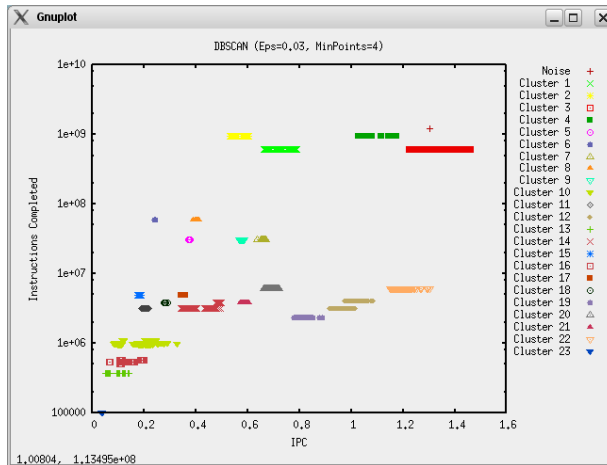


64 procs

256 procs

# Identifying main code regions (PARSEK)



duration vs. cluster

instr. vs. cluster

# Outline

**《 Computation Structure detection**

- Short intro
- **Aggregative Refinement**
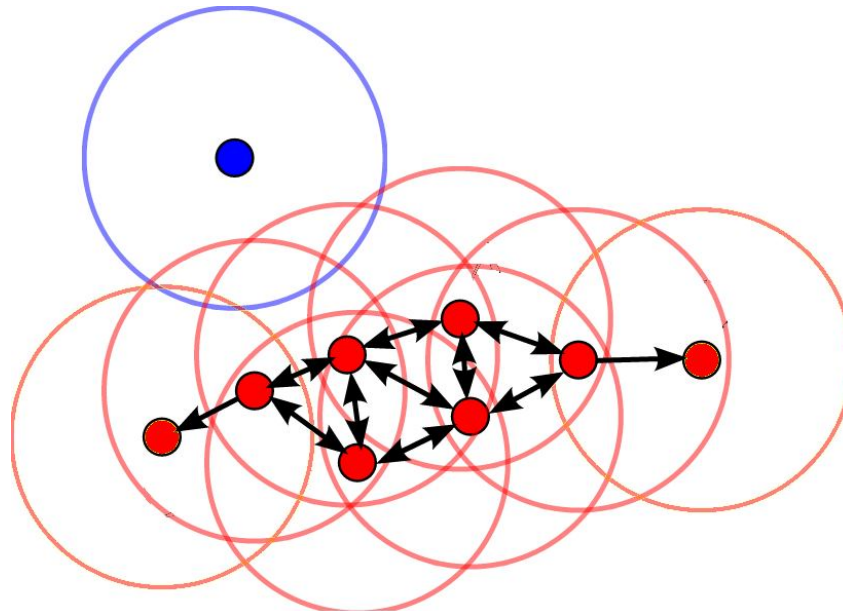- Tracking program evolution
- Scaling clustering algorithm

**《 Instantaneous performance metric**

- Clustering + Folding

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# DBSCAN characteristics

**Two parameters**

- Epsilon: search radius
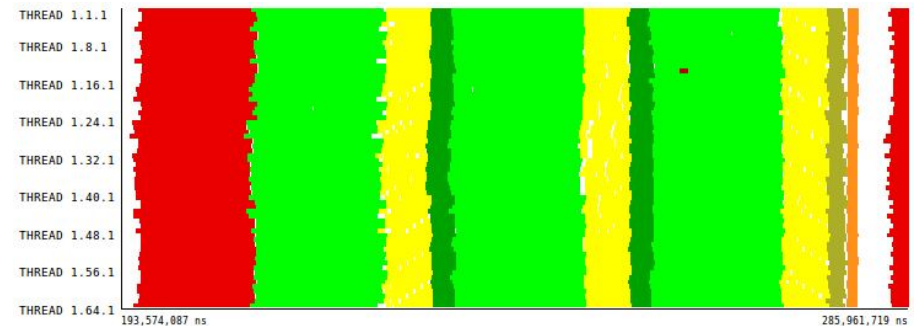- MinPoints: minimum cluster density



Noise point

Cluster points

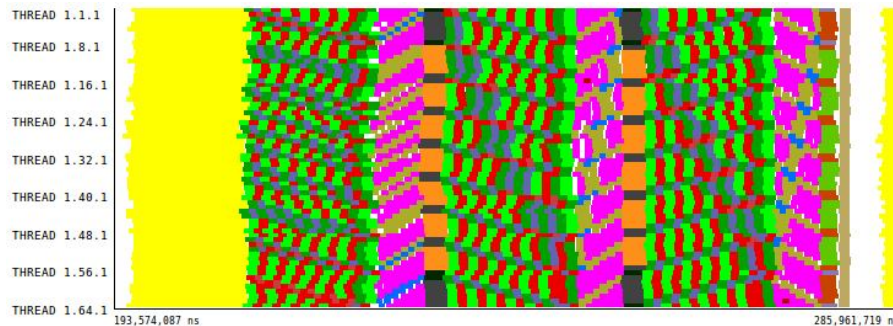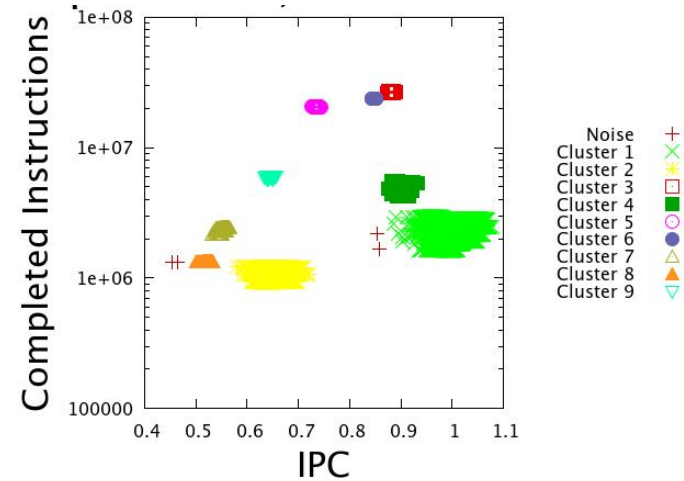# DBSCAN *Eps* selection

**《**  Which results are better?

### *Eps*=0.0140 (Low Value)



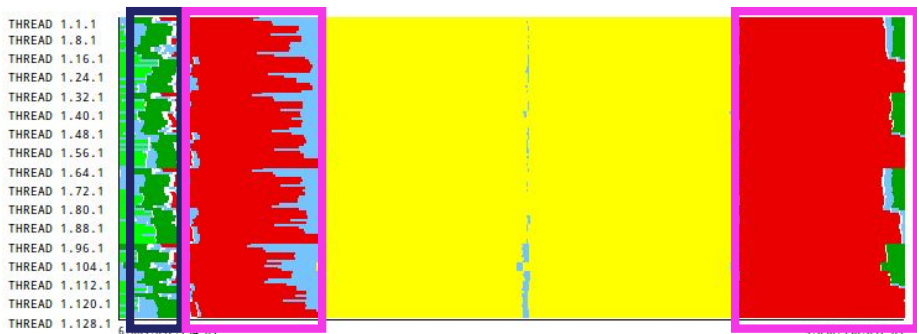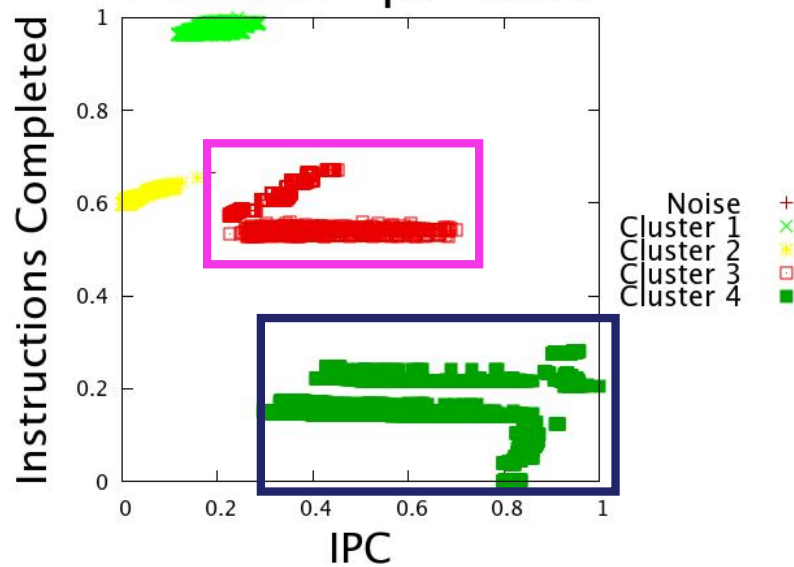### *Eps*=0.0400 (High Value)

# DBSCAN single *Eps* limitation

# Refinement Algorithm Approach
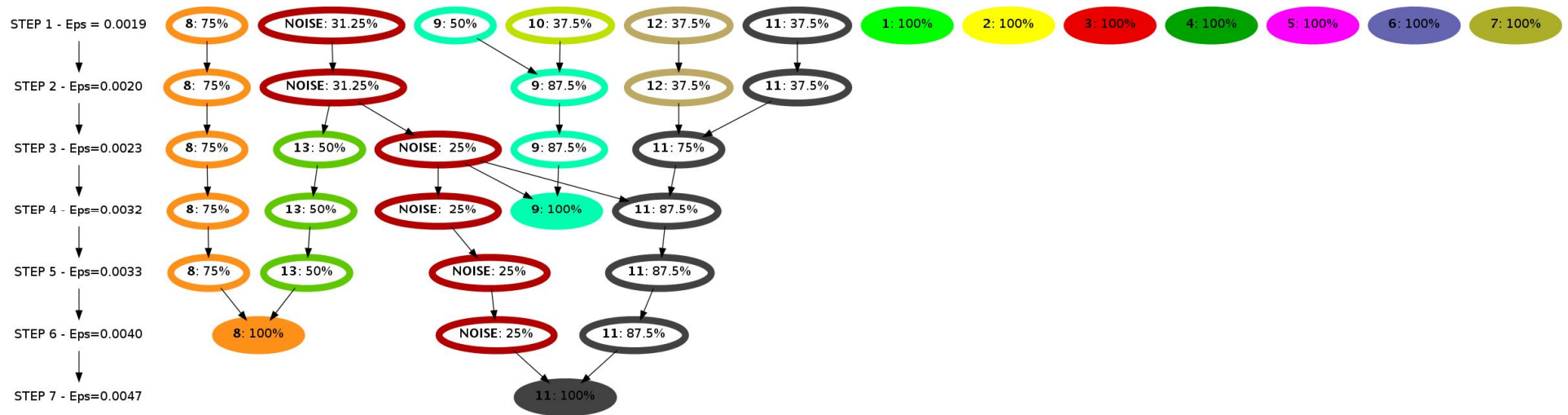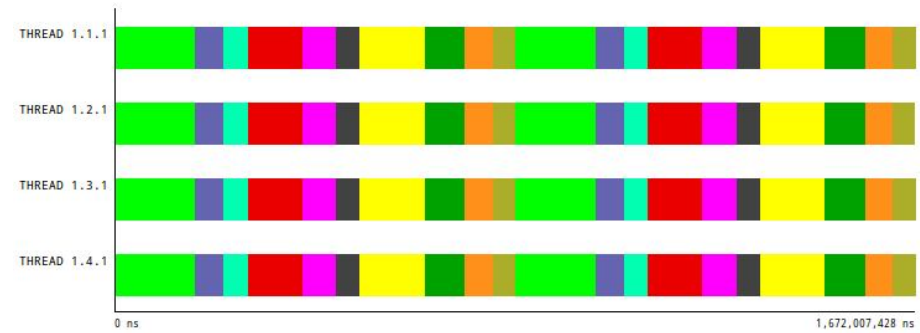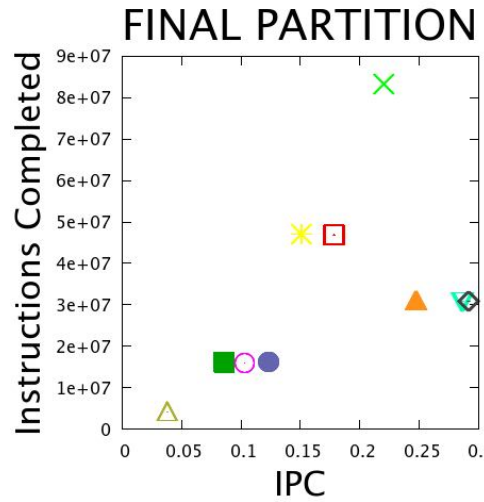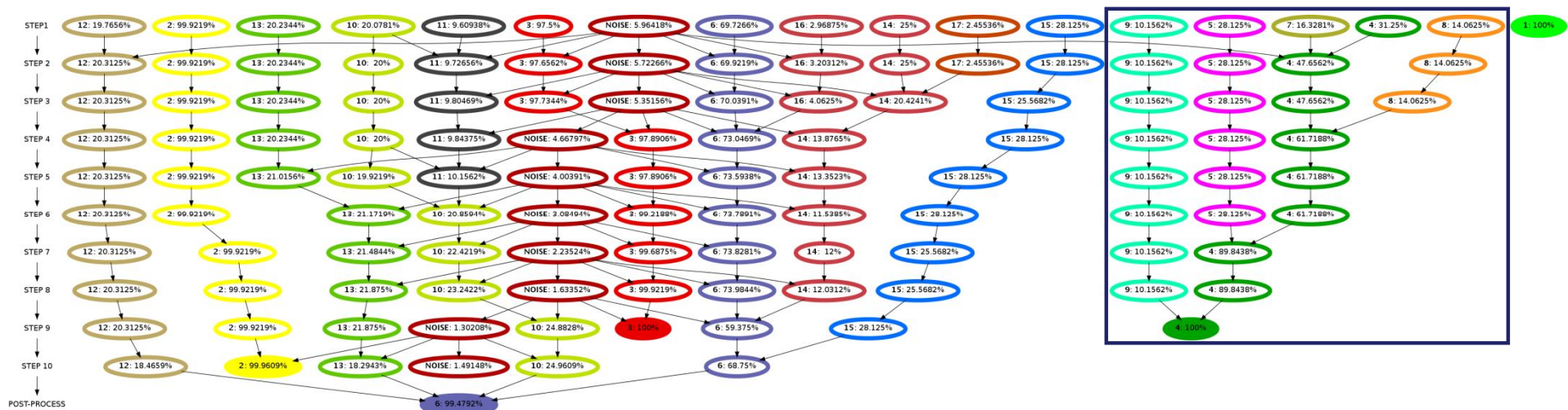
**《 Analogy between DBSCAN and hierarchical clustering**

    – Iterative bottom up construction of a *pseudo-dendogram*

**《 Cluster Sequence Score as target**

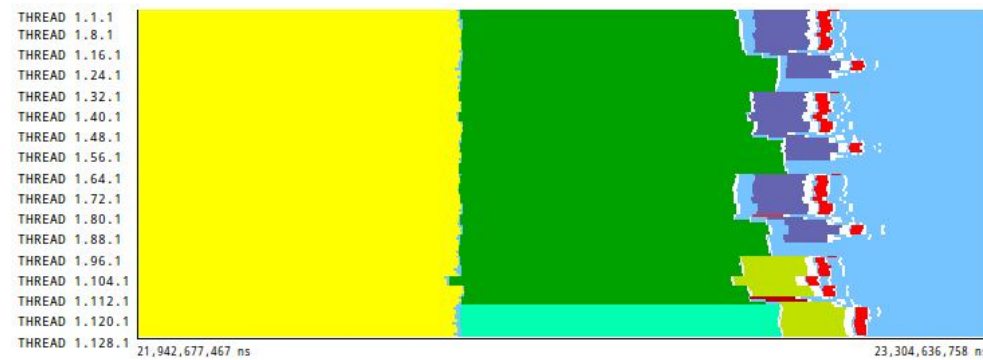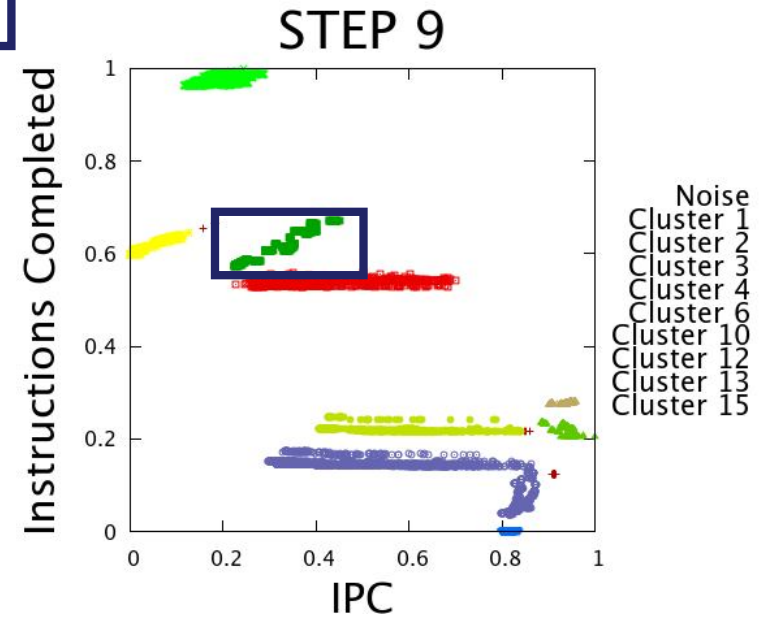    – Similar to *X*-means approach to decide *K*-means *k* parameter



*Automatic Refinement of Parallel Application Structure Detection (ICPADS 2011)*
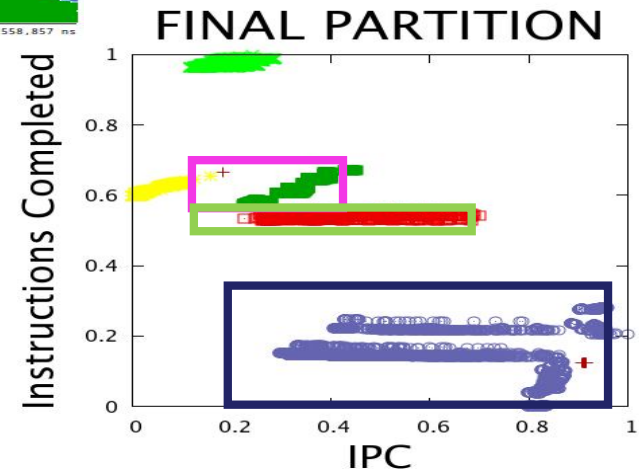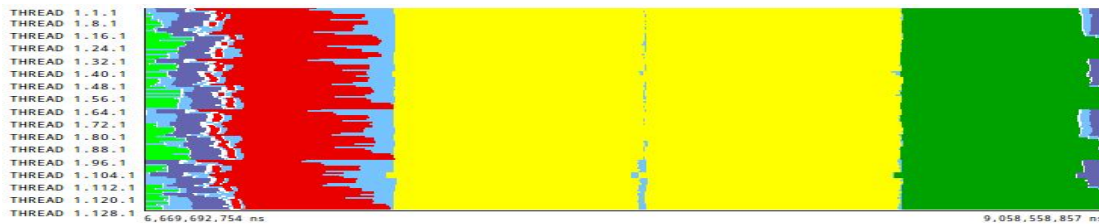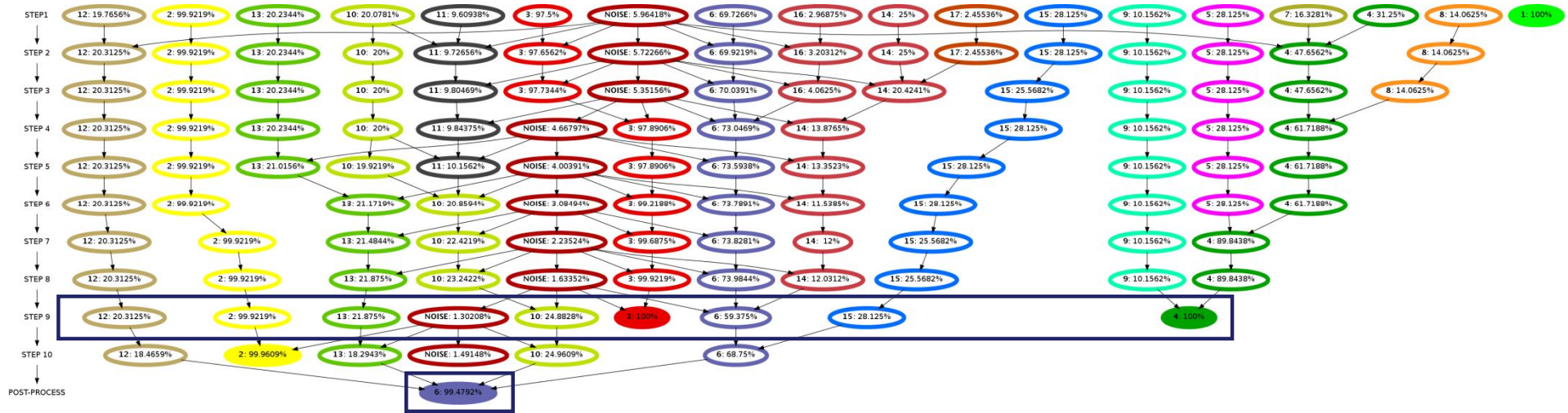
# BT A 4 tasks

# VAC4 128 tasks

# VAC4 128 tasks

# Outline

**《** Computation Structure detection

– Short intro

– Aggregative Refinement

– **Tracking program evolution**

– Scaling clustering algorithm

**《** Instantaneous performance metric

– Clustering + Folding

**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación

# Correlating multiple runs

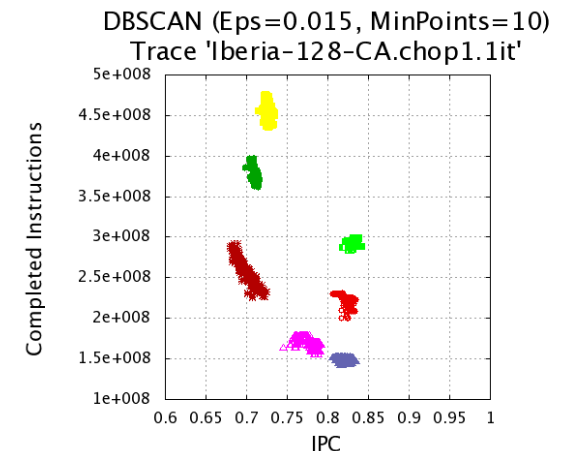**‖ Use and correlate information from different runs**

– Analysis of input parameters

– Code improvements

– Using different machines, compilers, flags, libraries

– Scalability studies

– Even for the same run: time evolution

**‖ Scatter plot = performance picture**

– Identifies objects and their weight

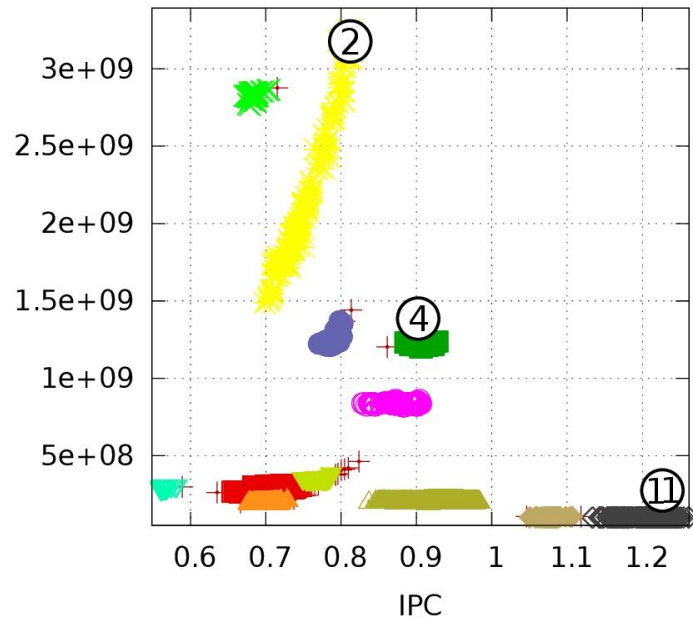– Correlation → image tracking

**‖ Based on heuristics**

– Code regions evolve smoothly (things keep closer)

– No common callstack means not the same region

– Time sequence identify regions within and between runs



DBSCAN (Eps=0.015, MinPoints=10)
Trace 'Iberia-128-CA.chop1.1it'
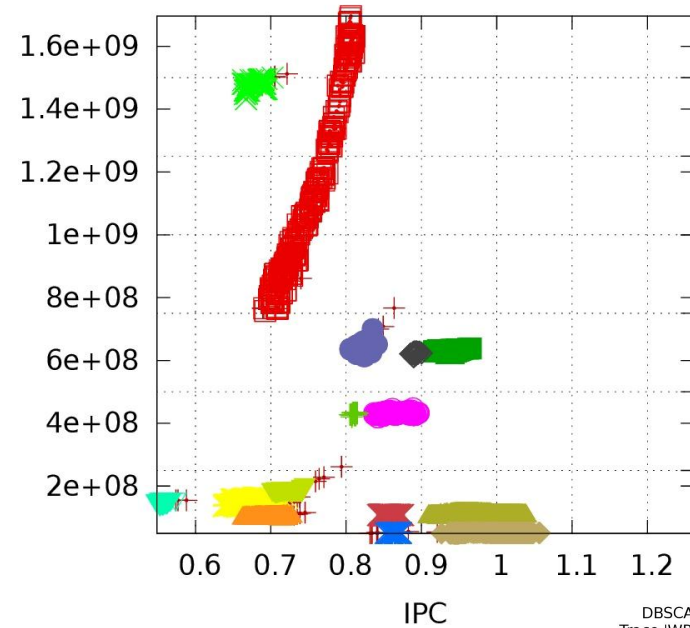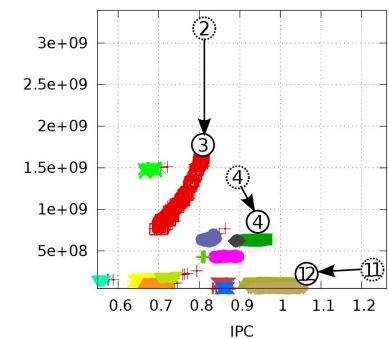
« 128 vs. 256



DBSCAN (Eps=0.018, MinPoints=10)
Trace 'WRF.MN.128p.chop2.clustered2.prv'

DBSCAN (Eps=0.018, MinPoints=10)
Trace 'WRF.MN.256p.chop2.clustered2.prv'

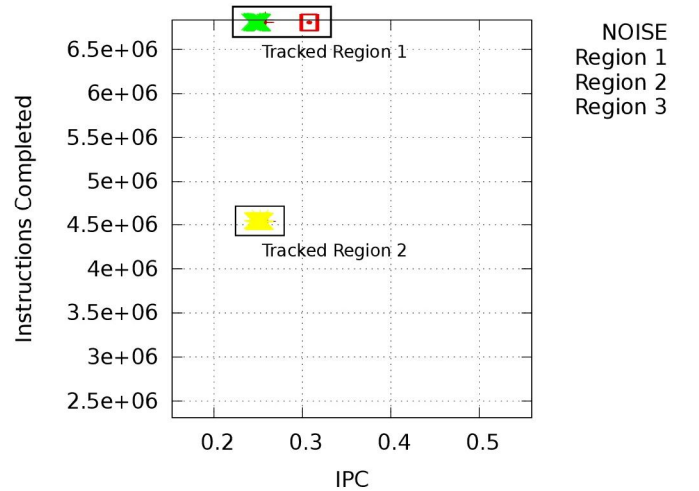DBSCAN (Eps=0.018, MinPoints=10)
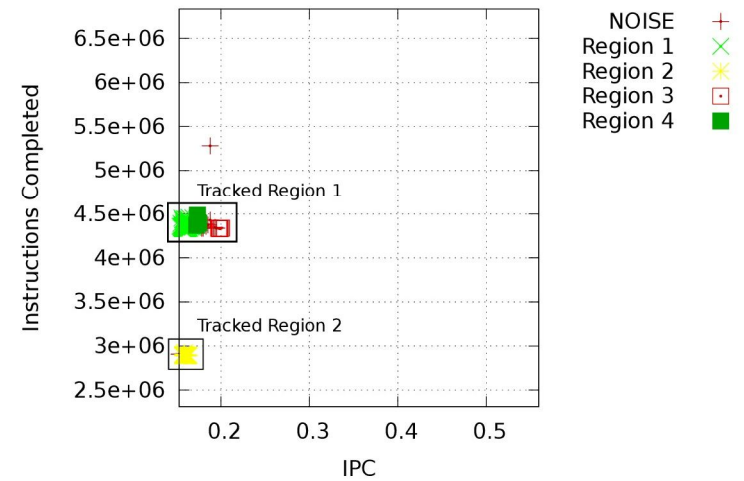Trace 'WRF.MN.256p.chop2.clustered2.prv'

# Scenario 2: Comparing machines & compilers (CG-POP)



## PowerPC, gfortran
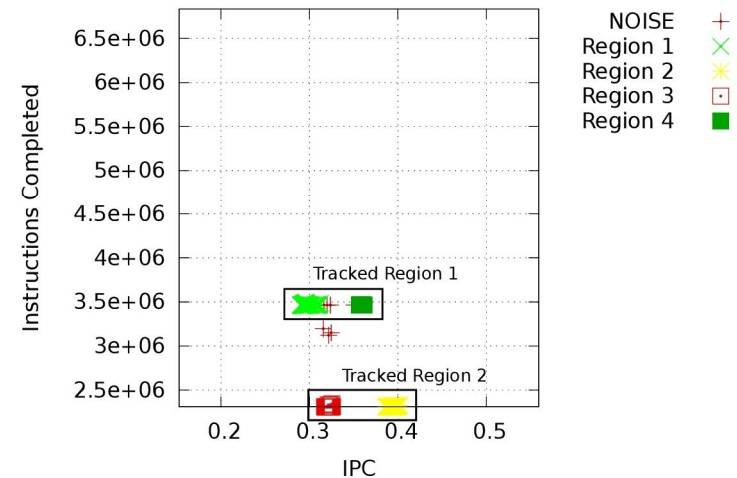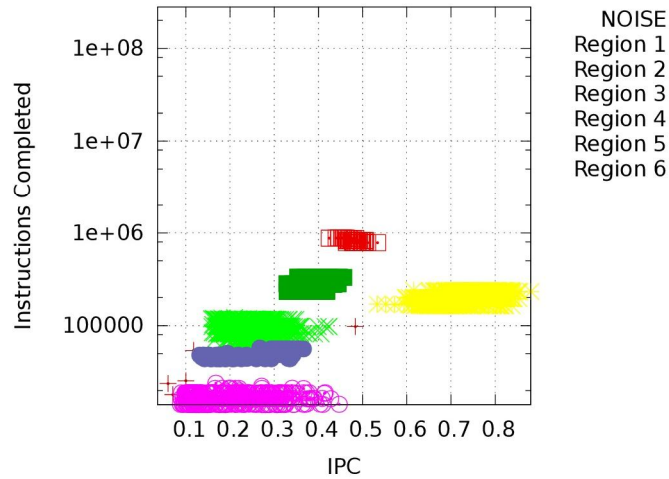
## PowerPC, xlc

## Intel, gfortran

## Intel, ifort

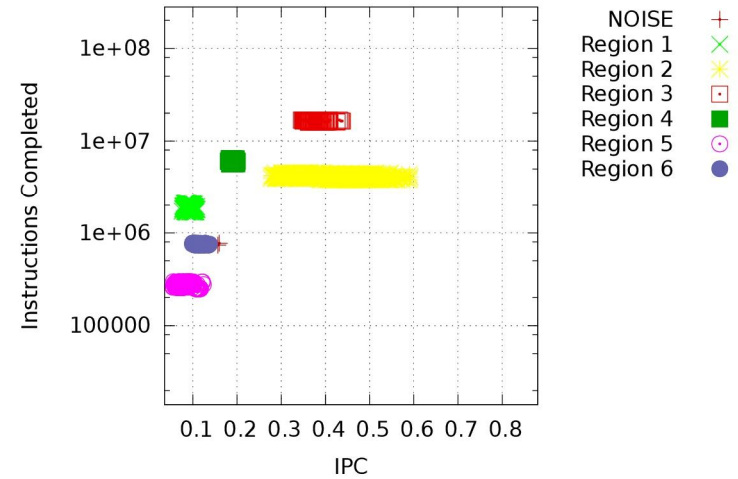# Scenario 3: Problem size impact (NAS-BT)

# Outline

**《** Computation Structure detection

- Short intro
- Aggregative Refinement
- Tracking program evolution
- **Scaling clustering algorithm**

**《** Instantaneous performance metric

- Clustering + Folding

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# Sampling input data

**《 First target: online clustering**

– Centralised approach (global clustering at the MRNet frontend)

– Data reduction trough sampling

– Data classification based on the samples clustering

All processes



32 representatives (50%)



15% random records



8 representatives + 15% random



25% random records



75% less data
6s down from 2m

# Hierarchical DBSCAN

1. **Local clustering**
   - Up to 20-30k points per local process
2. **Generate models**
   - Convex hull, medial axis…
3. **Merge the hulls over the MRNet**
   - Intersect?
4. **Broadcast the global model**
5. **Classify data locally using the global model**
   - Point inside the hull?



PEPC 4K tasks, 3095134 points, 273 tasks (16 way tree, 256 leaves, 12k points per local clustering) → clustering time 28.6 sec

# Comparison (parsek)

## 31515 points



Bar chart legend:
- PARSEK Total time
- PARSEK Clustering time

Categories: Sequential, Sampling 25%, Parallel 2+1 tasks, Parallel 4 +1tasks

## Sampling 25%

# Comparison (WRF)

74240 points
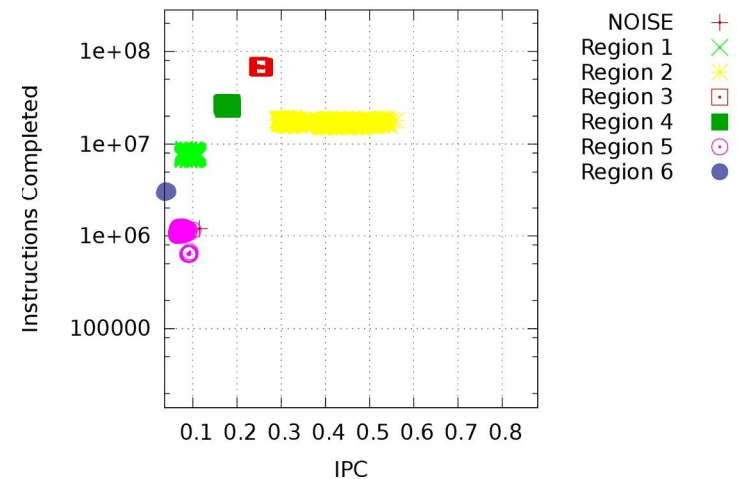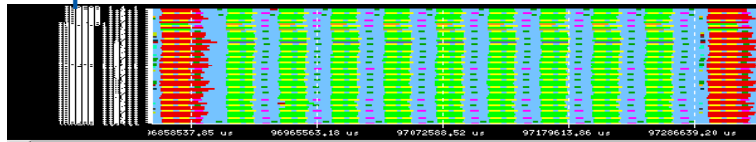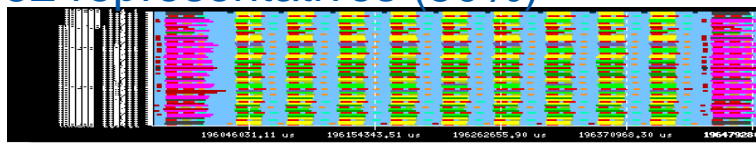


## Sequential /Par (4+1)



## Parallel (2+1)



## Sampling 25%

# Outline

**《** Computation Structure detection

– Short intro

– Aggregative Refinement

– Tracking program evolution

– Scaling clustering algorithm

**《 Instantaneous performance metric**

– **Clustering + Folding**

**Barcelona**
**Supercomputing**
**Center**
Centro Nacional de Supercomputación

# Can I get very detailed perf. data with low overhead?

**《 Application granularity vs. detailed granularity**

– Samples: hardware counters + callstack

**《 Folding: based on known structure: iterations, routines, clusters;**
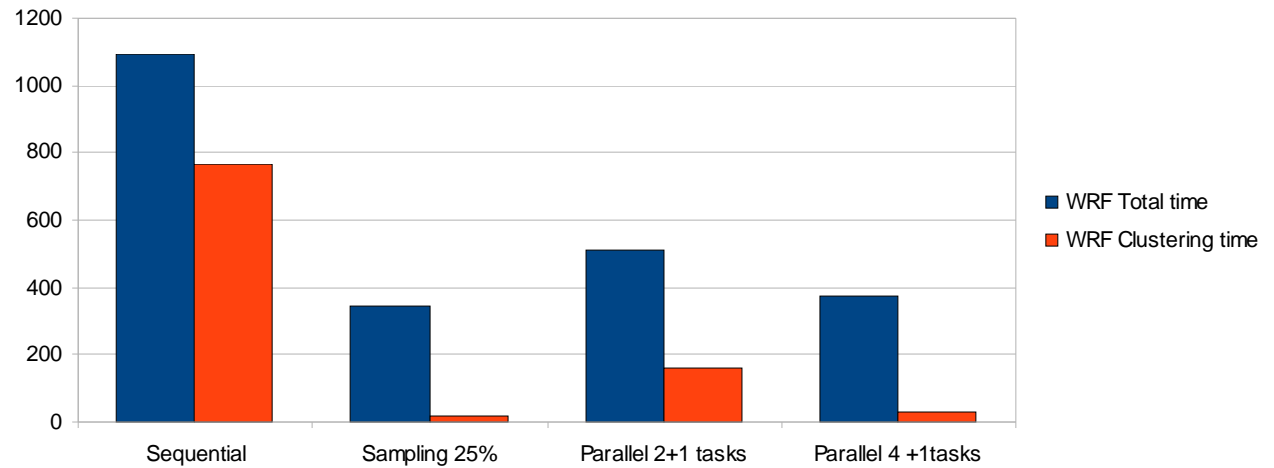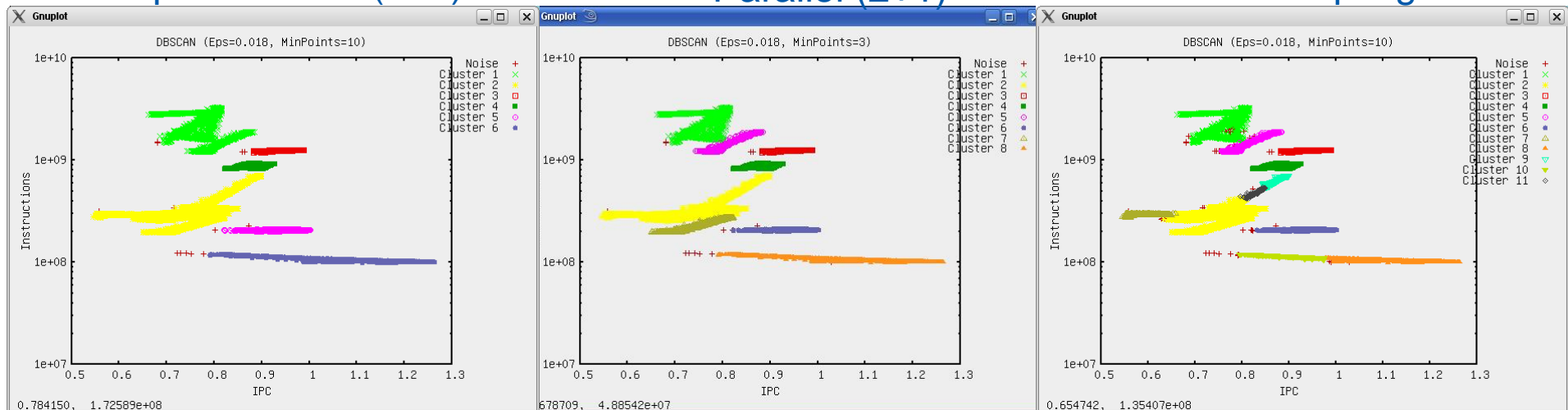
– Project all samples into one instance

**《 Extremely detailed time evolution of hardware counts, rates and callstack with minimal overhead**

– Correlate many counters
– Instantaneous CPI stack models



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

# Correlating with sources: which line should I look?



**Folded source code line**

**+**

**Folded instructions**

**=**

# The "benefits" of Fortran 90 intrinsic (PEPC)

DBSCAN (Eps=0.005, MinPoints=10)
Trace 'pepc.sorted.chop1.clustered.prv'



**96 MIPS**

Cluster 1
Cluster 2
Cluster 3
Cluster 4
Cluster 5
Cluster 6
Cluster 7
Cluster 8
Cluster 9
Cluster 10



THREAD 1.1.1

84.853.148 us          125.969.196 us

| Performance metrics |
|---|
| 16 MIPS |
| 2.3 M L2 misses/s |
| 0.1 M TLB misses/s |

Cumulative sample value within cluster 4 of PEPC



Total instructions
L2 cache misses
D-TLB cache misses

```
htable%node = 0
htable%key = 0
htable%link = -1
htable%leaves = 0
htable%childcode = 0
```

```
do i = 1, n
  htable(i)%node = 0
  htable(i)%key = 0
  htable(i)%link = -1
  htable(i)%leaves = 0
  htable(i)%childcode = 0
End do
```

| Changes |
|---|
| -70% time |
| -18% instructions |
| -63% L2 misses |
| -78% TLB misses |
| 253 MIPS (+163%) |

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Interchanging loops (MR. GENESIS)



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

# Pre-computing float data – loop split (PMEMD)



PowerPC model evolution for pmemd.pme.clustered.fused.extract.1.1.Cluster_1.0

PowerPC model evolution for pmemd.pme.optim.clustered.fused.extract.1.1.Cluster_1.0

Time (ms)

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

Barcelona
Supercomputing
Center
Centro Nacional de Supe

# Conclusions

**《 Performance analytics**

– Data analytics applied to raw performance data

– From data to insight

  • Information is on variability and distribution

– Huge room for research

**《 Showed results of some techniques**

– Clustering enables focusing the analysis and open many different uses on the analysis

– Folding makes possible to compute instantaneous performance metric functions with low overhead

– Tracking helps detecting movement in the performance space

  • Sequence of "frames" along many factors (not just time)

www.bsc.es/paraver

**Barcelona**
**Supercomputing**
**Center**
Centro Nacional de Supercomputación