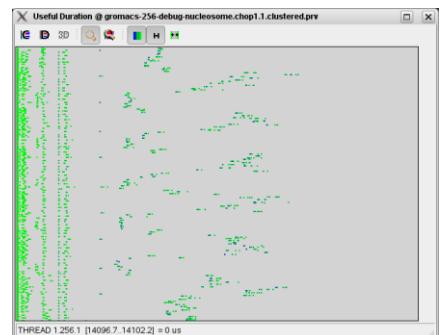
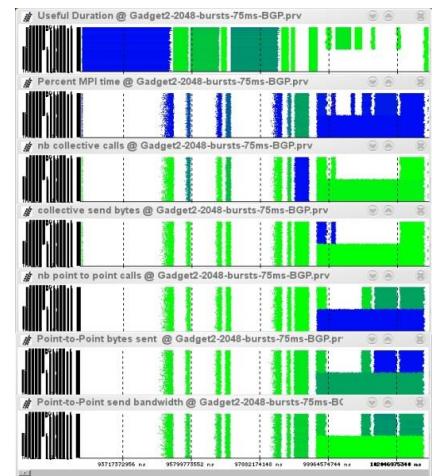
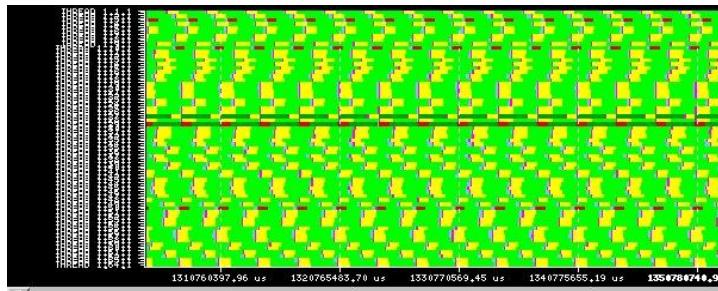


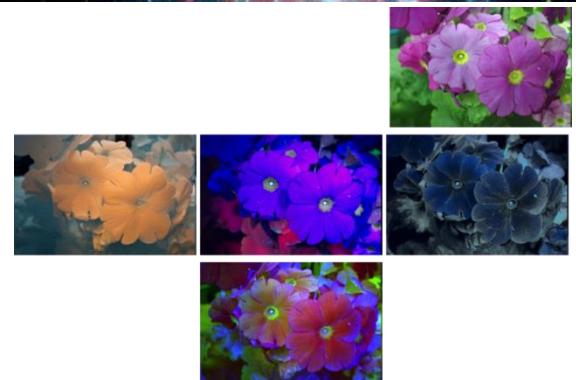
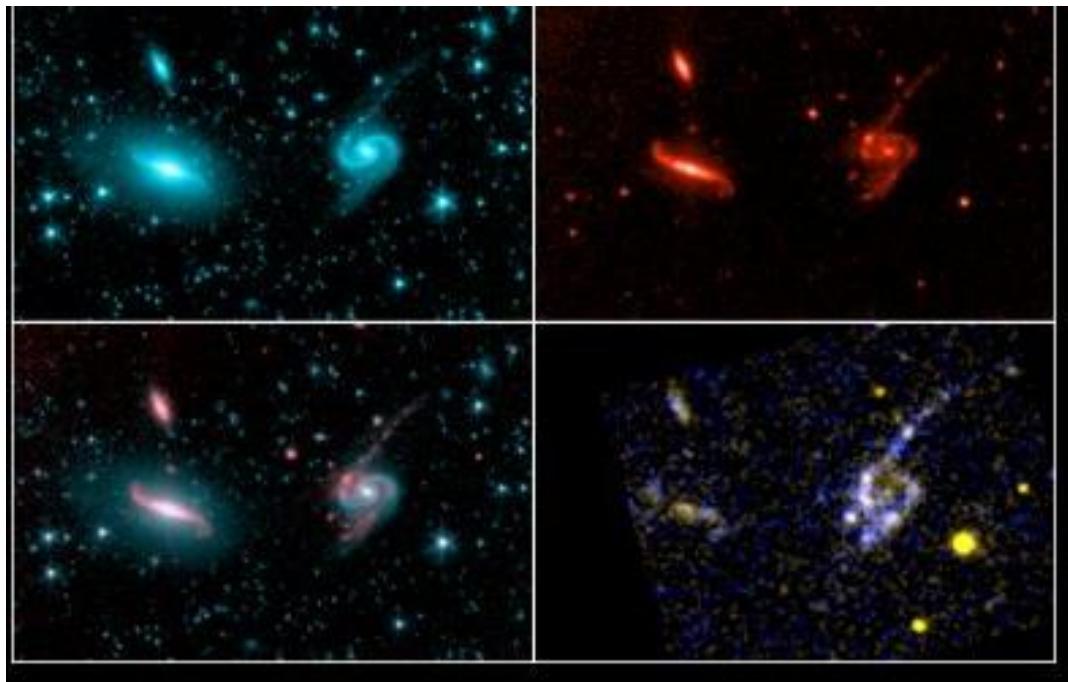
# Understanding applications with Paraver and Dimemas

Judit Giménez, Harald Servat  
[tools@bsc.es](mailto:tools@bsc.es)  
Oct 10th 2013

- Since 1991
- Based on traces
- Open Source
  - <http://www.bsc.es/paraver>
- Core tools:
  - Paraver (paramedir) – offline trace analysis
  - Dimemas – message passing simulator
  - Extrae – instrumentation
- Focus
  - Detail, flexibility, intelligence

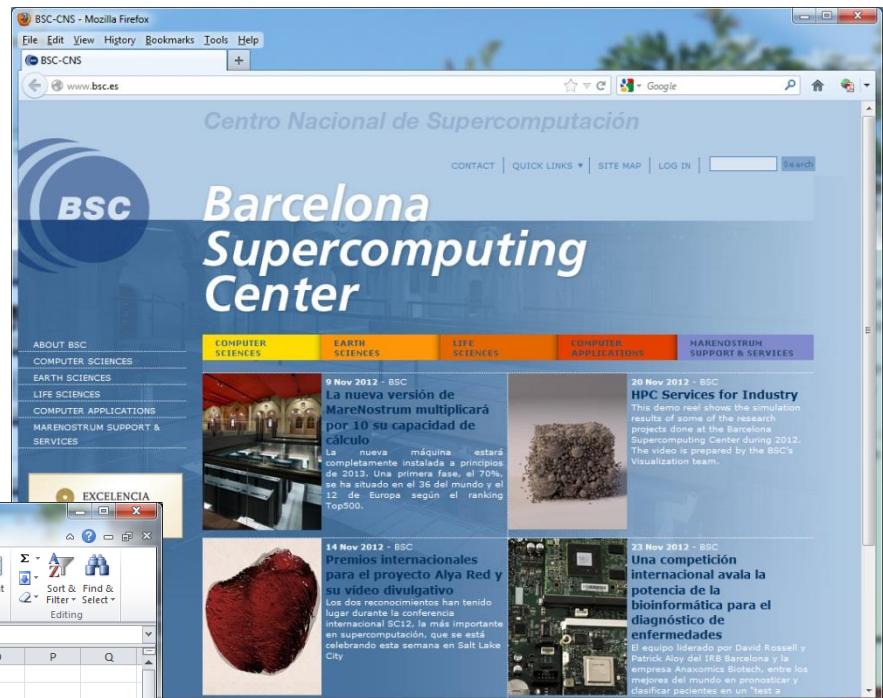
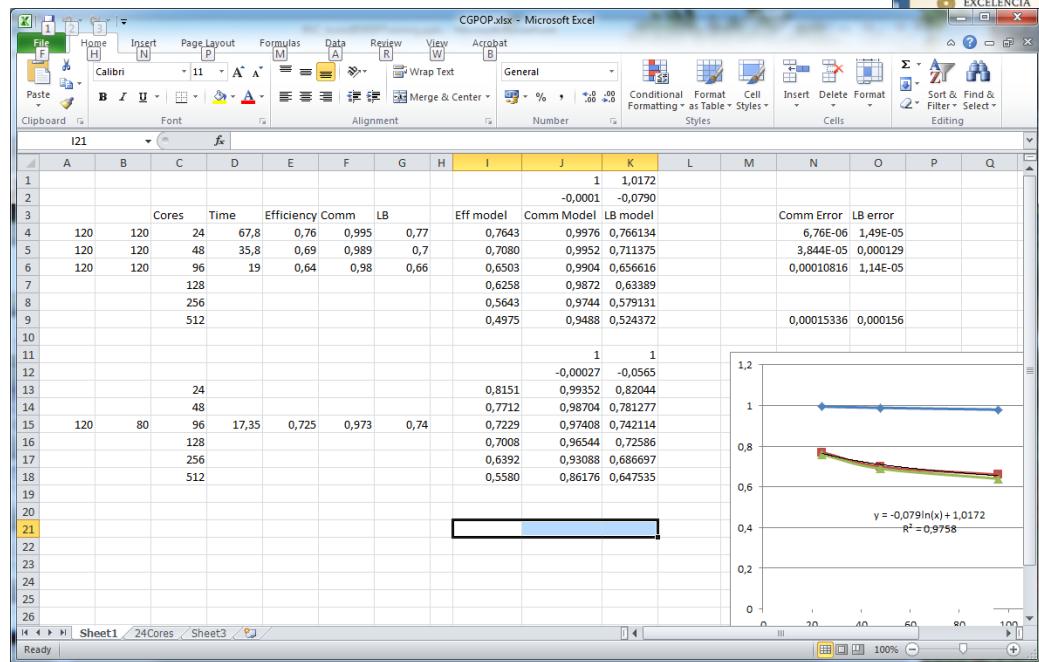


- Different looks at one reality
  - Different spectral bands (light sources and filters)
- Highlight different aspects
  - Can combine into false colored but highly informative images



# Spreadsheets and browsers

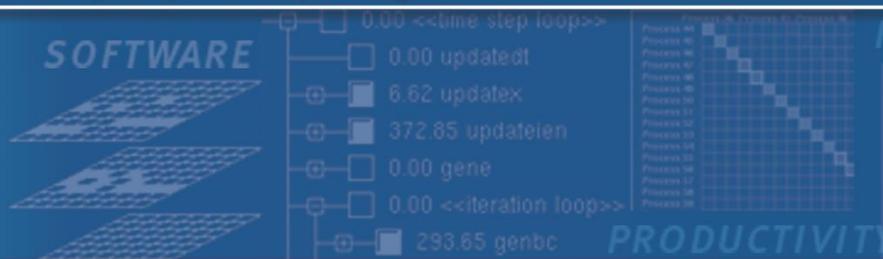
- Display, manipulate data
  - Dynamic content
  - User defined operations



- **Behavioral structure** vs. syntactic structure
  - Algorithmic and performance
  - In space and time
- **Variability**
  - Multimodal distributions
  - Variability + synchronization → critical non linear effects
- **Flexibility** to let analyst navigate the captured data and gain as much **insight** as possible from as **few application runs** as possible.



SOFTWARE

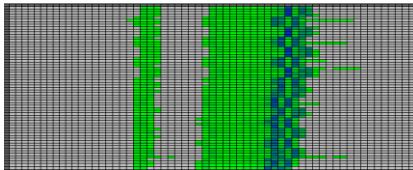
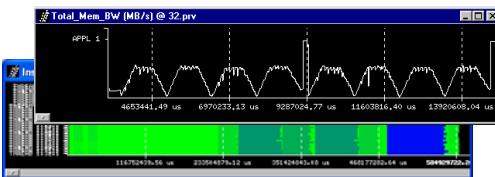
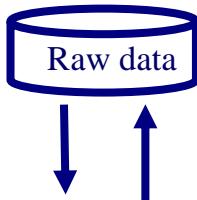


FAST SOLUTIONS

- PAPI\_L1\_DCM
- PAPI\_L1\_ICM
- PAPI\_L2\_DCM
- PAPI\_L2\_ICM
- PAPI\_L1\_TCM
- PAPI\_L2\_TCM

# 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**



## Timelines

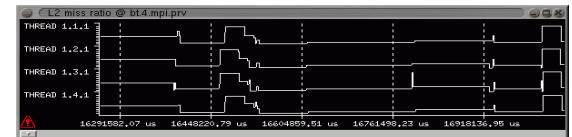
Trace visualization/analysis  
+ trace manipulation

Goal = Flexibility  
No semantics  
Programmable

## 2/3D tables (Statistics)

Comparative analyses  
Multiple traces  
Synchronize scales

- Each window displays one view
  - **Piecewise constant** function of time



$$s(t) = S_i, i \in [t_i, t_{i+1})$$

- Types of functions

- Categorical

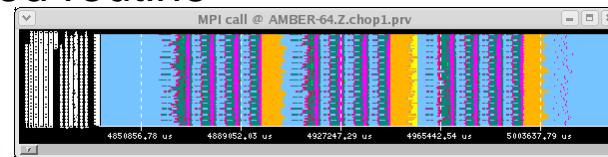
- State, user function, outlined routine

$$S_i \in [0, n] \subset N, \quad n <$$

- Logical

- In specific user function, In MPI call, In long MPI call

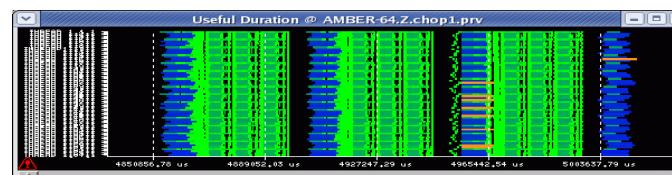
$$S_i \in \{0, 1\}$$



- Numerical

- IPC, L2 miss ratio, Duration of MPI call, duration of computation burst

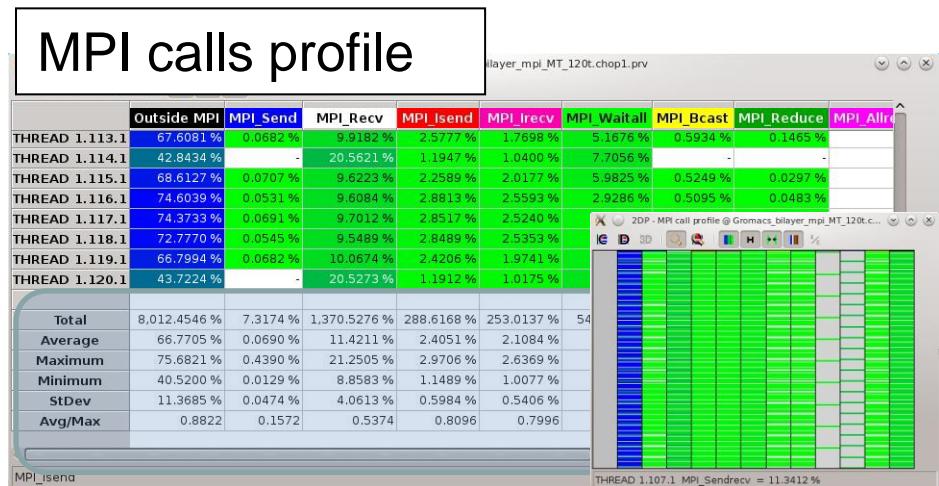
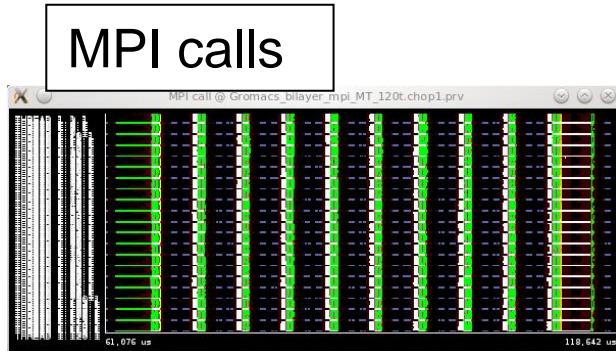
$$S_i \in R$$



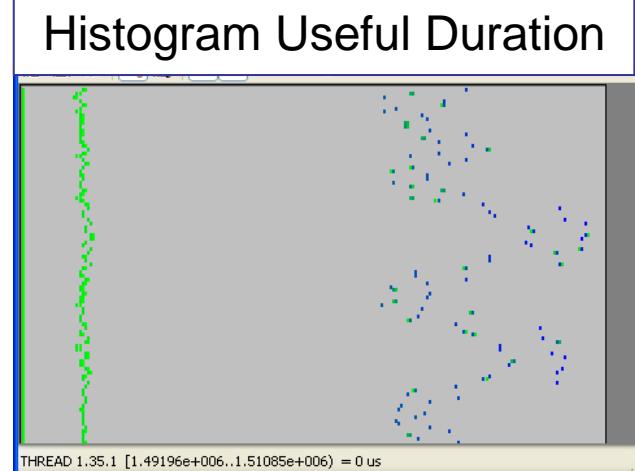
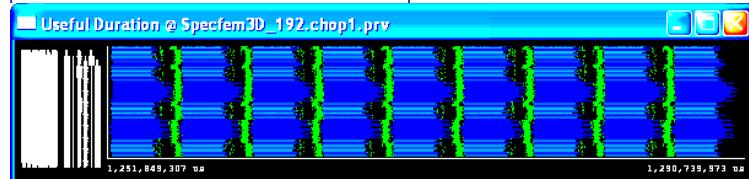
# Tables: Profiles, histograms, correlations

VI-HPS

- From timelines to tables

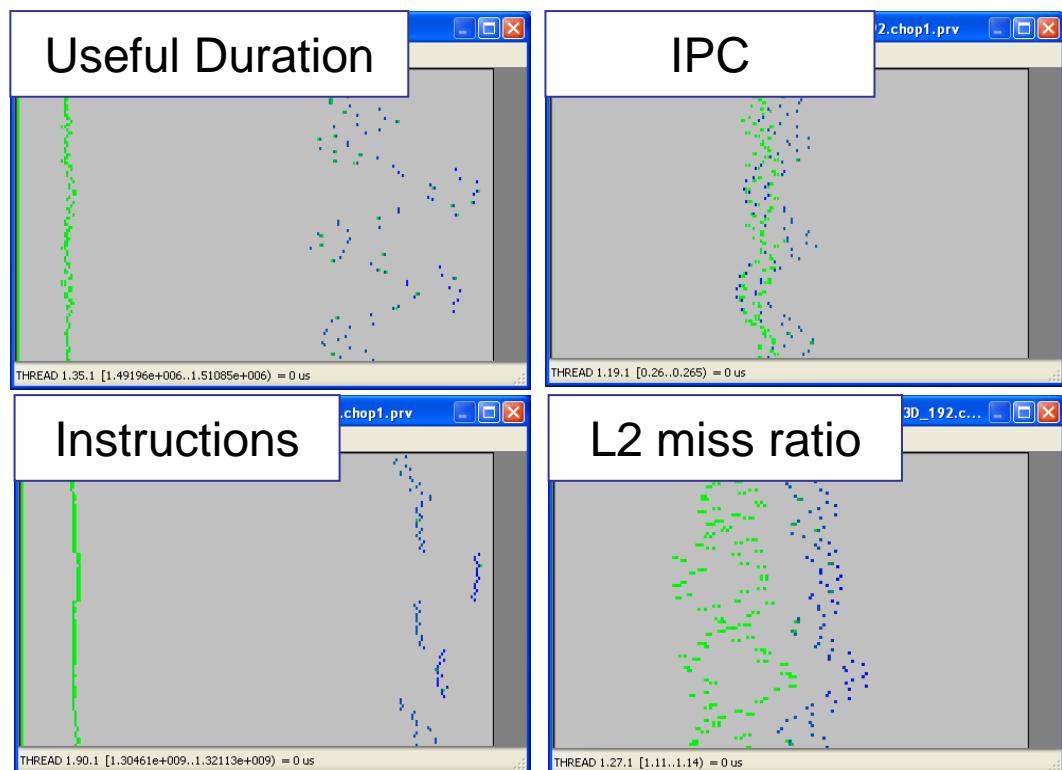
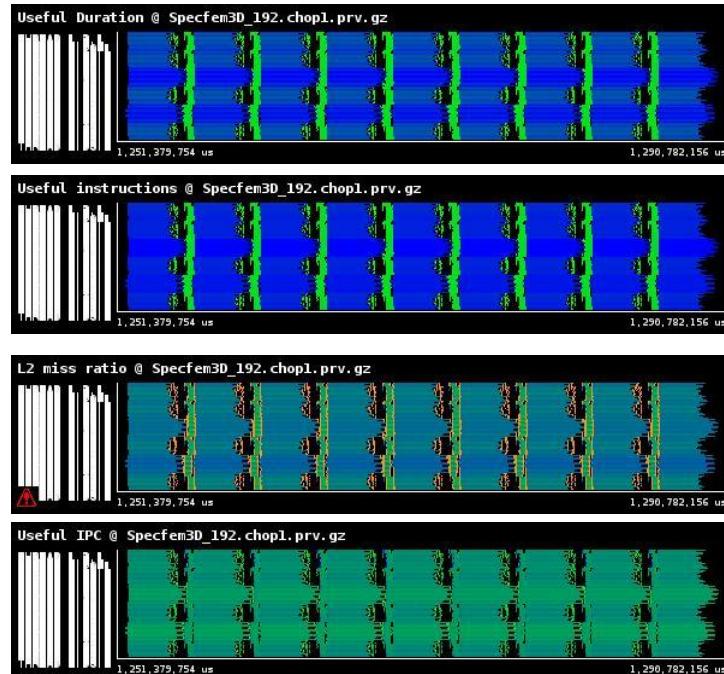


## Useful Duration



# Analyzing variability through histograms and timelines

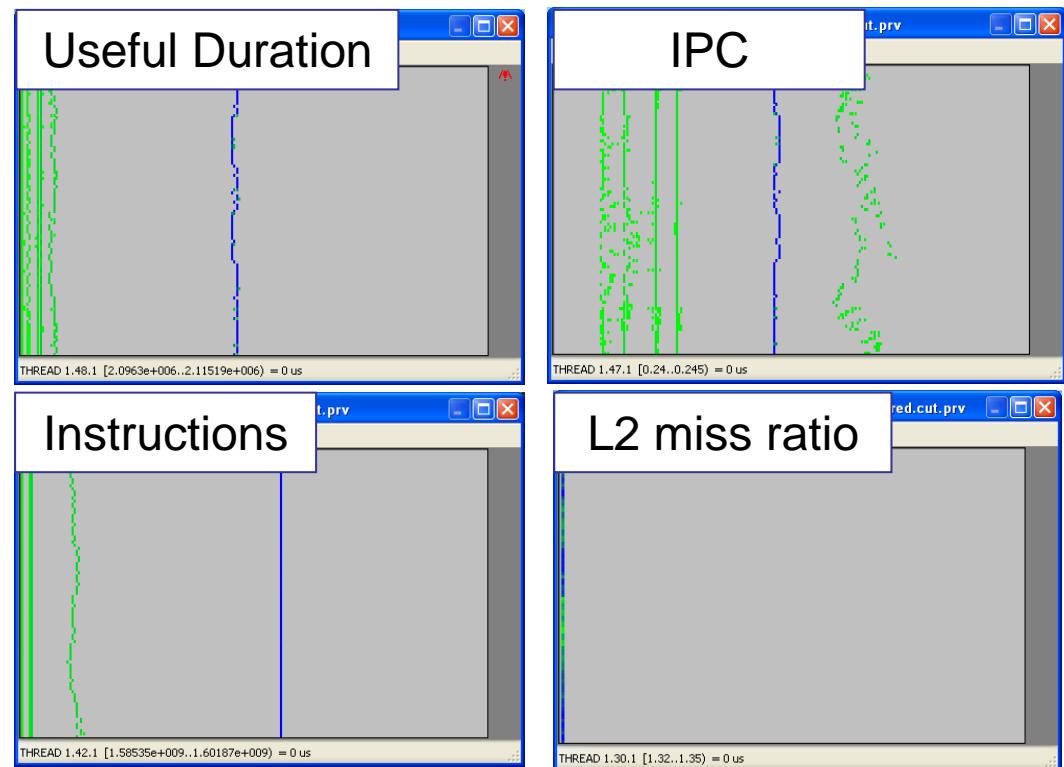
VI-HPS



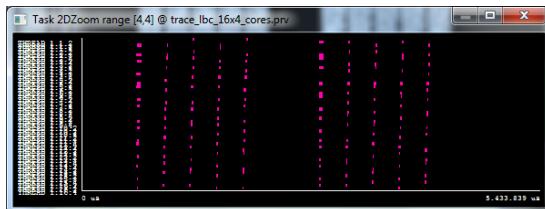
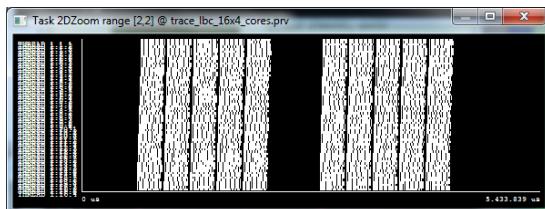
# Analyzing variability through histograms and timelines

VI-HPS

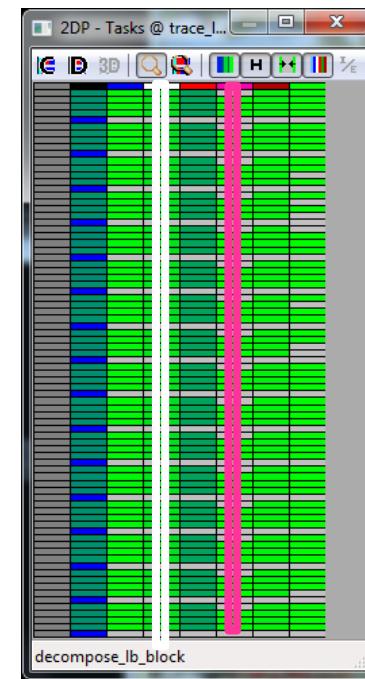
- By the way: six months later ....



- Where in the timeline do the values in certain table columns appear?
  - ie. want to see the time distribution of a given routine?

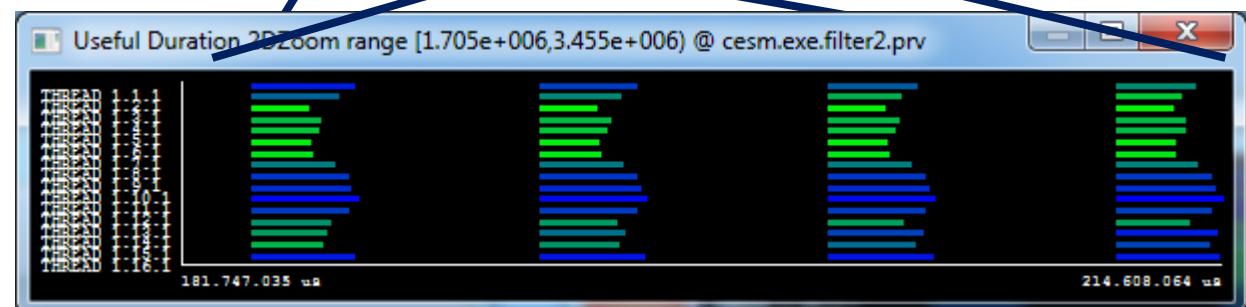
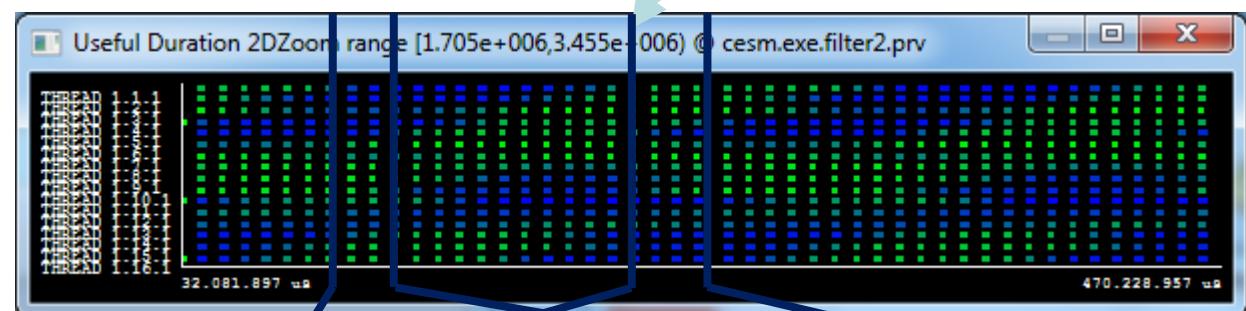
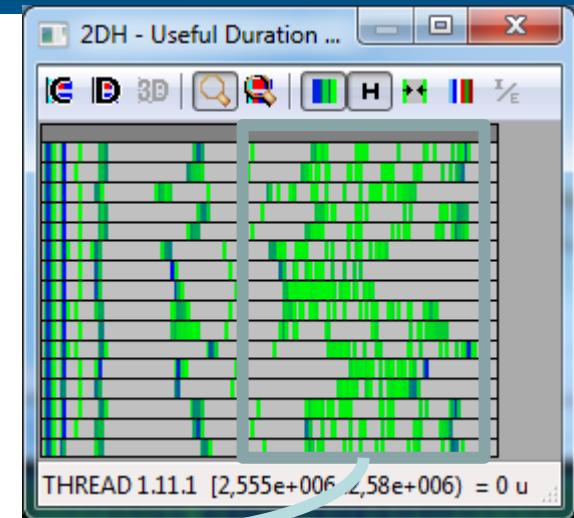


Only showing when a  
given value happens

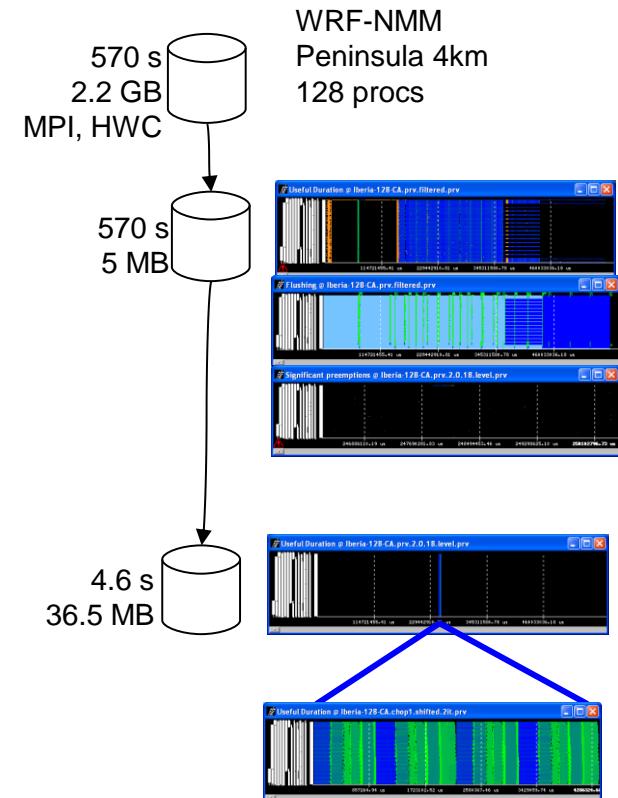


# Variability ... is everywhere

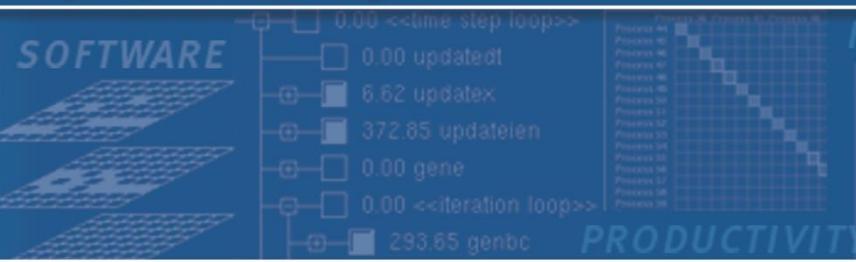
- CESM: 16 processes, 2 simulated days
- Histogram useful computation duration shows high variability
- How is it distributed?
- Dynamic imbalance
  - In space and time
  - Day and night.
  - Season ? ☺



- Data handling/summarization capability
  - Filtering
    - Subset of records in original trace
    - By duration, type, value,...
    - Filtered trace IS a paraver trace and can be analysed with the same cfgs (as long as needed data kept)
  - Cutting
    - All records in a given time interval
    - Only some processes
  - Software counters
    - Summarized values computed from those in the original trace emitted as new even types
    - #MPI calls, total hardware count,...



See slides at end of presentation for details



FAST SOLUTIONS

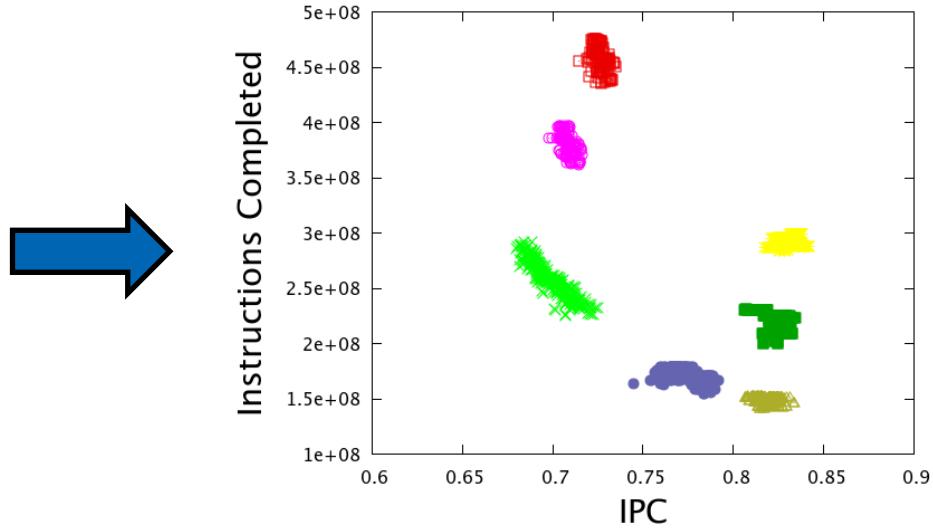
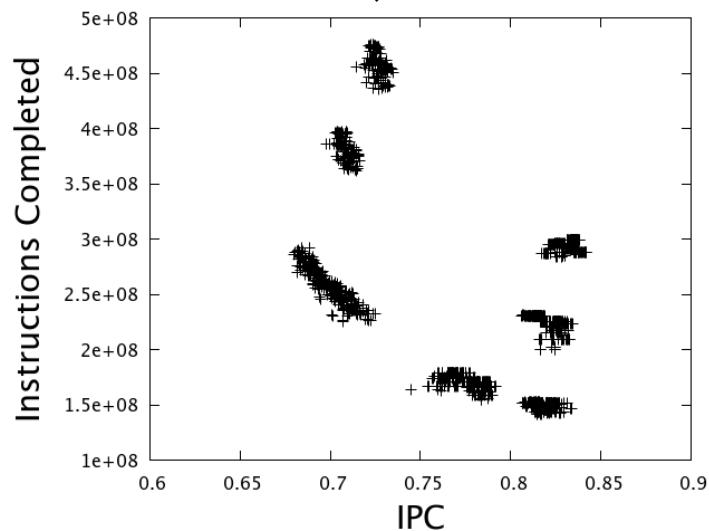
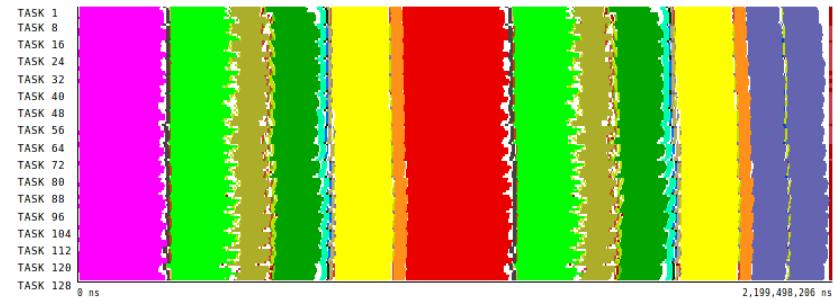
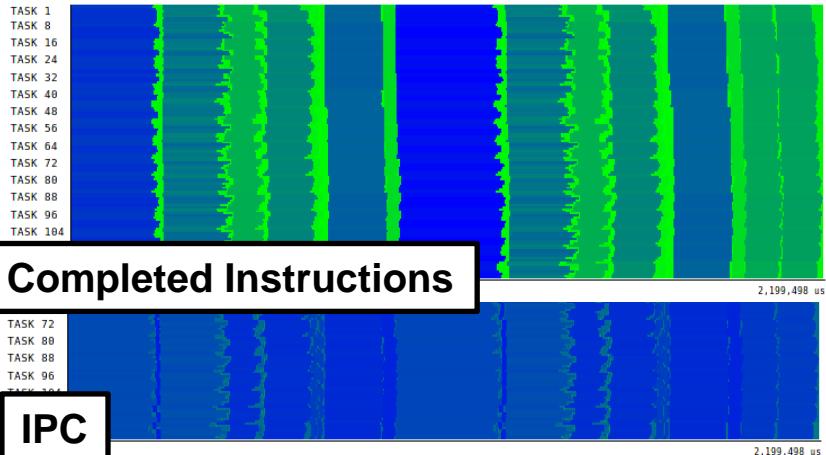
- PAPI\_L1\_DCM
- PAPI\_L1\_ICM
- PAPI\_L2\_DCM
- PAPI\_L2\_ICM
- PAPI\_L1\_TCM
- PAPI\_L2\_TCM

# Performance Analytics

# Clustering

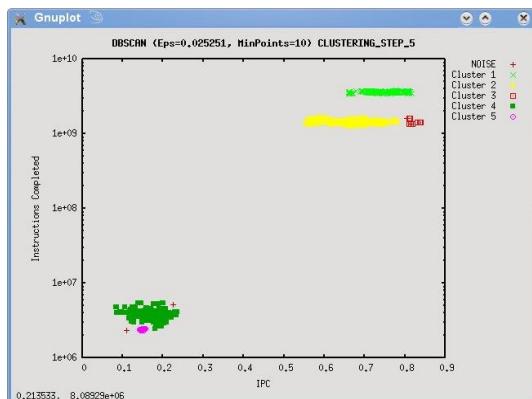
# Using Clustering to identify structure

VI-HPS

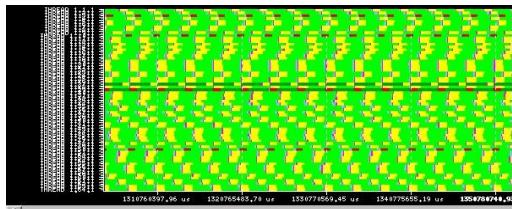


# Performance @ serial computation bursts

VI-HPS

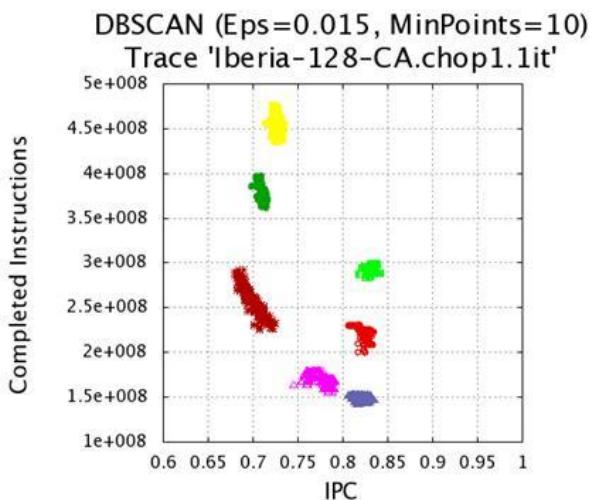


SPECFEM3D

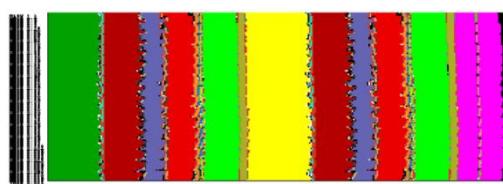


Asynchronous SPMD

Balanced #instr  
variability in IPC

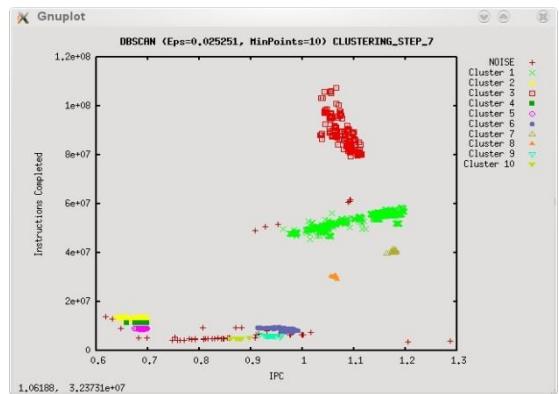


WRF 128 cores

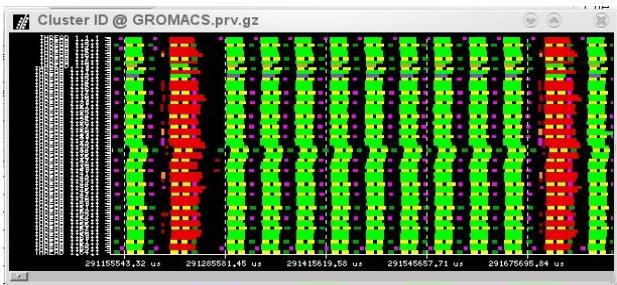


SPMD  
Repeated substructure

Coupled imbalance



GROMACS



MPMD structure

Different coupled  
imbalance trends

# Example PARSEK (DEEP)

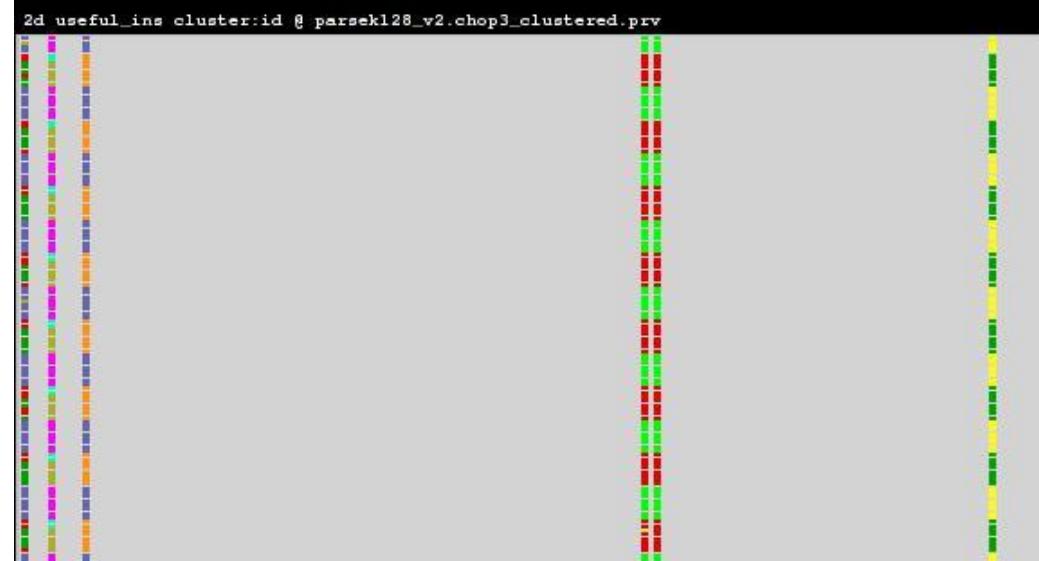
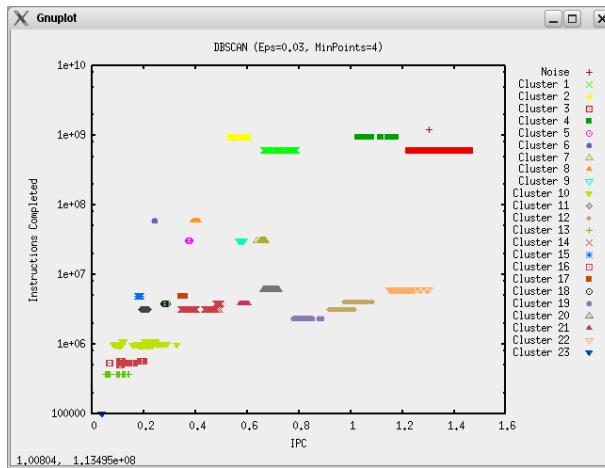
VI-HPS



duration vs. cluster



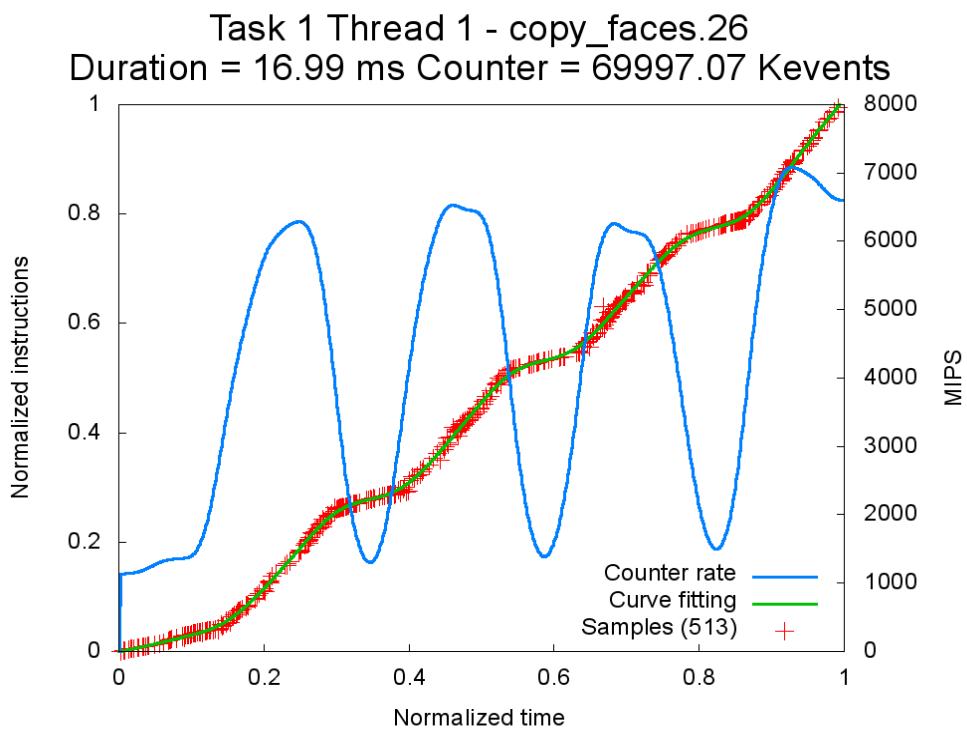
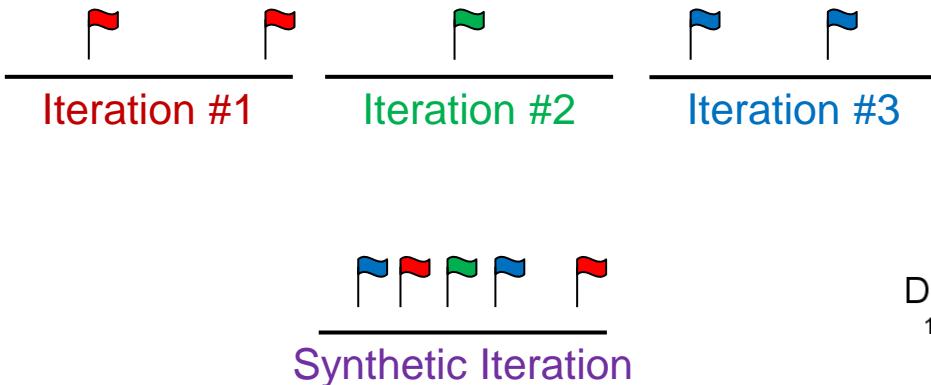
instr. vs. cluster



# Folding

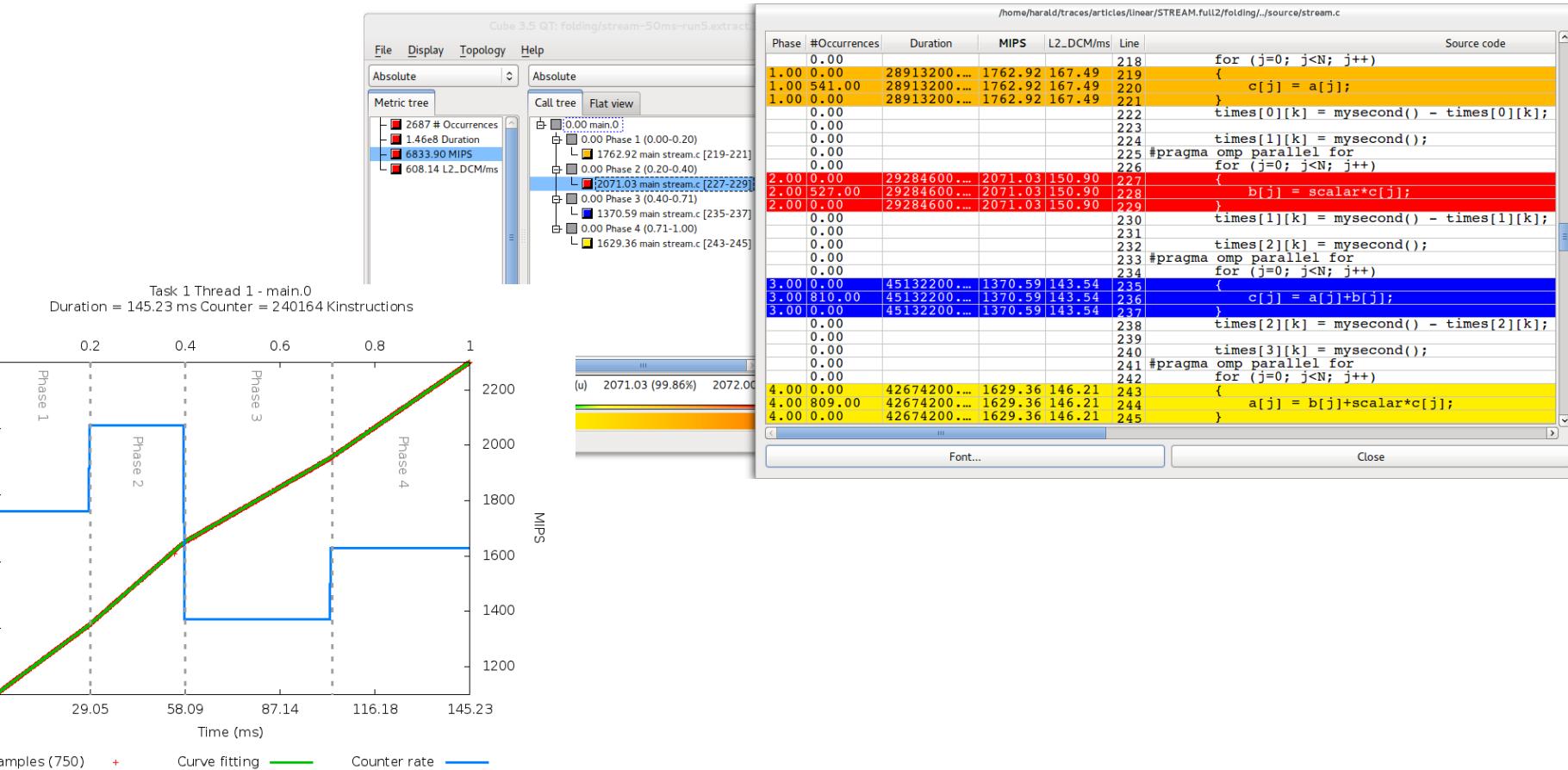
# Folding: Detailed time evolution

- Benefit from applications' repetitiveness



# Folding → profiles of rates and ratios

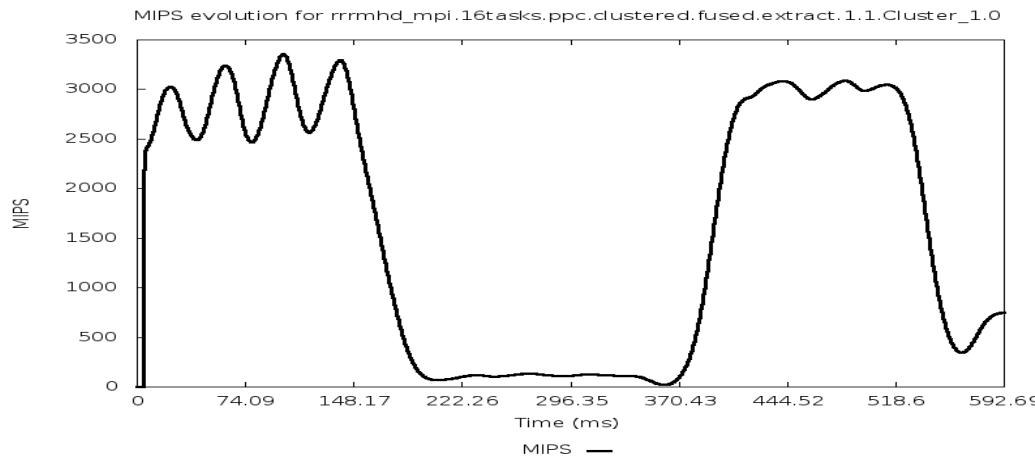
- Call-site sampling information is folded
  - Correlation between hwc and call-sites
  - GVIM/CUBE add-on to show performance within source code
    - Timeless but useful to point performance issues



- Performance of a sequential region = 2000 MIPS

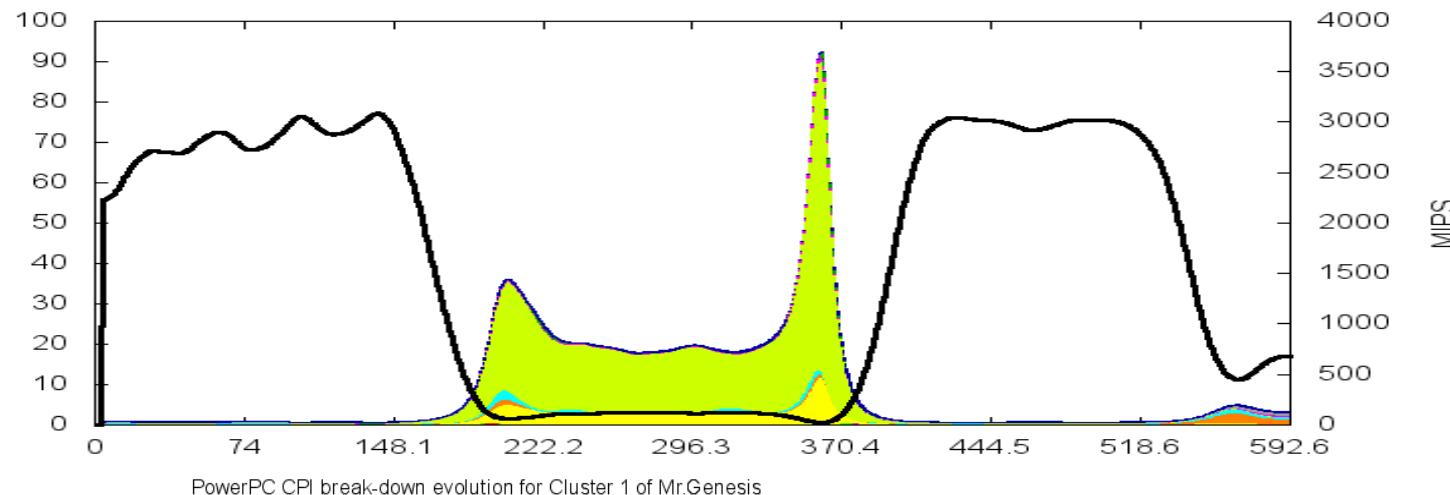
Is it good enough?

Is it easy to improve?

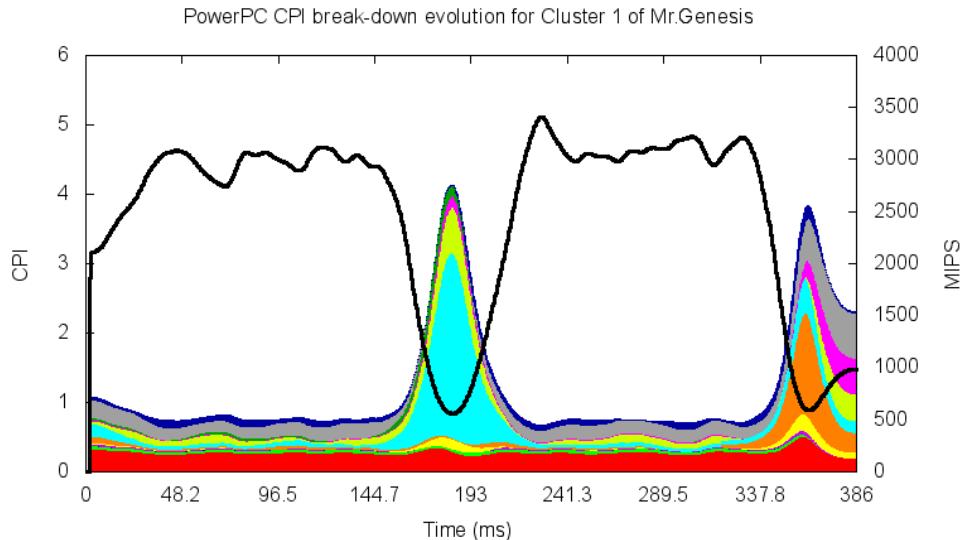


# Instantaneous CPI stack

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



MRGENESIS



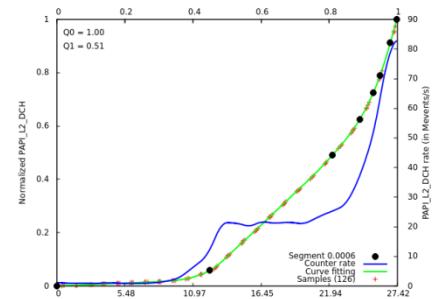
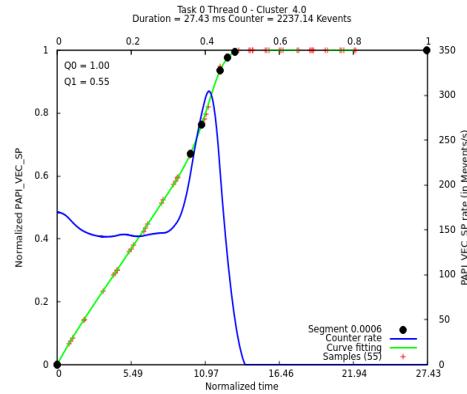
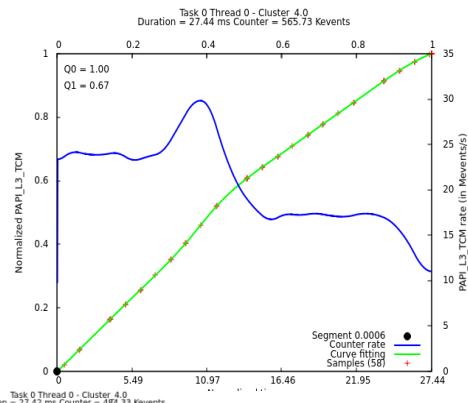
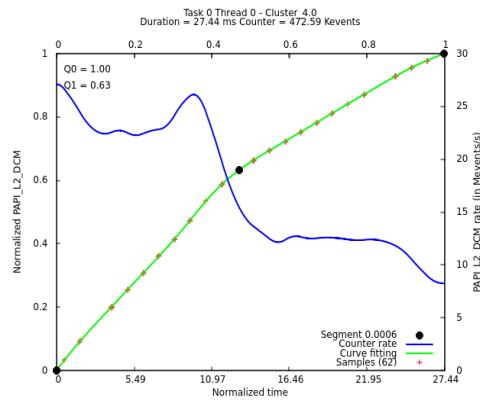
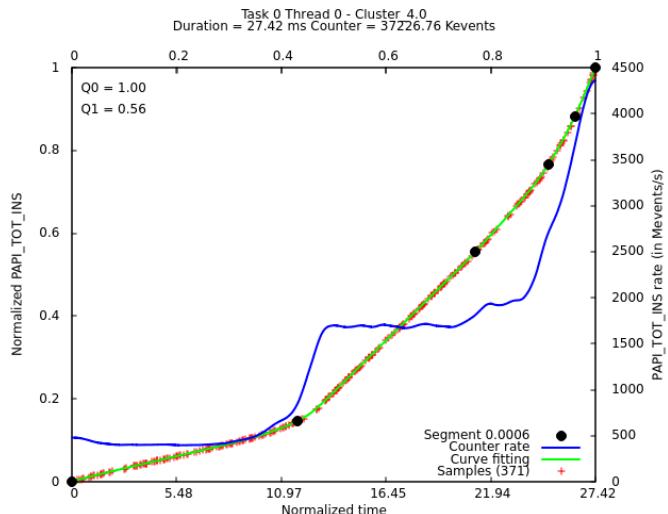
- |                         |                                      |                      |   |
|-------------------------|--------------------------------------|----------------------|---|
| Useful cycles           | <span style="color:red">█</span>     | LSU: Basic latency   | <span style="background-color:yellow;">█</span> |
| I-cache miss            | <span style="color:green">█</span>   | FXU: Div/MSTPR/MSFPR | <span style="background-color:grey;">█</span>   |
| Branch mispredict       | <span style="color:blue">█</span>    | FXU: Basic latency   | <span style="color:magenta;">█</span>           |
| Flush penalties, etc    | <span style="color:purple">█</span>  | FPU: FDIV/FSQRT      | <span style="color:green;">█</span>             |
| LSU: Translation lookup | <span style="color:yellow;">█</span> | FPU: Basic latency   | <span style="background-color:grey;">█</span>   |
| LSU: Other reject       | <span style="color:orange;">█</span> | Other stall cycles   | <span style="color:darkblue;">█</span>          |
| LSU: D-cache miss       | <span style="color:cyan;">█</span>   | MIPS                 | <span style="background-color:black;">█</span>  |

- Trivial fix.(loop interchange)
- Easy to locate?
- Next step?
- Availability of CPI stack models for production processors?
  - Provided by manufacturers?

# Correlating counters

VI-HPS

CG-POP

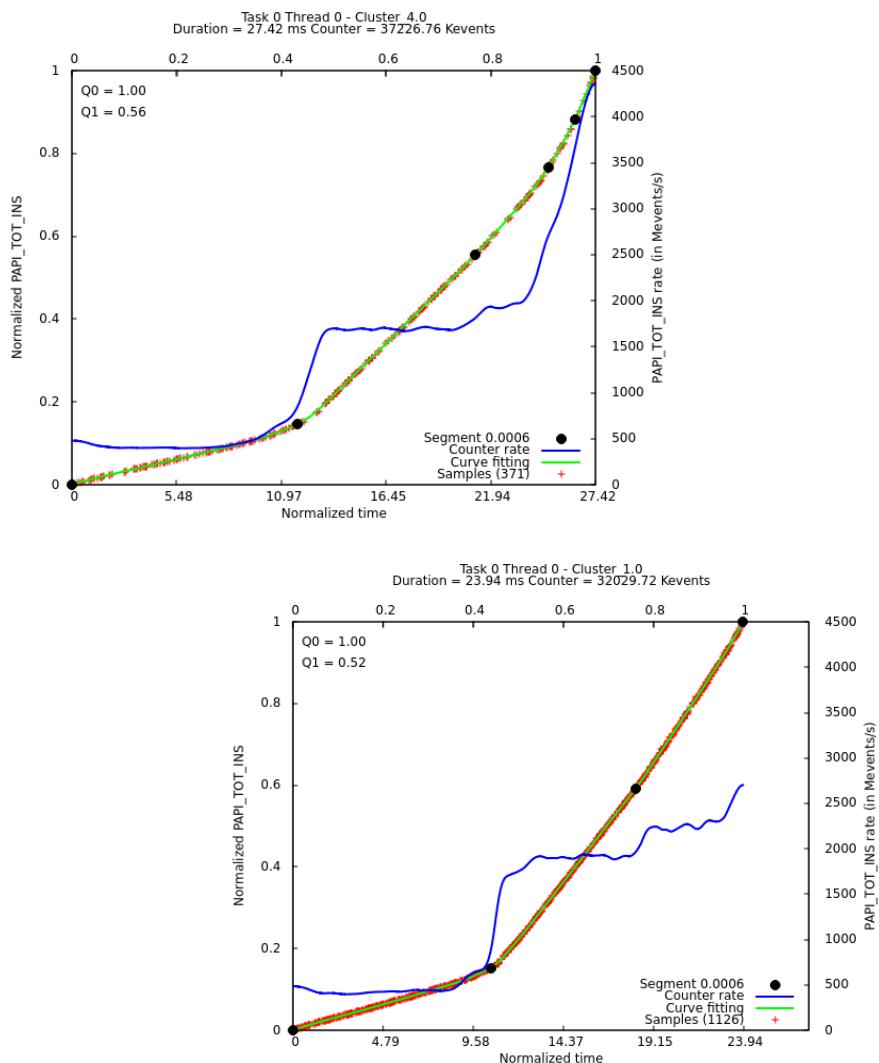


- Within a process
- 3 algorithmic phases
- Impact of multicore sharing

# Correlating counters

VI-HPS

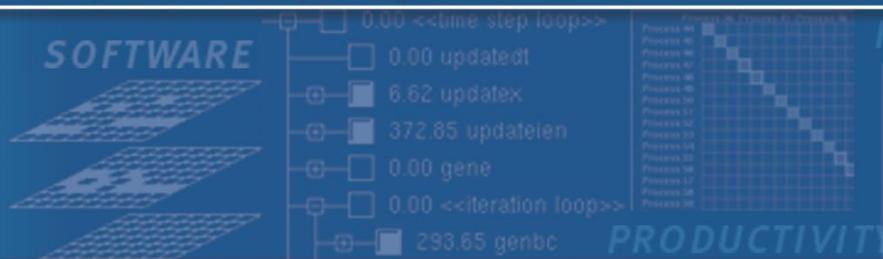
CG-POP



- Between processes
- 3 Algorithmic phases
- Impact of multicore sharing



SOFTWARE



FAST SOLUTIONS

- PAPI\_L1\_DCM
- PAPI\_L1\_ICM
- PAPI\_L2\_DCM
- PAPI\_L2\_ICM
- PAPI\_L1\_TCM
- PAPI\_L2\_TCM

PRODUCTIVITY

# Dimemas

- Key factors influencing performance

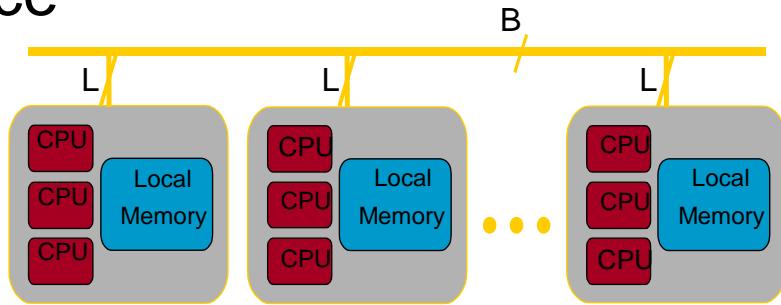
- Abstract architecture
- Basic MPI protocols
- No attempt to model details
- Objectives
  - Simple / general, Fast simulations

- Linear components

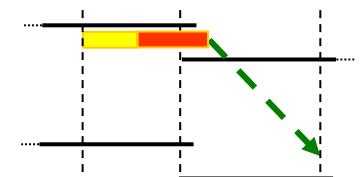
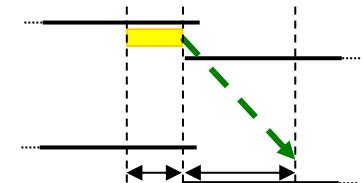
- Point to point communication
- Sequential processor performance (global CPU speed, per block/subroutine)

- Non-linear components

- Synchronization semantics (blocking receives, rendezvous)
- Resources contention (CPU, links half/full duplex, busses)



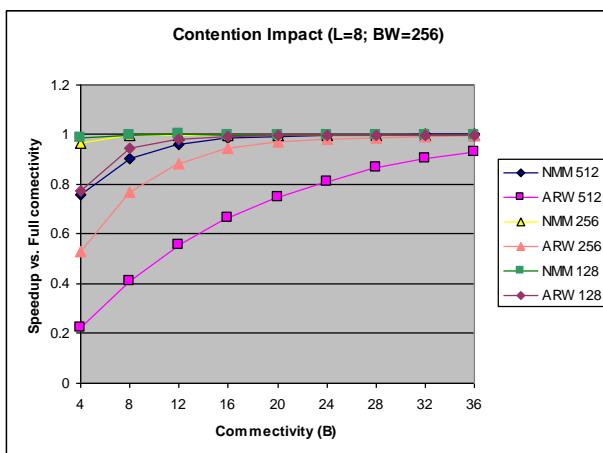
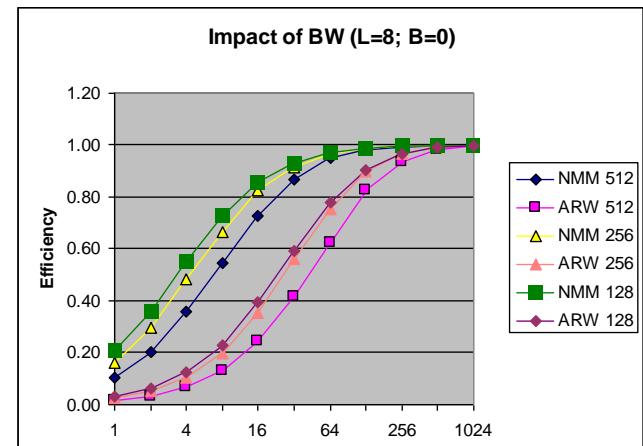
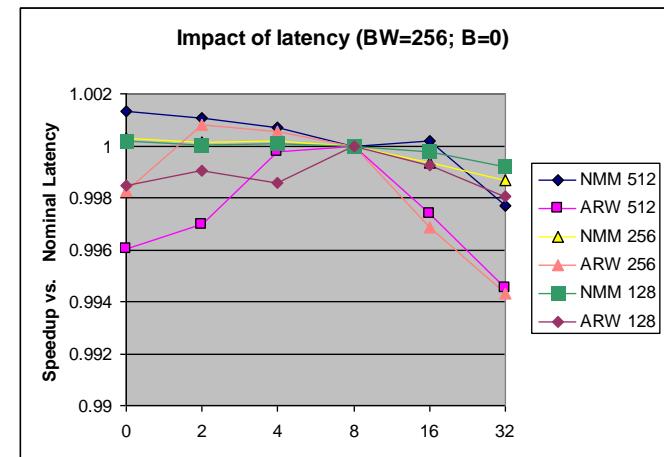
$$T = \frac{\text{MessageSize}}{\text{BW}} + L$$



- Paraver trace: what happens when
  - Actual wall clock time of events
- Dimemas trace: sequence of resource demands
  - Duration of computation bursts
  - Type of communication, partners and bytes
- Can be generated from Paraver trace
  - `prv2dim input.prv output.dim`
- Dimemas generates as output a Paraver file of the simulated run

# Parametric studies – network sensitivity

- WRF, Iberia 4Km, 4 procs/node
  - No sensitive to latency
  - NMM
    - BW – 256MB/s
    - 512 – sensitive to contention
  - ARW
    - BW - 1GB/s
    - Sensitive to contention



# Ideal network machine (instantaneous communication)

VI-HPS

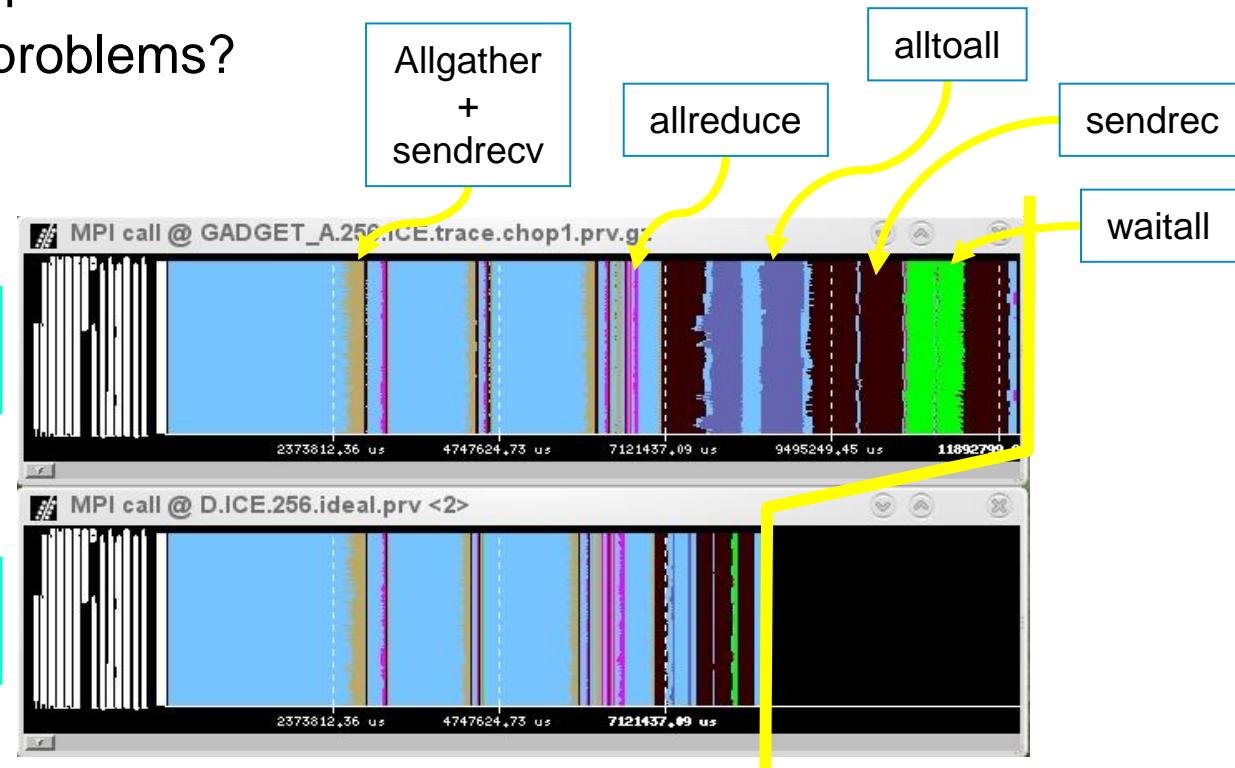
## « The impossible machine: $BW = \infty$ , $L = 0$

- Actually describes/characterizes Intrinsic application behavior
  - Load balance problems?
  - Dependence problems?

GADGET @ Nehalem cluster  
256 processes

Real run

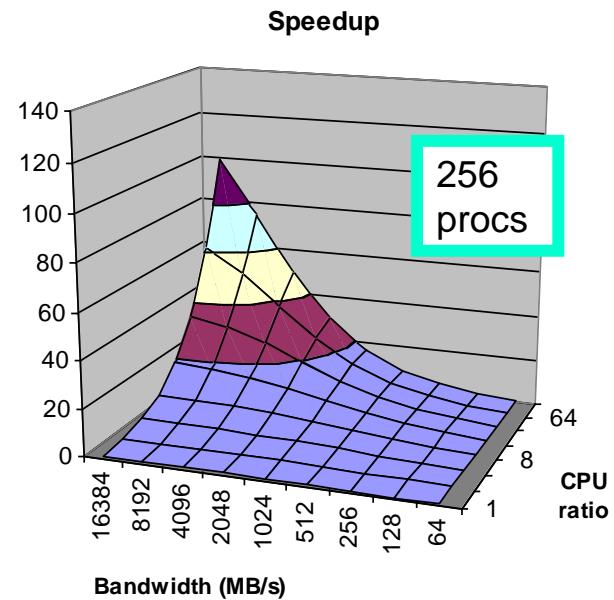
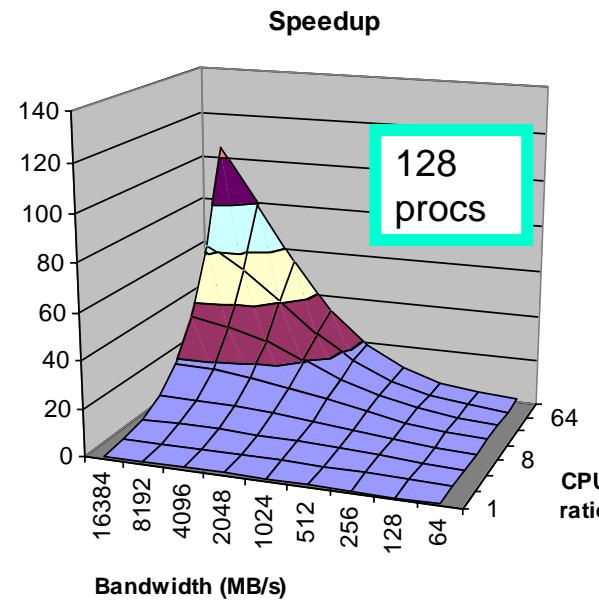
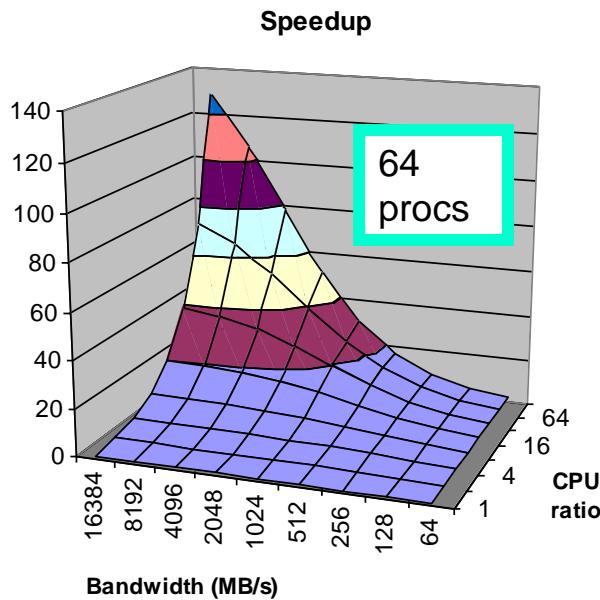
Ideal network



Impact on practical machines?

- **Ideal speeding up ALL the computation bursts by the CPUratio factor**
  - The more processes the less speedup (higher impact of bandwidth limitations) !!

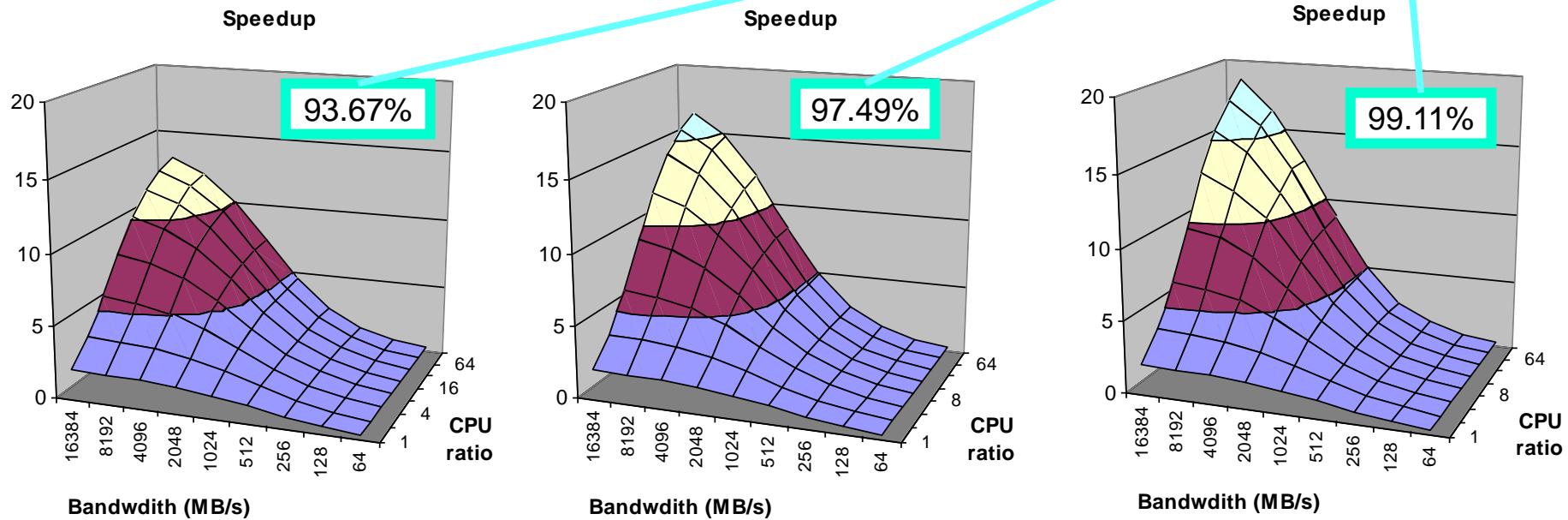
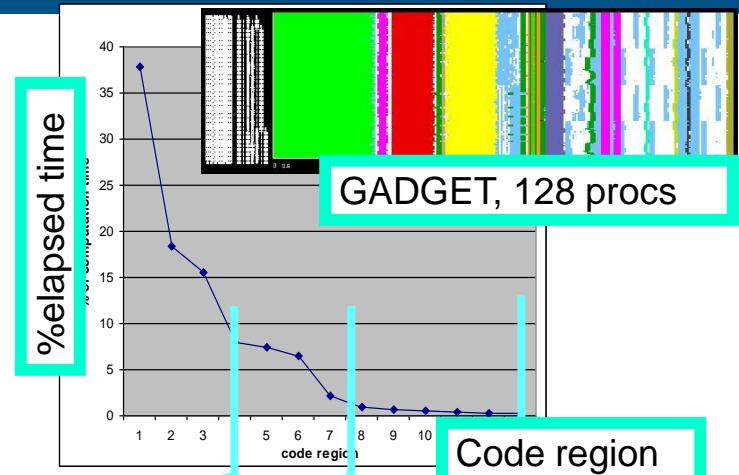
GADGET

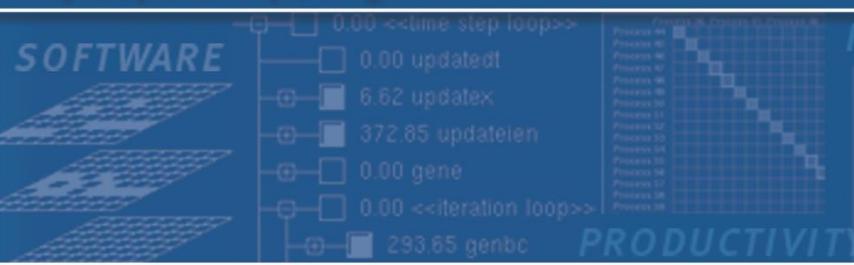


# The potential of hybrid/accelerator parallelization

VI-HPS

- **Hybrid parallelization**
  - Speedup SELECTED regions by the CPUratio factor
- We do need to overcome the **hybrid Amdahl's law**
  - → **asynchrony + Load balancing mechanisms !!!**





FAST SOLUTIONS

PRODUCTIVITY

# Methodology

**Help generate hypotheses**

**Help validate hypotheses**

**Qualitatively**

**Quantitatively**



- 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?

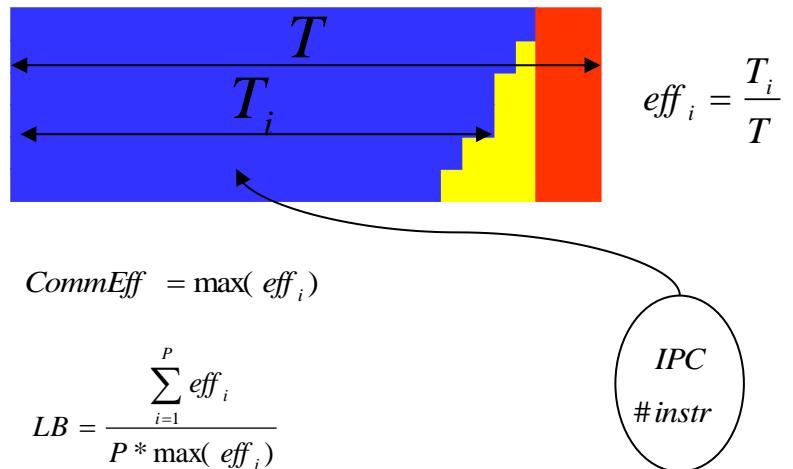
Paraver Tutorial:  
Introduction to Paraver and Dimemas methodology

$$\eta_{\parallel} = LB * CommEff$$

Directly from real execution metrics

2DP - MPI call profile @ trace\_24h\_atmos\_symbols.cho...

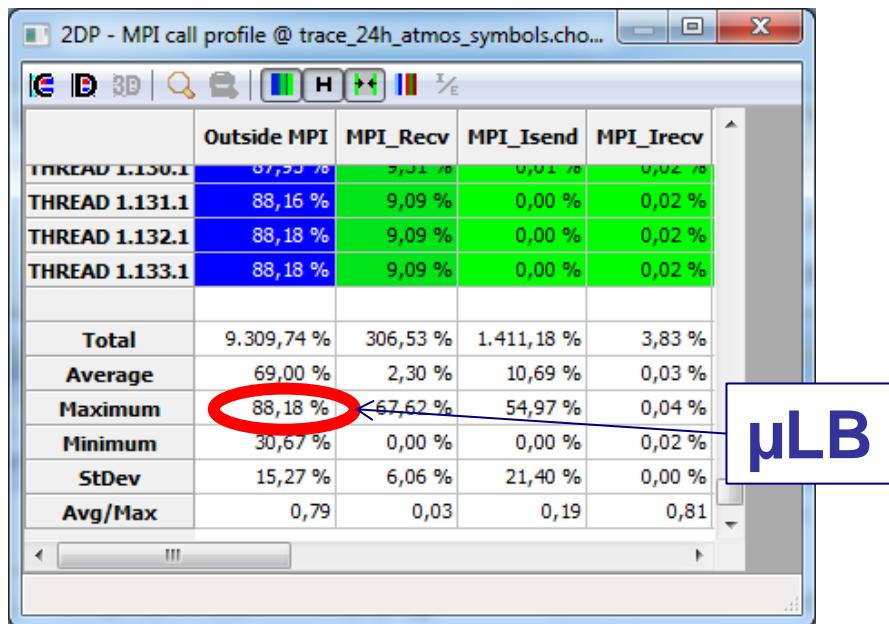
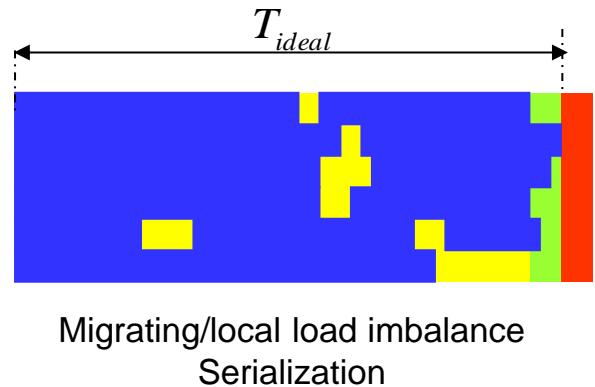
	Outside MPI	MPI_Recv	MPI_Isend	MPI_Irecv
THREAD 1.150.1	67,93 %	2,31 %	0,01 %	0,02 %
THREAD 1.131.1	88,16 %	9,09 %	0,00 %	0,02 %
THREAD 1.132.1	88,18 %	9,09 %	0,00 %	0,02 %
THREAD 1.133.1	88,18 %	9,09 %	0,00 %	0,02 %
<b>Total</b>	<b>9.309,74 %</b>	<b>306,53 %</b>	<b>1.411,18 %</b>	<b>3,83 %</b>
<b>Average</b>	<b>69,00 %</b>	<b>2,30 %</b>	<b>10,69 %</b>	<b>0,03 %</b>
<b>Maximum</b>	<b>88,18 %</b>	<b>9,62 %</b>	<b>54,97 %</b>	<b>0,04 %</b>
<b>Minimum</b>	<b>30,67 %</b>	<b>0,00 %</b>	<b>0,00 %</b>	<b>0,02 %</b>
<b>StDev</b>	<b>15,27 %</b>	<b>6,06 %</b>	<b>21,40 %</b>	<b>0,00 %</b>
<b>Avg/Max</b>	<b>0,79</b>	<b>0,03</b>	<b>0,19</b>	<b>0,81</b>



# Scaling model

- Dimemas simulation with ideal target
  - Latency =0; BW =  $\infty$

$$CommEff = \mu LB * Transfer$$



$$\mu LB = \frac{\max(T_i)}{T_{ideal}}$$

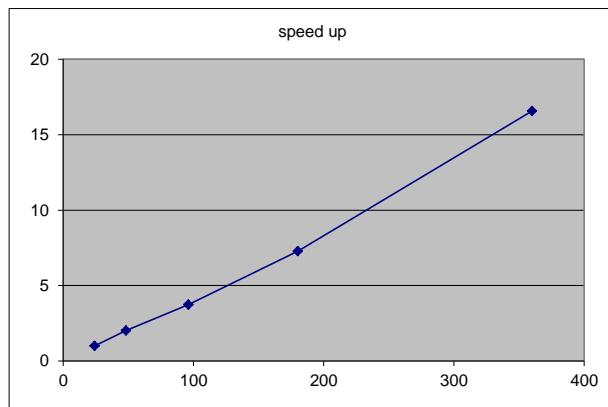
$$Transfer = \frac{T_{ideal}}{T}$$

# Modelling efficiency

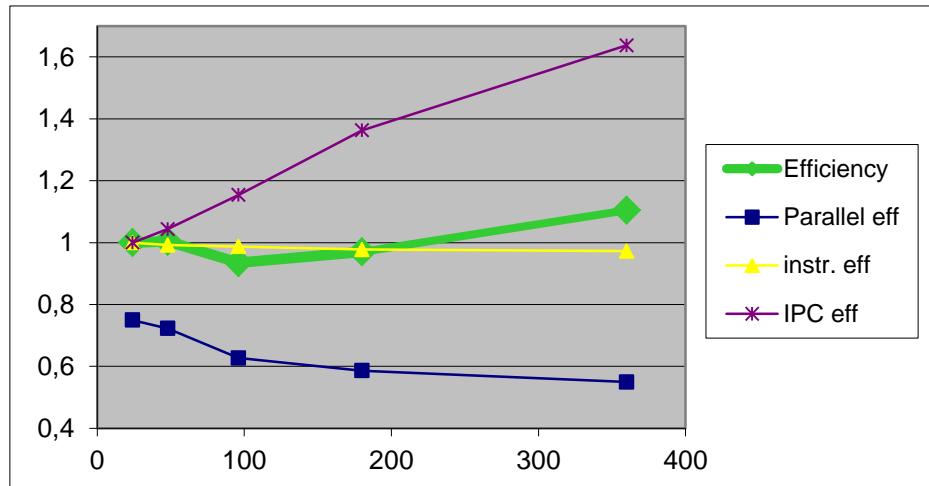
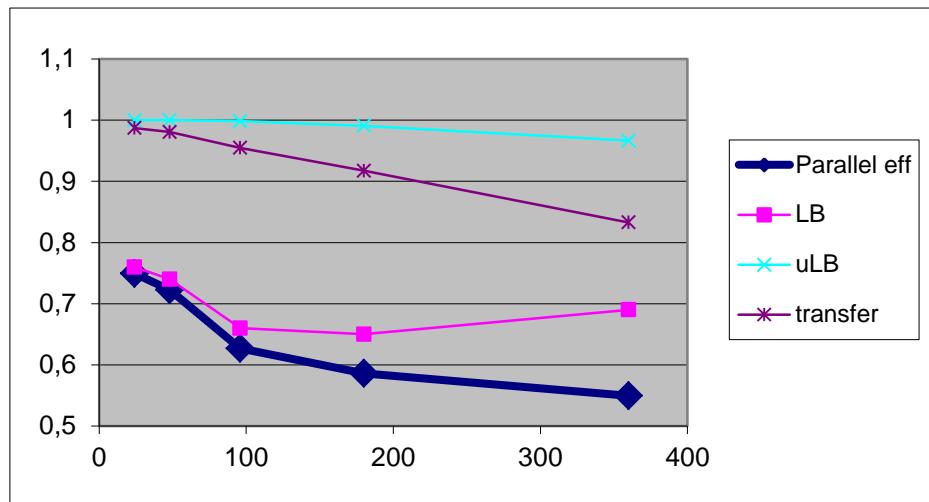
$$\eta_{\parallel} = LB * \mu LB * Transfer$$

CG-POP mpi2s1D - 180x120

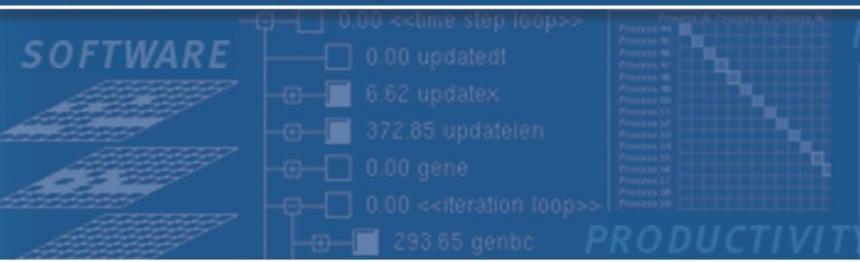
Good scalability !!  
Should we be happy?



$$\eta = \eta_{\parallel} * \eta_{instr} * \eta_{IPC}$$



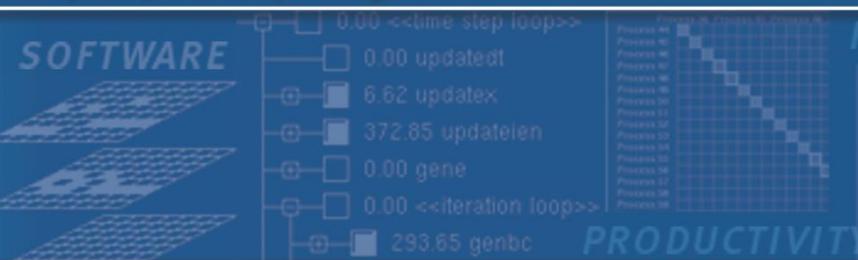
- [www.bsc.es/paraver](http://www.bsc.es/paraver)
  - downloads
    - Sources / Binaries
    - Linux / windows / MAC
  - documentation
    - Training guides
    - Tutorial slides
- Getting started
  - Start wxparaver
  - Help → tutorials and follow instructions
  - Follow training guides
    - Paraver introduction (MPI): Navigation and basic understanding of Paraver operation



FAST SOLUTIONS

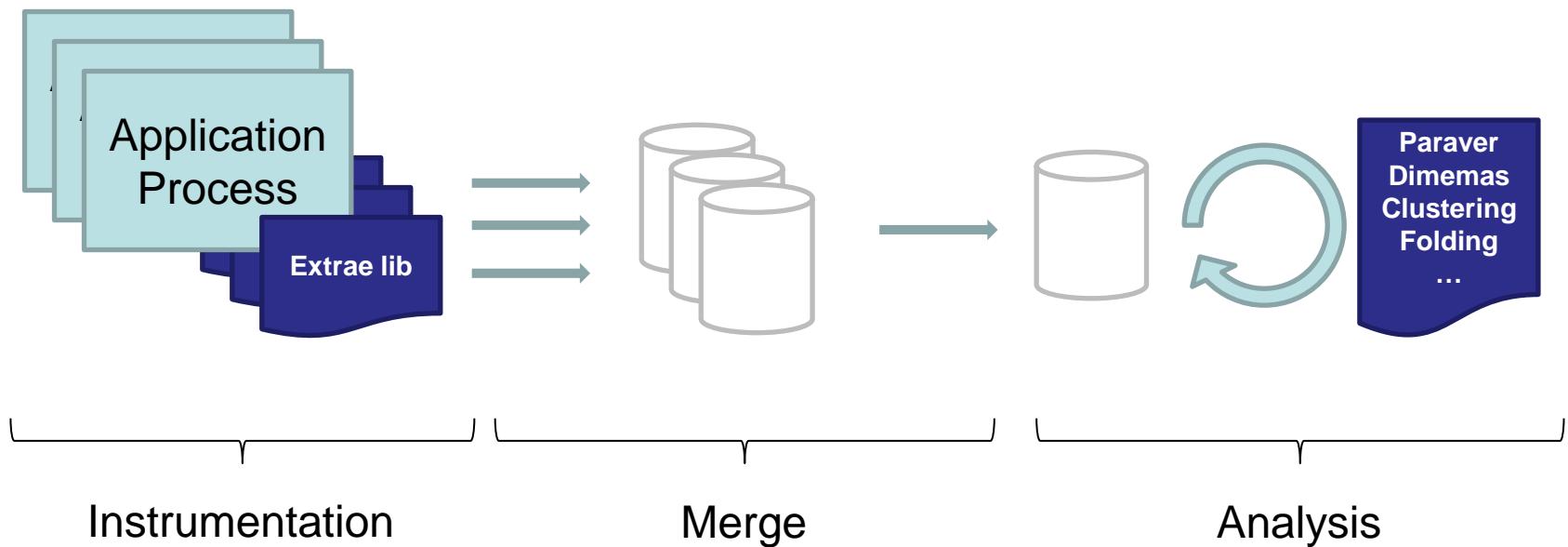
- PAPI\_L1\_DCM
- PAPI\_L1\_ICM
- PAPI\_L2\_DCM
- PAPI\_L2\_ICM
- PAPI\_L1\_TCM
- PAPI\_L2\_TCM

# Short Paraver Demo



# BSC tools info for Juqueen & Juropa

# Trace Generation Workflow



- USB or BSC Tools account `~train036`
  - in Juqueen and Juropa
- Content
  - Binaries for your laptops
    - Paraver (Linux, Windows, and Mac)
    - Dimemas (Linux only)
    - ClusteringSuite (Linux only)
  - Tutorials
  - Extrae Documentation
  - Slides
  - Traces
  - 2 application examples: CG-POP, GTC

- BSC Tools account is `~train036`
  - in Juqueen and Juropa
  
- To use the tools
  - source `~train036/env-bsc-tools`

- Currently, only MPI statically linked programs can be traced
- To facilitate compilation
  - Add to your link stage
    - \$(EXTRAE\_LIBS) for C/C++ apps
    - \$(EXTRAE\_LIBS\_F) for Fortran apps
- We provide 2 full-examples
  - GTC2, CGPOP
    - In ~train036/analysis-apps

- Within `~train036/sample-exrae-job` directory
  - `job.ll` – Execution of the instrumented application
    - Load-Leveler Submission script.
      - Executes the application, and,
      - Submits the merging process
  - `merge.ll` – Merging process to generate a Paraver trace-file
    - Automatically launched through `job.ll`
  - `extrae/extrae.xml` – Extrae configuration

- job.ll

```
#@ job_name = TRACE  
#@ ...  
#@ queue
```

*export BINARY= # Write the binary here*

```
export EXTRAE_CONFIG_FILE=extrae/extrae.xml  
runjob -n N -exe ${BINARY} -env-all
```

**llsubmit merge.ll**

- merge.ll

```
#@ job_name = MERGE  
#@ ...  
#@ queue
```

*export TRACEFILE= # Write here the name of the resulting tracefile  
export BINARY= #Write here the name of the binary*

```
source ~train036/env-bsc-tools  
runjob -n 16 : $EXTRAE_HOME/bin/mpimpi2prv -f TRACE.mpits  
-e ${BINARY} -o ${TRACEFILE}
```

- extrae.xml

```
<counters enabled="yes">
  <cpu enabled="yes" starting-set-distribution="cyclic">

    <set enabled="yes" domain="all" changeat-globalops="125">
      PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_L1_DCM,PAPI_TLB_DM,PAPI_LD_INS,PAPI_BR_INS,PAPI_BR_MSP,PAPI_FP_OPS
    </set>
    <set enabled="yes" domain="all" changeat-globalops="125">
      PAPI_TOT_INS,PAPI_TOT_CYC,PEVT_AXU_FP_DS_ACTIVE_CYC,PEVT_IU_IL1_MISS_CYC,PEVT_IU_RELOAD_COLLISIONS_CYC,PEVT_IU_IU0_REDİR_CYC,PEVT_IU_IERAT_MISS_CYC
    </set>
    <set enabled="yes" domain="all" changeat-globalops="125">
      PAPI_TOT_INS,PAPI_TOT_CYC,PEVT_IU_ICACHE_FETCH_CYC,PEVT_IU_INSTR_FETCHED_CYC,PEVT_IU_FX_U_ISSUE_PRI_LOSS_CYC,PEVT_IU_RAW_DEP_HIT_CYC
    </set>
    ...
    <set enabled="yes" domain="all" changeat-globalops="125">
      PAPI_TOT_INS,PAPI_TOT_CYC,PEVT_L1P_BAS LU_STALL_SRT_CYC,PEVT_L1P_BAS LU_STALL_MMIO_DCR_CYC,PEVT_L1P_BAS LU_STALL_STRM_DET_CYC,PEVT_L1P_BAS LU_STALL_LIST_RD_CYC,PEVT_L1P_BAS LU_STALL_LIST_WRT_CYC
    </set>
  </cpu>
  <network enabled="no" />
  <resource-usage enabled="no" />
  <memory-usage enabled="no" />
</counters>
```

Hardware counters

- **extrae.xml**

```
<storage enabled="no">
    <trace-prefix enabled="yes">TRACE</trace-prefix>
    <size enabled="no">5</size>
    <temporal-directory enabled="yes">/scratch</temporal-directory>
    <final-directory enabled="yes">/gpfs/scratch/bsc41/bsc41273</final-directory>
    <gather-mpits enabled="no" />
</storage>
```

Storage options

```
<buffer enabled="yes">
    <size enabled="yes">500000</size>
    <circular enabled="no" />
</buffer>
```

Buffering options

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

MPI section

```
<sampling enabled="no" type="default" period="50m" variability="20m" />
```

Sampling

```
<callers enabled="yes">
    <mpi enabled="yes">1-3</mpi>
    <sampling enabled="yes">1-5</sampling>
</callers>
```

Callstack information

- Within `~train036/sample-exrae-job` directory
  - `job.moab` – Execution of the instrumented application
    - Submission script.
      - Executes the application,
  - `merge.moab` – Merging process to generate a Paraver trace-file
    - It is NOT automatically launched through `job.moab`
  - `extrae/extrae.xml` – Extrae configuration
  - `extrae/trace.sh` – intercept execution with tracing library

```
#!/bin/bash
source ~/train036/env-bsc-tools
export EXTRAE_CONFIG_FILE=extrae/extrae.xml
export LD_PRELOAD=$EXTRAE_HOME/lib/libmpitracef.so
$@
```

- Instrumentation through shared library preloading
  - Specific libraries for each combination of runtimes
    - MPI
    - OpenMP
    - OpenMPI+MPI
    - ...

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

<sup>1</sup> for Fortran codes

- job.moab

```
#!/bin/bash -x  
#MSUB -e trace.err  
#MSUB -o trace.out
```

..

```
export BINARY= # Write here the binary you want to trace  
export NP= # Write here the number of processes
```

```
mpiexec -np ${NP} extrae/trace.sh ${BINARY}
```

- merge.moab

```
#!/bin/bash -x  
#MSUB -e merge.err  
#MSUB -o merge.out
```

..

*export TRACEFILE= # Write here the tracefile name*

```
source ~train036/env-bsc-tools  
mpiexec -np ${NP} ${EXTRAE_HOME}/bin/mpimpi2prv -f  
TRACE.mpits -o ${TRACEFILE}
```

1. Follow guidelines for paraver basic navigation (Tut #1)  
Tracefile provided
2. (if Juqueen) Build your binary with extrae library
3. Adapt scripts to instrument your application
4. Obtain a tracefile
5. Follow first steps of methodology guidelines (Tut #5)
  - Parallel efficiency
  - Distribution of computations
  - Instruction balance
  - ...
6. Depending on the diagnosis use clustering, folding, dimemas...



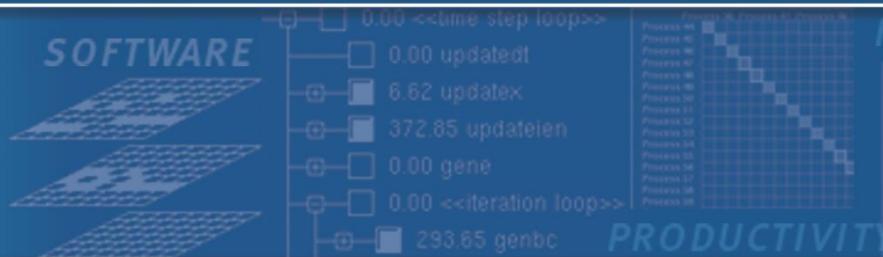
**Barcelona  
Supercomputing  
Center**

*Centro Nacional de Supercomputación*

**THANKS**



SOFTWARE



# Detailed material

# Semantic Module

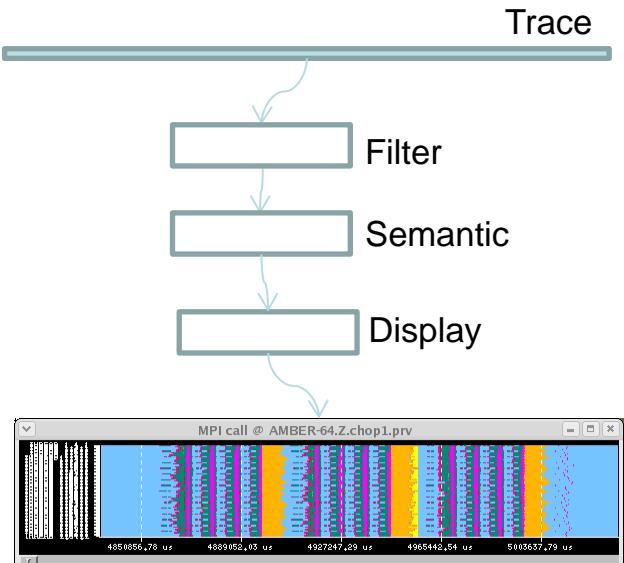
# Basic functions of time

- The **filter module** presents a subset of the trace to the semantic module. Each thread  $th$  is described by
  - A sequence of events  $Ev_i, i \in N$ , states  $St_i, i \in N$  and communications  $C_i, i \in N$
  - For each event let  $T(Ev_i)$  be its time and  $V(Ev_i)$  its value
  - For each state let  $T_s(St_i)$  be its start time  $T_e(St_i)$  its stop time and  $V(St_i)$  its value
  - For each Communication let  $T_s(C_i)$  be its send time,  $T_r(C_i)$  its receive time,  $Sz(C_i)$  its size.
  - $Partner(C_i)$  and  $Dir(C_i) \in \{send, recv\}$  identify the partner process and direction of the transfer

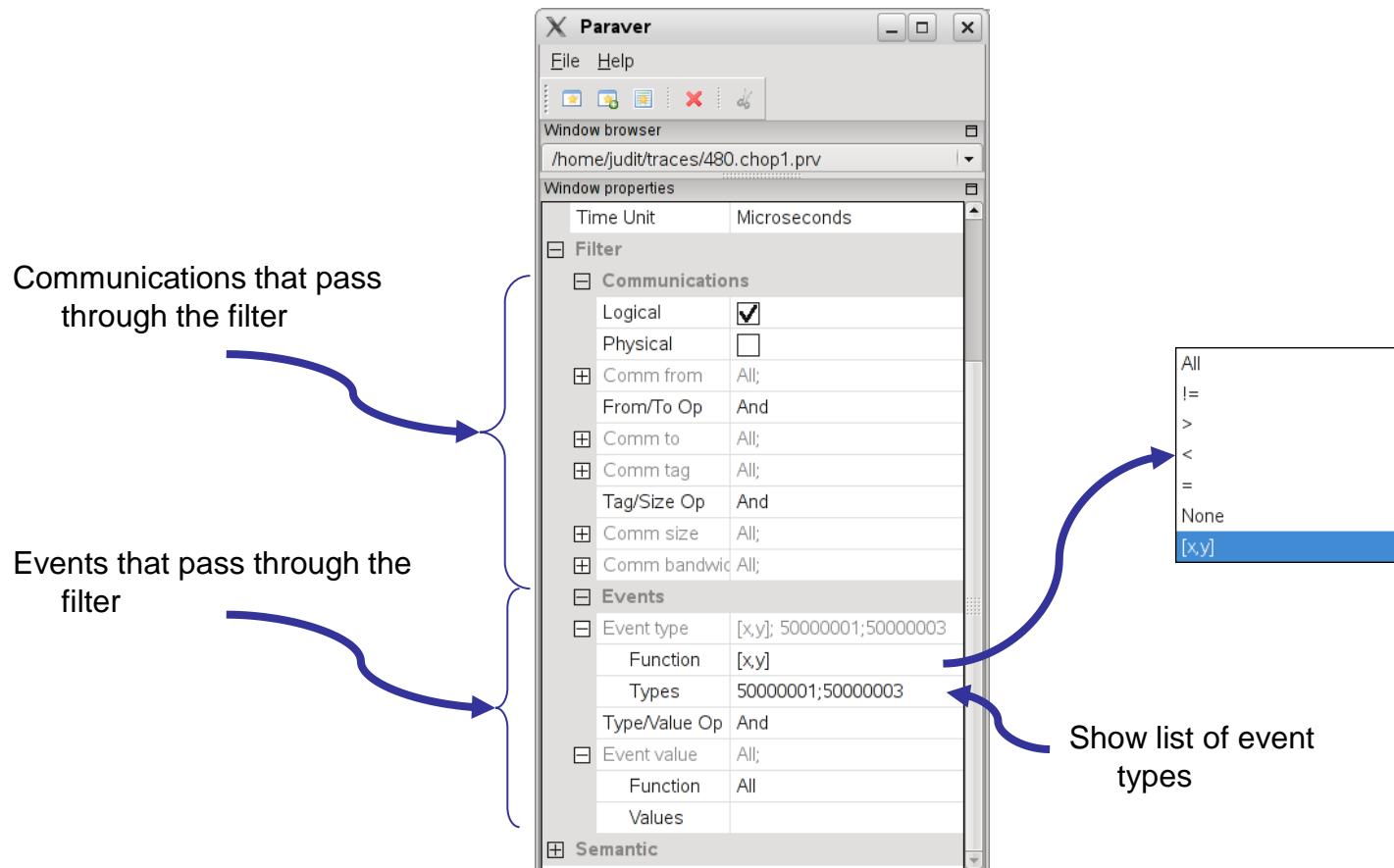
- Semantic module** builds

$$s(t) = S(i), t \in [t_i, t_{i+1}), i \in \mathbb{N}$$

Function of time      Series of values

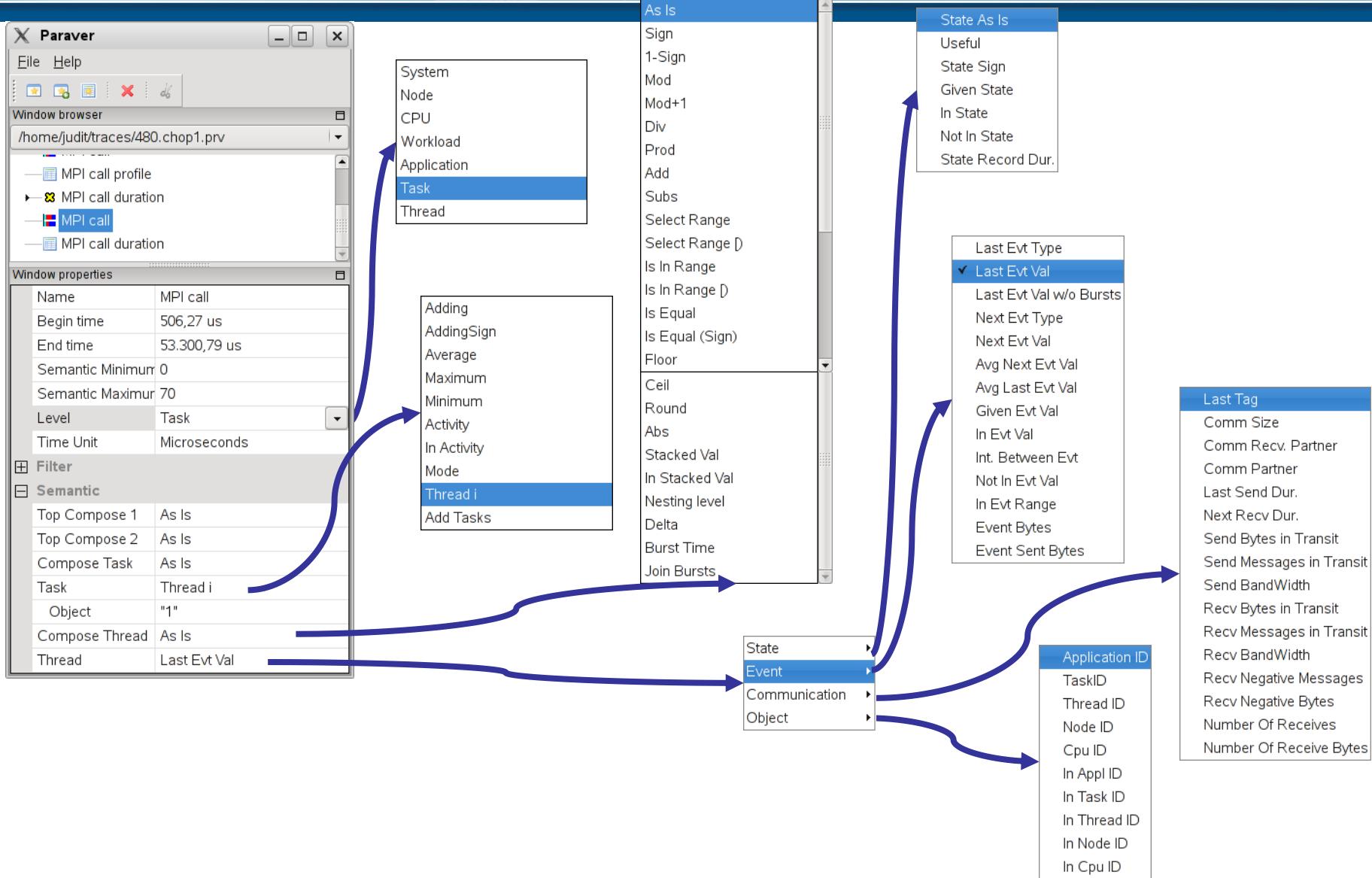


# Filter module



# Semantic module: Control

VI-HPS



- From Events to functions of time

- Last event value

$$S(i) = V(Ev_i)$$

- Next event value

$$S(i) = V(Ev_{i+1})$$

- Average Next Event Value

$$S(i) = \frac{V(Ev_{i+1})}{T(Ev_{i+1}) - T(Ev_i)}$$

- Interval btw. Events

$$S(i) = T(Ev_{i+1}) - T(Ev_i)$$

- From communication records to functions of time

- Send Bytes

$$s(t) = \sum_j Sz(C_j), j \mid (T_s(C_j) < t) \wedge (T_r(Cj) > t) \wedge (Dir(Cj) == send)$$

- Send Bandwidth

$$s(t) = \sum_j \frac{Sz(C_j)}{T_r(C_j) - T_s(C_j)}, j \mid (T_s(C_j) < t) \wedge (T_r(Cj) > t) \wedge (Dir(Cj) == send)$$

- Msgs in transit

$$s(t) = \sum_j sign(j), j \mid (T_s(C_j) < t) \wedge (T_r(Cj) > t) \wedge (Dir(Cj) == send)$$

- Recv. Bandwidth

$$s(t) = \sum_j \frac{Sz(C_j)}{T_r(C_j) - T_s(C_j)}, j \mid (T_s(C_j) < t) \wedge (T_r(Cj) > t) \wedge (Dir(C_j) == recv)$$

- Rec. Negative Msgs

$$s(t) = \sum_j sign(j), j \mid (T_r(C_j) < t) \wedge (T_s(Cj) > t) \wedge (Dir(Cj) == recv)$$

- Comm. Partner

$$s(t) = Partner(C_j), j \mid (T_s(C_j) < t) \wedge (T_r(Cj) > t)$$

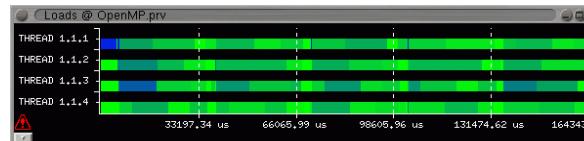
- Bytes btw. Events

$$S(i) = \sum_j Sz(C_j), j \mid T_s(C_j) \in [T(Ev_i), T(Ev_{i+1})) \vee T_r(C_j) \in [T(Ev_i), T(Ev_{i+1}))$$

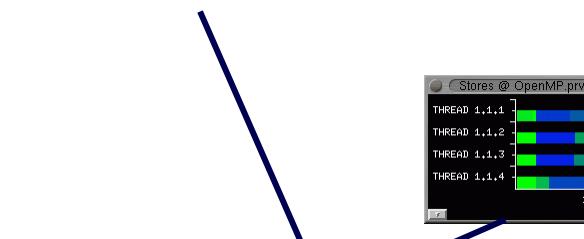
- $S'(t) = f(S(t))$        $S' = f \circ S$ 
  - Sign       $S'(t) = sign(S(t))$
  - 1-sign       $S'(t) = 1 - sign(S(t))$
  - Select range       $S'(t) = S(t) \in [a, b] ? S(t) : 0$
  - Sign  $\circ$  Is equal       $S'(t) = sign(S(t) = a ? S(t) : 0)$
  - Delta       $S'(t) = S_{i+1} - S_i$
  - Stacked value

# Semantic module

- Derived windows
  - Point wise operation
    - $S = \alpha * S^a <\text{op}> \beta * S^b$
    - $<\text{op}> : +, -, *, /, \dots$

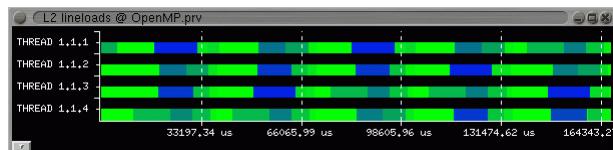


Loads

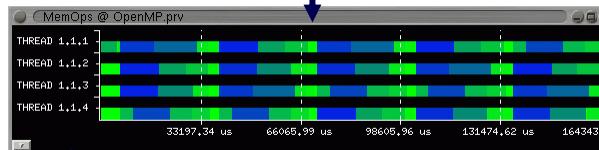


Stores

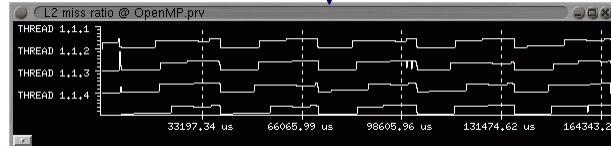
## L2 Line Loads



x100



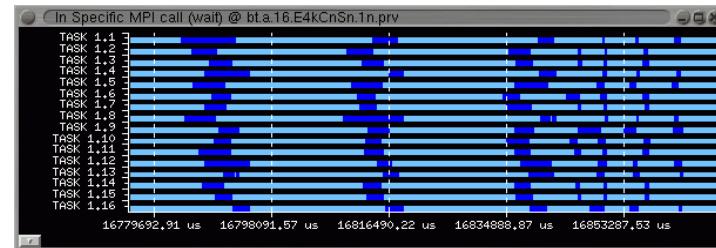
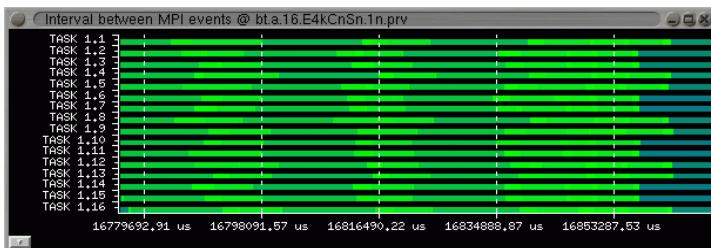
Mem Ops



L2 miss ratio

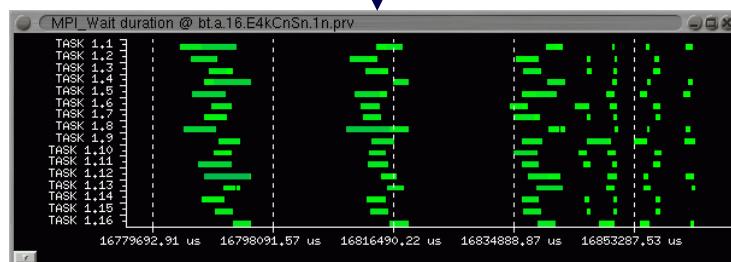
# Semantic module

- Derived windows
  - Point wise operation
    - $S = \alpha * S^a <\text{op}> \beta * S^b$
    - $<\text{op}> : +, -, *, /, \dots$

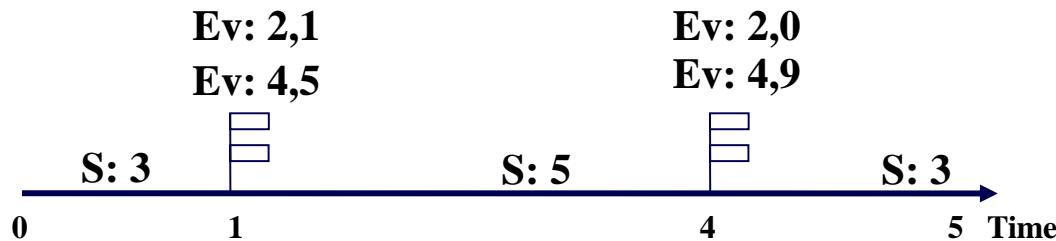


Interval between MPI events

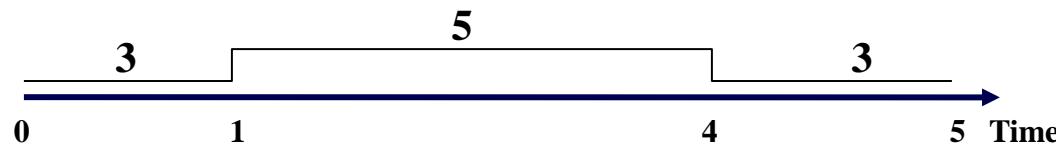
In MPI call



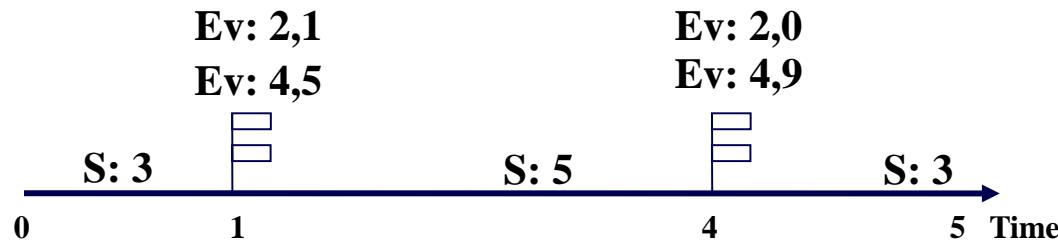
MPI call duration



- Thread function: State as is



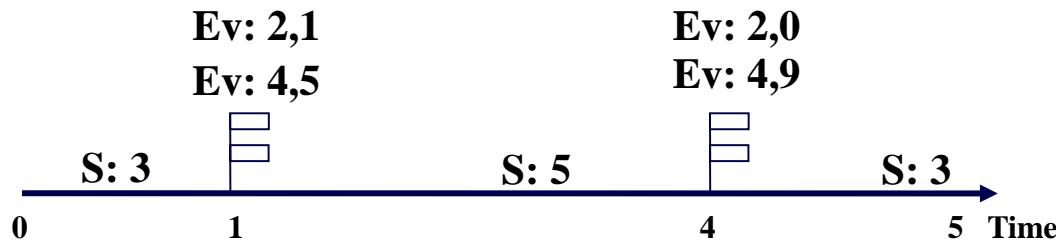
- Useful for
  - Global thread activity: computing, idle, fork/join, waiting,.....



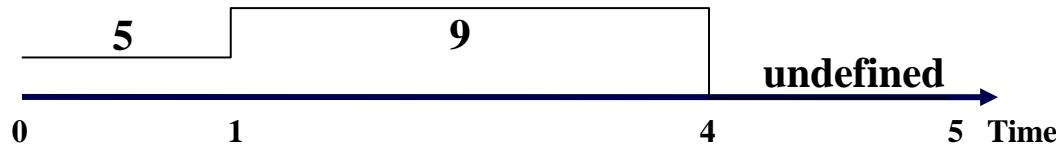
- Filter: type == 2
  - Thread function: Last event value



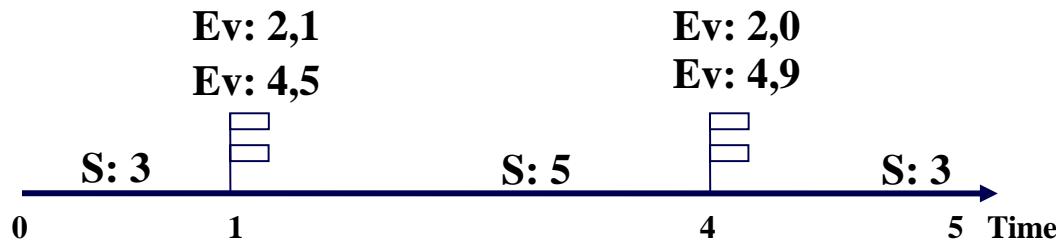
- Useful for
  - In parallel region
  - Mutual exclusion
  - Variable values: iteration,....



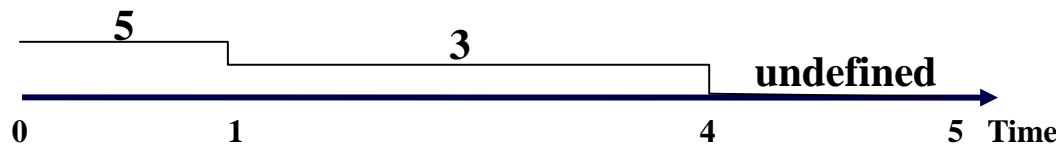
- Filter: type == 4
  - Thread function: Next event value



- Useful for
  - Hwc events (TLB, L1 misses,...) within interval



- Filter: type == 4
  - Thread function: Average next event value

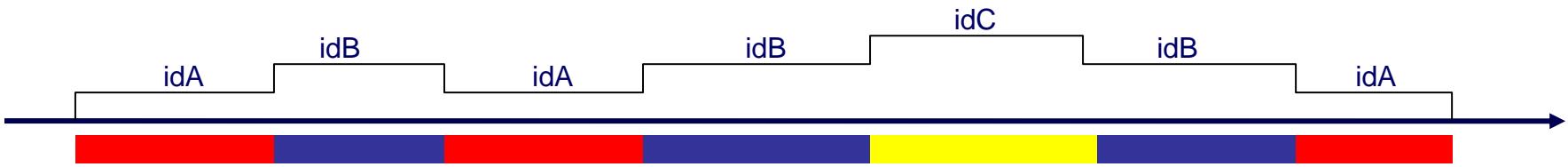


- Useful for
  - Hwc events (TLB, L1 misses,...) per time unit within interval

# Semantic module: Examples



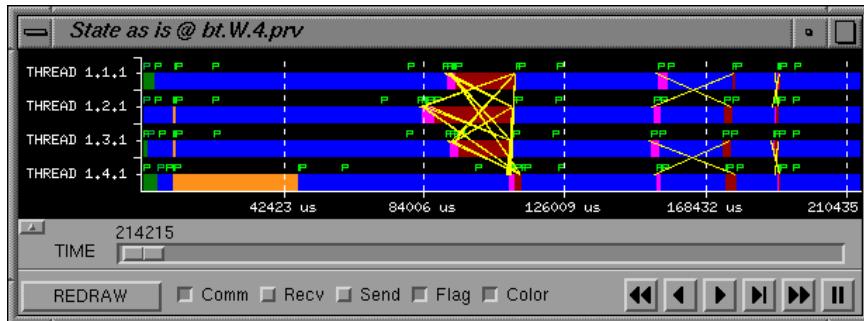
- Filter: type == USR\_FCT  
Thread function: Last event value  
Compose: Stacked value



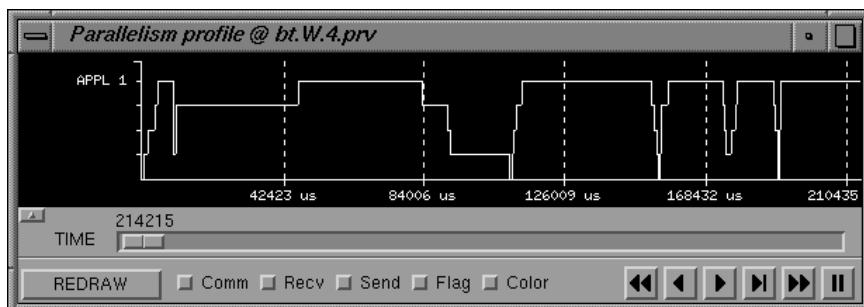
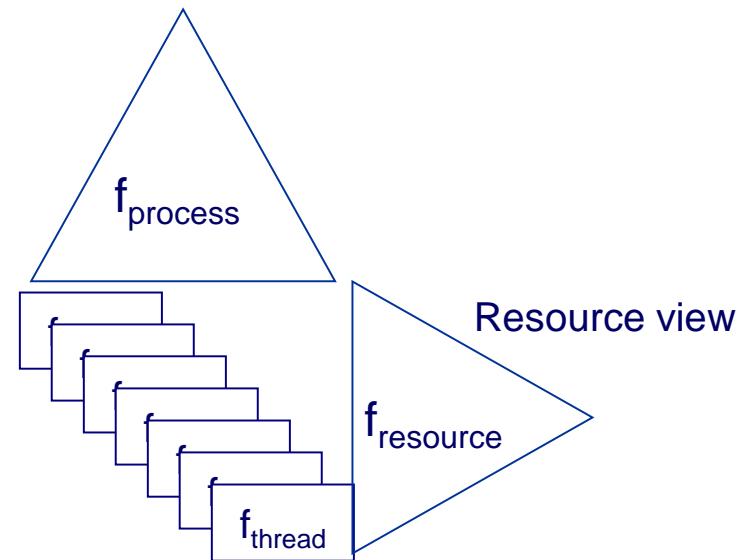
- Useful for
  - Routine

- Process model
  - Thread, task, application, workload
- Resource model
  - CPU, node, system

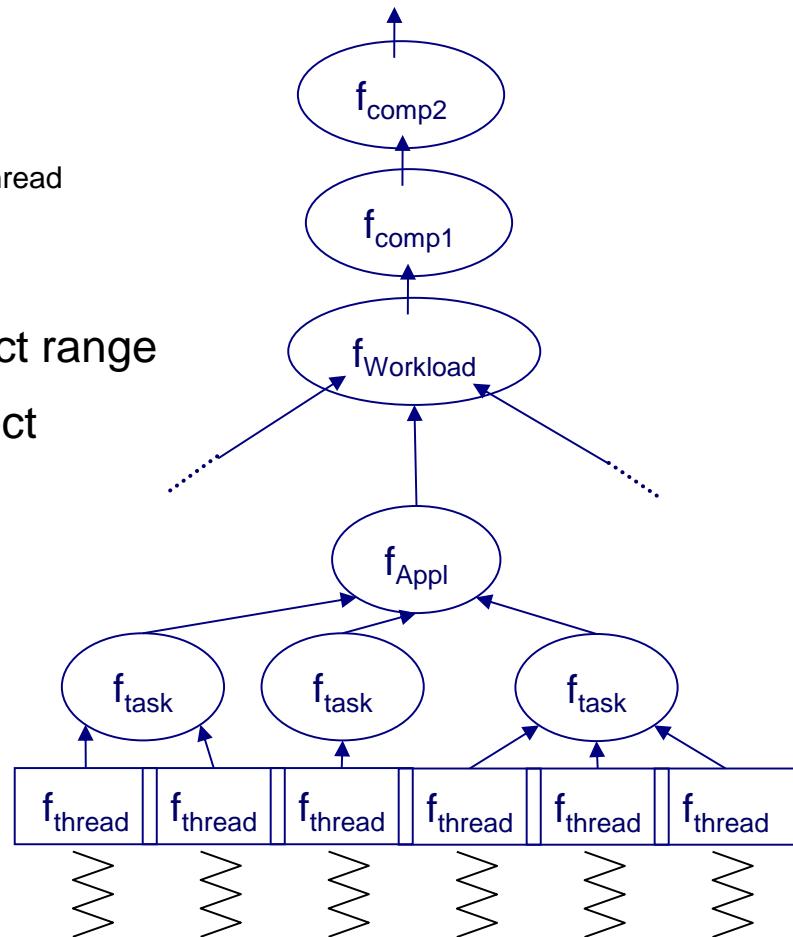
Process view



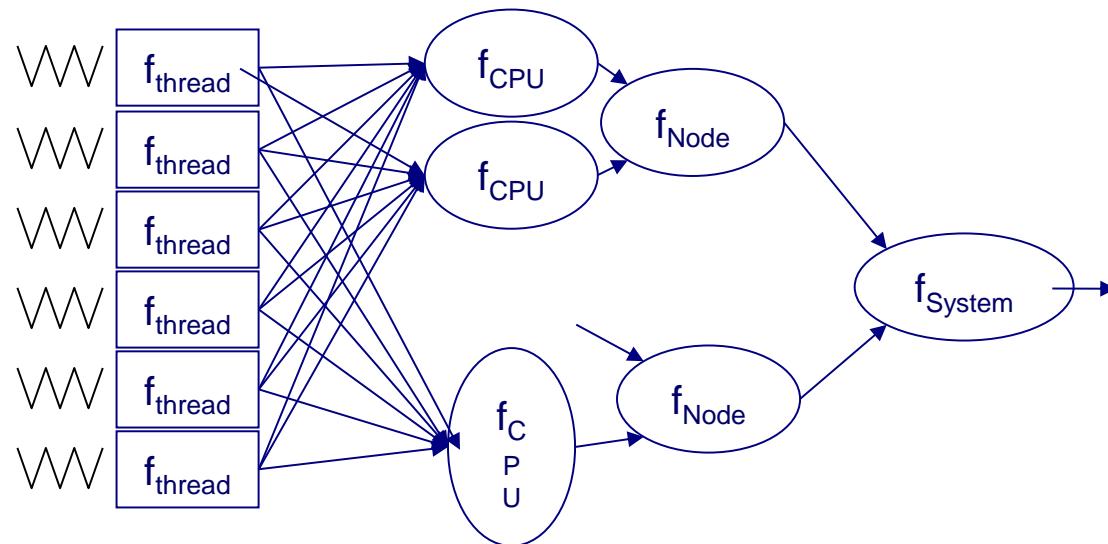
Resource view



- Semantic value:  $S(t)$
- $S = f_{comp2} \circ f_{comp1} \circ f_{Workload} \circ f_{Application} \circ f_{task} \circ S_{thread}$
- Semantic functions
  - $f_{comp2}, f_{comp1}$ : sign, mod, div, in range, select range
  - $f_{Application}, f_{Workload}$ : add, average, max, select
  - $f_{task}$ : add, average, max, select
  - $S_{thread}$ : in state, useful, given state,
  - last event value,
  - next event value,
  - average next event value
  - interval between events, ...



- $Sf_{resource} = f_{comp2} \circ f_{comp1} \circ f_{System} \circ f_{Node} \circ f_{CPU} \circ S_{thread}$
- Semantic functions
  - $f_{System}$  : add, average, max, select
  - $f_{Node}$  : add, average, max, select
  - $f_{CPU}$  : active thread, select
  - $S_{thread}$ : in state, useful, given state, next event value, thread\_id



# Analysis Module

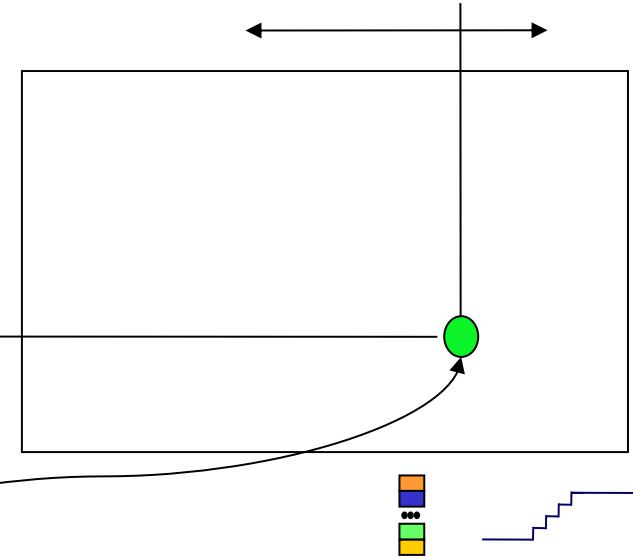
# How to read profiles

One columns per specific value of categorical **Control window**

	End	MPI_Isend	MPI_Irecv	MPI_Wait	MPI_Allreduce	MPI_Comm
THREAD 1.1.1	86,98 %	0,06 %	0,08 %	11,12 %	1,75 %	
THREAD 1.2.1	88,29 %	0,10 %	0,10 %	9,95 %	1,56 %	
THREAD 1.3.1	88,33 %	0,13 %	0,10 %	9,92 %	1,51 %	
THREAD 1.4.1	89,75 %	0,10 %	0,09 %	8,62 %	1,44 %	
THREAD 1.5.1	89,47 %	0,11 %	0,10 %	8,85 %	1,46 %	
THREAD 1.6.1	88,76 %	0,12 %	0,09 %	9,54 %	1,48 %	
THREAD 1.7.1	91,77 %	0,13 %	0,10 %	6,51 %	1,49 %	
THREAD 1.8.1	90,23 %	0,06 %	0,08 %	8,13 %	1,50 %	
THREAD 1.9.1	91,88 %	0,13 %	0,09 %	6,73 %	1,17 %	
THREAD 1.10.1	93,24 %	0,18 %	0,11 %	5,41 %	1,05 %	
THREAD 1.11.1	93,25 %	0,18 %	0,11 %	5,45 %	1,00 %	
THREAD 1.12.1	94,63 %	0,17 %	0,11 %	4,16 %	0,93 %	
THREAD 1.13.1	93,40 %	0,17 %	0,11 %	5,35 %	0,96 %	
THREAD 1.14.1	94,99 %	0,20 %	0,11 %	3,77 %	0,93 %	
THREAD 1.15.1	96,80 %	0,22 %	0,11 %	1,92 %	0,95 %	
THREAD 1.16.1	95,73 %	0,12 %	0,09 %	2,99 %	1,06 %	

MPI call, user function,...

Thread



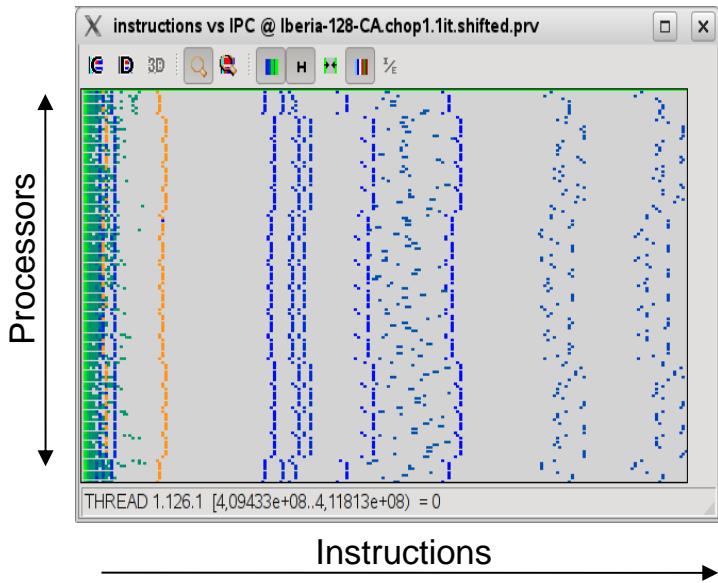
Value/color is a statistic computed for the specific thread  
when control window had the value corresponding to the column

## Relevant statistics:

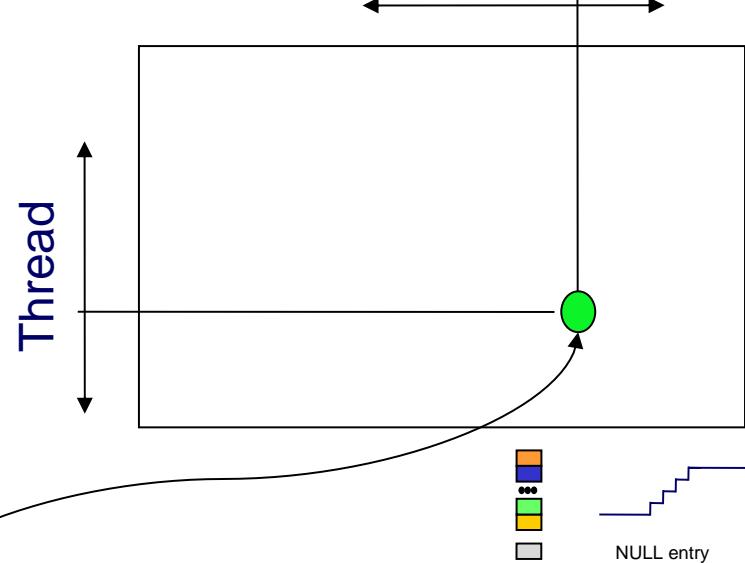
Time, %time, #bursts, Avg. burst time  
Average of **Data window**

# How to read histograms

Columns correspond to bins of values of a numeric **Control window**



duration, instructions, BW, IPC, ...



Value/color is a statistic computed for the specific thread  
when control window had the value corresponding to the column

## Relevant statistics:

Time, %time, #bursts, Avg. burst time  
Average of **Data window**

- Single flexible quantitative analysis mechanism

- Let

- $cw_1$  and  $cw_2$  two views we will call control views
- $dw$  a view we will call data window

For each window w

$$S_{th}^w(t) = S_{th}^w(i), t \in [t_i^w, t_{i+1}^w)$$

- For each control window we define a set of bins

$$bin_j^{cw} = [range_j^{cw}, range_{j+1}^{cw}) \quad range_{j+1}^{cw} = range_j^{cw} + delta^{cw}$$

- And the discriminator functions

$$\delta_j^{cw}(t) = ((S^{cw}(t) \in bin_j^{cw}) ? 1 : 0)$$

Identify regions with cw's within the (j,k) bin

$$\delta_{j,k}(t) = \delta_j^{cw_1}(t) * \delta_k^{cw_2}(t)$$

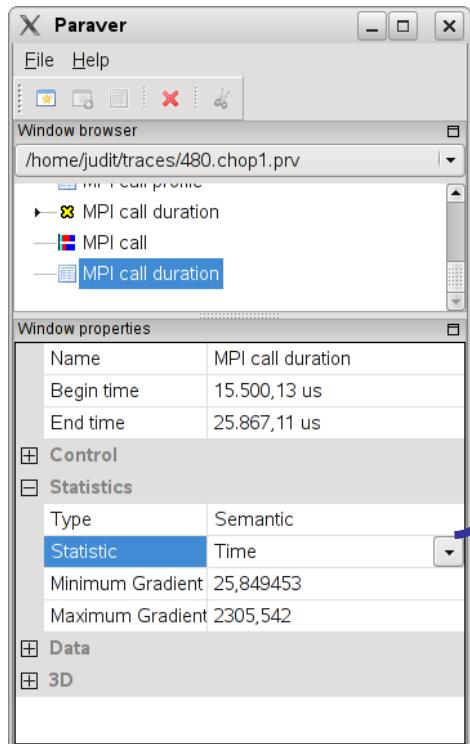
- The 3D analysis module computes a cube (or plane in the case of 2D) of statistics

$$M(thread, j, k) = statistic(S_{th}^{dw}(t) * \delta_{th,j,k}(t))$$

- Where the statistic can represent the average value, the number of intervals,....

# 2D analysis module

VI-HPS



**Time** ( $th, j, k$ ) =  $\int_{t_{start}}^{t_{end}} \delta_{th,j,k}(t) dt$

**% Time** ( $th, j, k$ ) =  $\frac{\int_{t_{start}}^{t_{end}} \delta_{th,j,k}(t) dt}{t_{end} - t_{start}}$

**% TimeNotZero** ( $th, j, k$ ) =  $\frac{\int_{t_{start}}^{t_{end}} \delta_{th,j,k}(t) dt}{\sum_j \int_{t_{start}}^{t_{end}} \delta_{th,j,k}(t) dt}$

**NumBurst** ( $th, j, k$ ) =  $i_{end} - i_{start} + 1$   
 $i_{start} = \min(i) | t_i > t_{start}, i_{end} = \max(i) | t_i < t_{end}$

**Integral** ( $th, j, k$ ) =  $\int_{t_{start}}^{t_{end}} S_{th}^{dw}(t) \delta_{th,j,k}(t) dt$

**Average** ( $th, j, k$ ) =  $\frac{\int_{t_{start}}^{t_{end}} S_{th}^{dw}(t) \delta_{th,j,k}(t) dt}{\int_{t_{start}}^{t_{end}} \delta_{th,j,k}(t) dt}$

**Maximum** ( $th, j, k$ ) =  $\max(S_{th}^{dw}(t) \delta_{th,j,k}(t), t = [t_{start}, t_{end}])$

**SumBurst** ( $th, j, k$ ) =  $\sum_{i_{start}}^{i_{end}} S_{th}^{dw}(i) \delta_{th,j,k}(t_i), i_{start} = \min(i) | t_i > t_{start}, i_{end} = \max(i) | t_i < t_{end}$

# Distributed Configurations

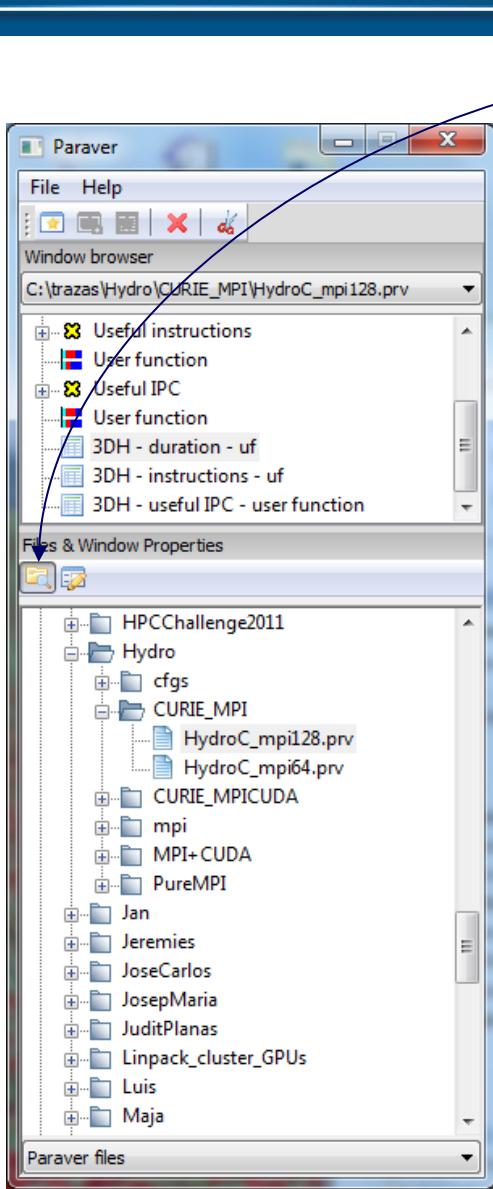
# Distribution of cfg directories

- CFG \$PARAVER\_HOME/cfgs
  - General
    - including basic views (timelines) and analyses (2/3D profiles), including views of the user functions and call-stack
  - Counters\_PAPI
    - Hardware counter derived metrics. Grouped in directories for
      - Program: related to algorithmic/compilation (i.e. instructions, FP ops,...)
      - Architecture: related to execution on specific architectures (i.e. cache misses,...)
      - Performance: metrics reporting rates per time (i.e. MFlops, MIPS, IPC,...)
  - MPI
    - Grouped in directories displaying views and analysis. Further separated into point to point and collectives.
  - OpenMP
    - Grouped in directories displaying views and analysis

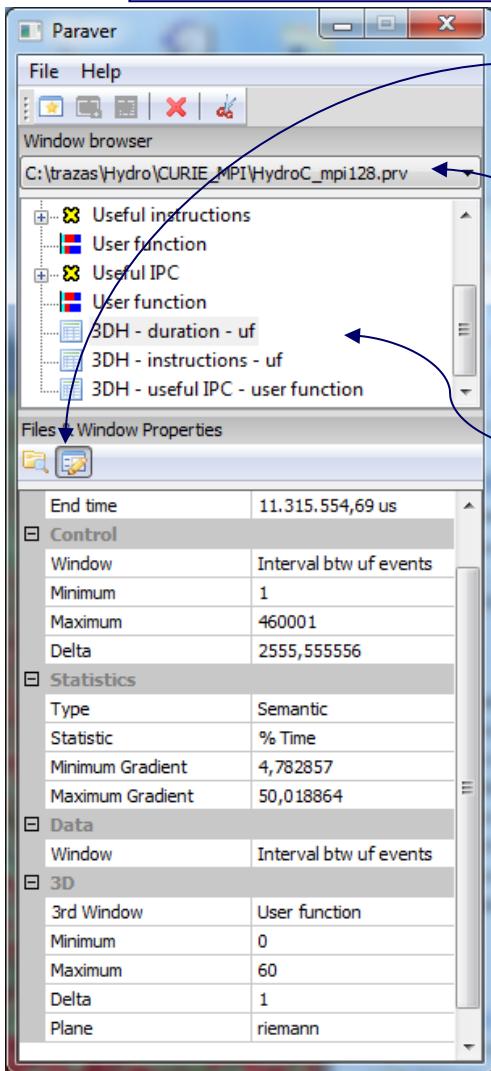
How to ...

# Main Paraver window

VI-HPS



Select to browse in lower panel for traces or cfs



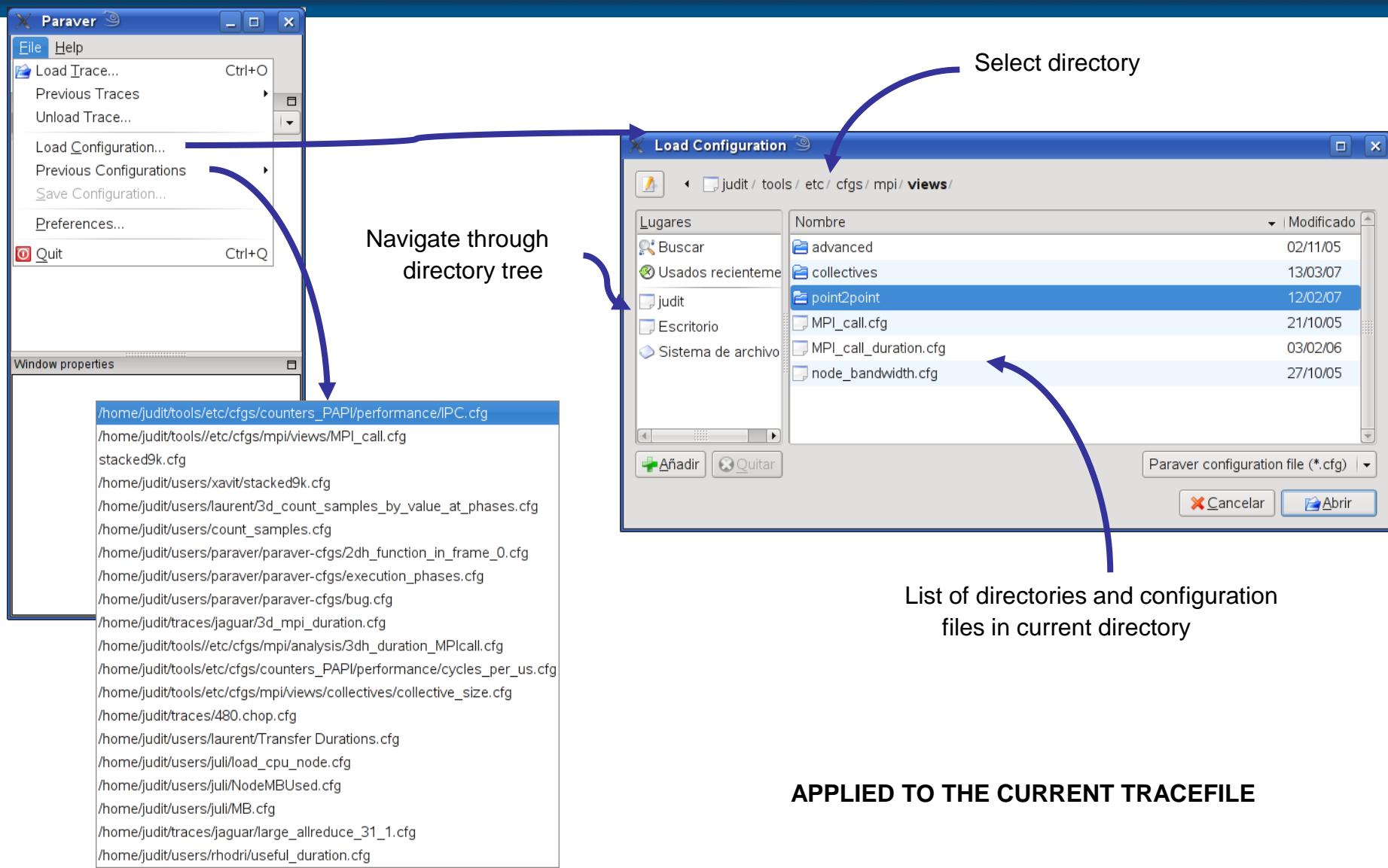
Select to browse characteristics of active view or table

Active trace

Available views and tables  
Active view or table highlighted

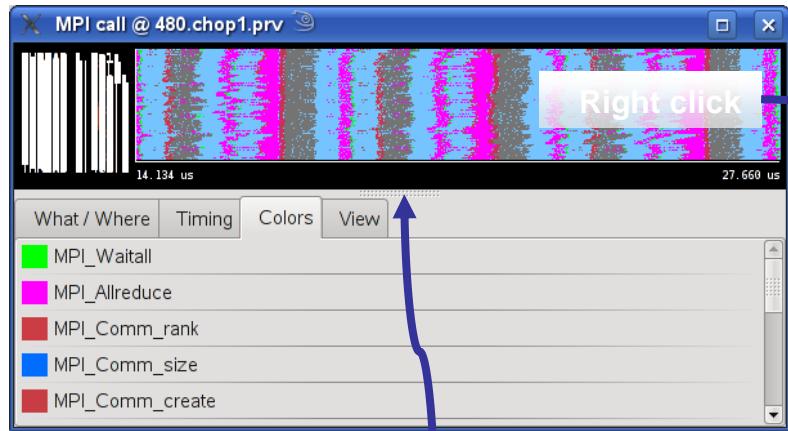
# Load configuration files

VI-HPS

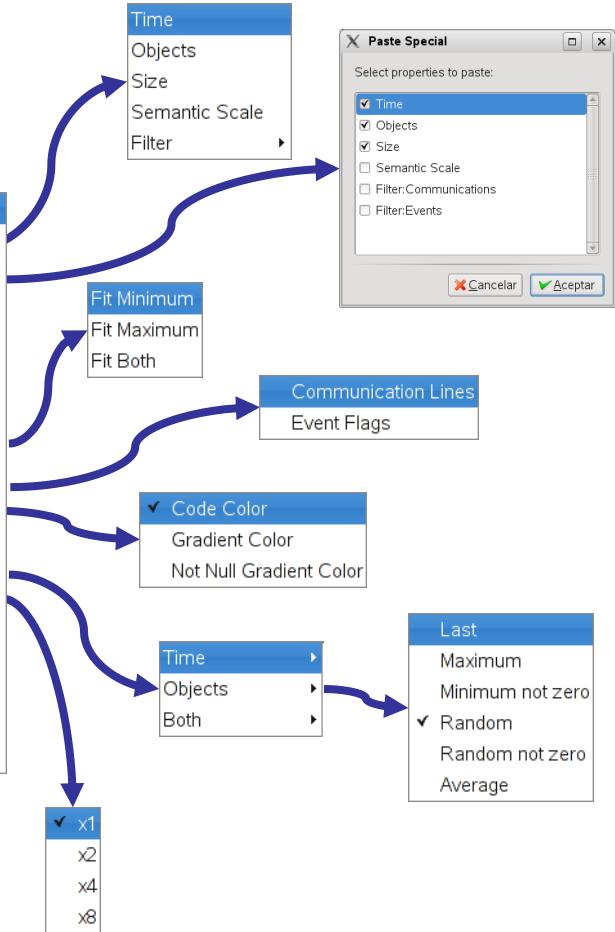
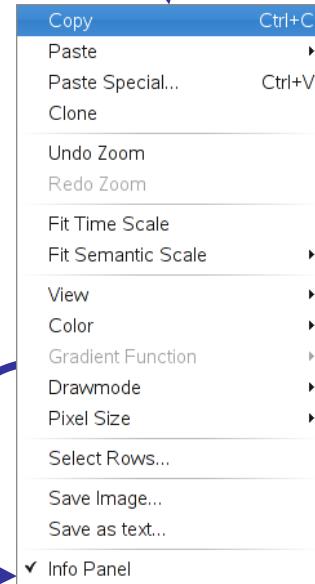


# Navigation

VI-HPS



Hide lower panel  
(double click)



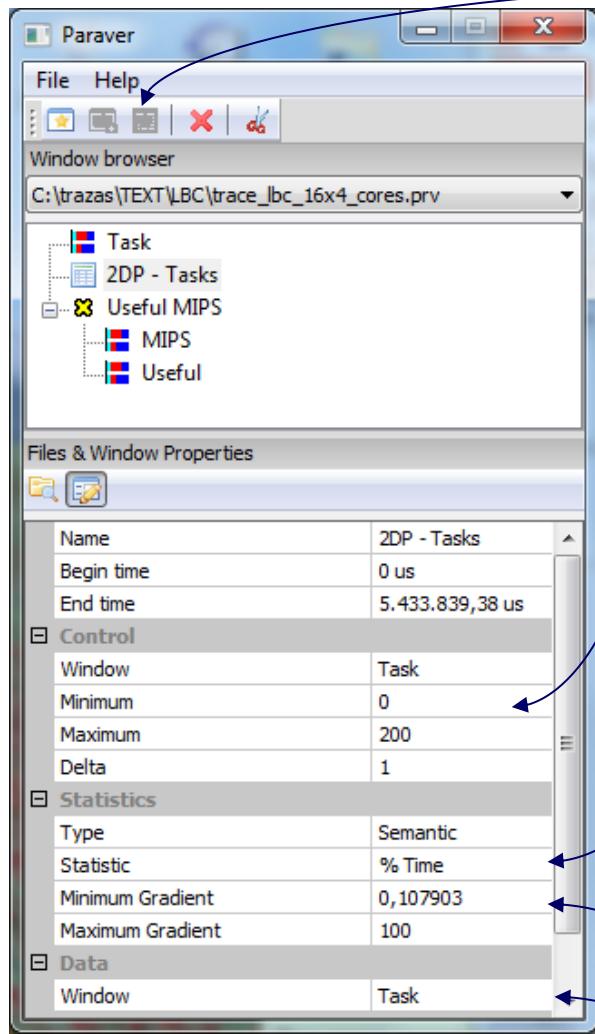
## Shortcuts:

Drag and move (D&M) – Zoom

Control D&M – Zoom XY

Shift D&M – Timing

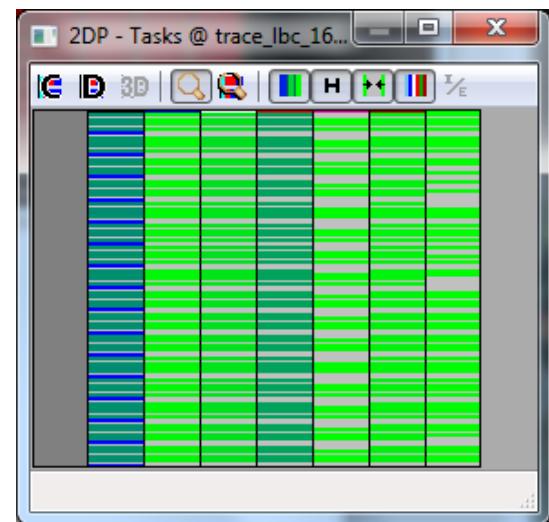
# How to generate table and change statistic



To generate table: click button and select region of the window whose values will determine the columns of the table

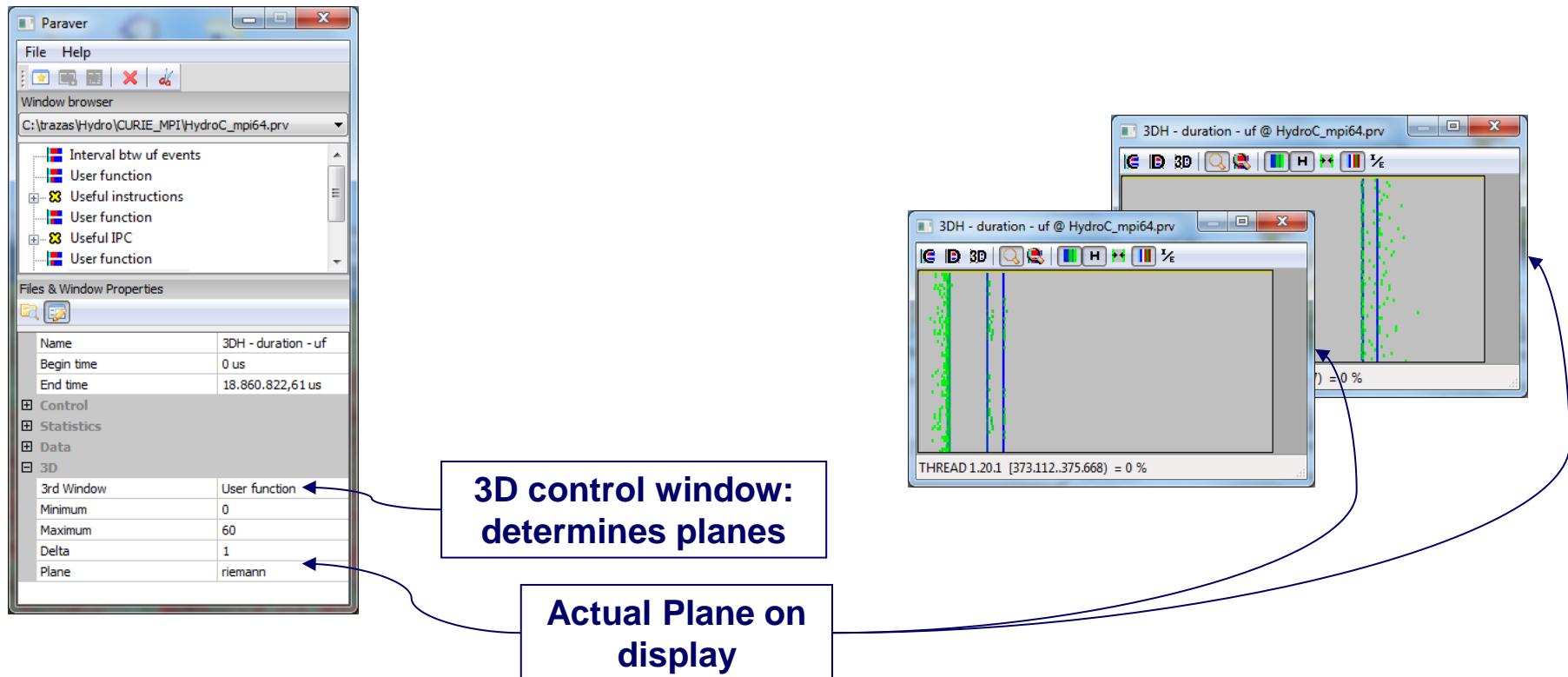
Range and bin width (delta) represented by each column. By default is automatically selected, but can be manually changed

Selection of statistic to appear in each cell



Cell coloring gradient control  
Window used to compute statistic (only used by some statistics)

- One additional dimension
  - One plane per value of a 3D control window
- Useful to categorize histograms
  - i.e. histogram of duration of specific user function



# Table information and control

VI-HPS

Create a new table

Display whole table / cell text

Color/not cells

Transpose

Hide null columns

Region analyzed

Bin definition

Change Data window

Activate 3D analysis

Color encoding

max

min

The screenshot shows the Paraver interface with two main windows. On the left is the 'Window browser' showing a file path and two selected items: 'MPI call' and 'MPI call profile'. Below it are 'Window properties' and 'Control' sections. The 'Control' section includes fields for 'Window' (set to 'MPI call'), 'Minimum' (0), 'Maximum' (124), and 'Delta' (1). The 'Statistics' section shows 'Type' as 'Semantic' and 'Statistic' as '% Time'. The 'Data' section is expanded, showing 'Minimum Gradient' (1,188944) and 'Maximum Gradient' (49,410556). The '3D' section is also visible. On the right is a table titled 'MPI call profile @ 481.ch01.prv'. The table has columns: End, MPI\_Waitall, MPI\_Allreduce, MPI\_Comm\_rank, and MPI\_Waitany. The rows list various threads from THREAD 1.1.1 to THREAD 1.16.1. The table uses color coding where most cells are green, except for some blue ones. The top row is highlighted in blue. The bottom of the table shows the text 'MPI\_Waitany'.

	End	MPI_Waitall	MPI_Allreduce	MPI_Comm_rank	MPI_Waitany
THREAD 1.1.1	44,03 %	2,08 %	29,03 %	2,03 %	9,59 %
THREAD 1.2.1	48,61 %	1,81 %	6,23 %	2,40 %	11,68 %
THREAD 1.3.1	48,62 %	2,04 %	6,60 %	1,99 %	11,25 %
THREAD 1.4.1	48,59 %	1,83 %	6,41 %	2,58 %	11,48 %
THREAD 1.5.1	48,30 %	1,83 %	6,36 %	2,61 %	11,57 %
THREAD 1.6.1	48,40 %	1,82 %	6,55 %	2,60 %	11,37 %
THREAD 1.7.1	48,37 %	2,23 %	7,82 %	1,90 %	10,67 %
THREAD 1.8.1	48,54 %	2,08 %	7,08 %	2,13 %	10,89 %
THREAD 1.9.1	47,89 %	3,10 %	10,69 %	1,44 %	8,49 %
THREAD 1.10.1	48,09 %	2,82 %	8,62 %	1,52 %	9,93 %
THREAD 1.11.1	48,60 %	2,51 %	9,02 %	1,50 %	9,80 %
THREAD 1.12.1	48,76 %	2,00 %	6,76 %	2,26 %	10,00 %
THREAD 1.13.1	44,08 %	3,73 %	28,53 %	2,51 %	8,17 %
THREAD 1.14.1	48,94 %	1,91 %	8,35 %	2,29 %	12,02 %
THREAD 1.15.1	48,81 %	1,94 %	8,27 %	2,46 %	11,91 %
THREAD 1.16.1	49,07 %	1,95 %	8,38 %	2,42 %	11,74 %

# Table information and control

VI-HPS

Open Data window

Open Control window

Open 3D window

Generate a timeline, derived from control window with the range of values selected clicking in the table (zoom mode only)

Right click

Generate ASCII file with table data

Shortcuts (zoom mode only):  
Drag and move (D&M) – Zoom  
Control D&M – Zoom XY

Selected plane

Copy Ctrl+C

Paste Ctrl+V

Paste Special...

Clone

Undo Zoom

Redo Zoom

Fit Time Scale

✓ Auto Fit Control Scale

✓ Auto Fit 3D Scale

✓ Auto Fit Data Gradient

Gradient Function

Drawmode

Save as text...

Time

Objects

Size

Semantic Scale

Control scale

3D scale

Paste Special

Select properties to paste:

Time

Objects

Size

Semantic Scale

Cancel Acceptar

Linear

Steps

Logarithmic

Exponential

Semantic

Objects

Both

File Help

Window browser

/home/judit/traces/480.chop1.prv

Window properties

Name MPI call duration

Begin time 15.500,13 us

End time 25.867,11 us

Control

Window MPI call duration

Minimum 25,719

Maximum 4024,328

Delta 19,993045

Statistics

Type Semantic

Statistic Time

Minimum Gradient 25,849453

Maximum Gradient 2305,542

Data

Window MPI call duration

3D

3rd Window MPI call

Minimum 6

Maximum 124

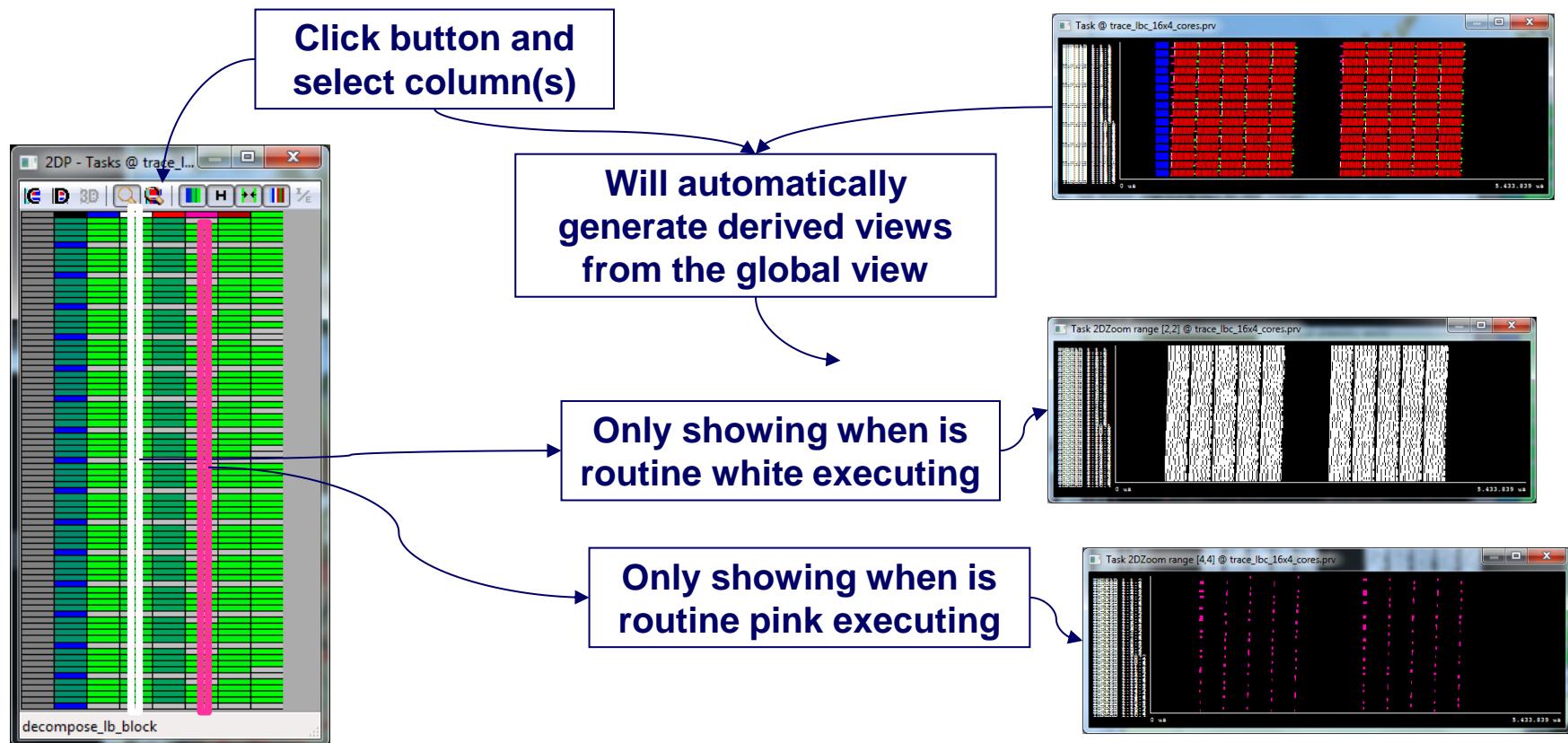
Delta 1

Plane MPI\_Allreduce

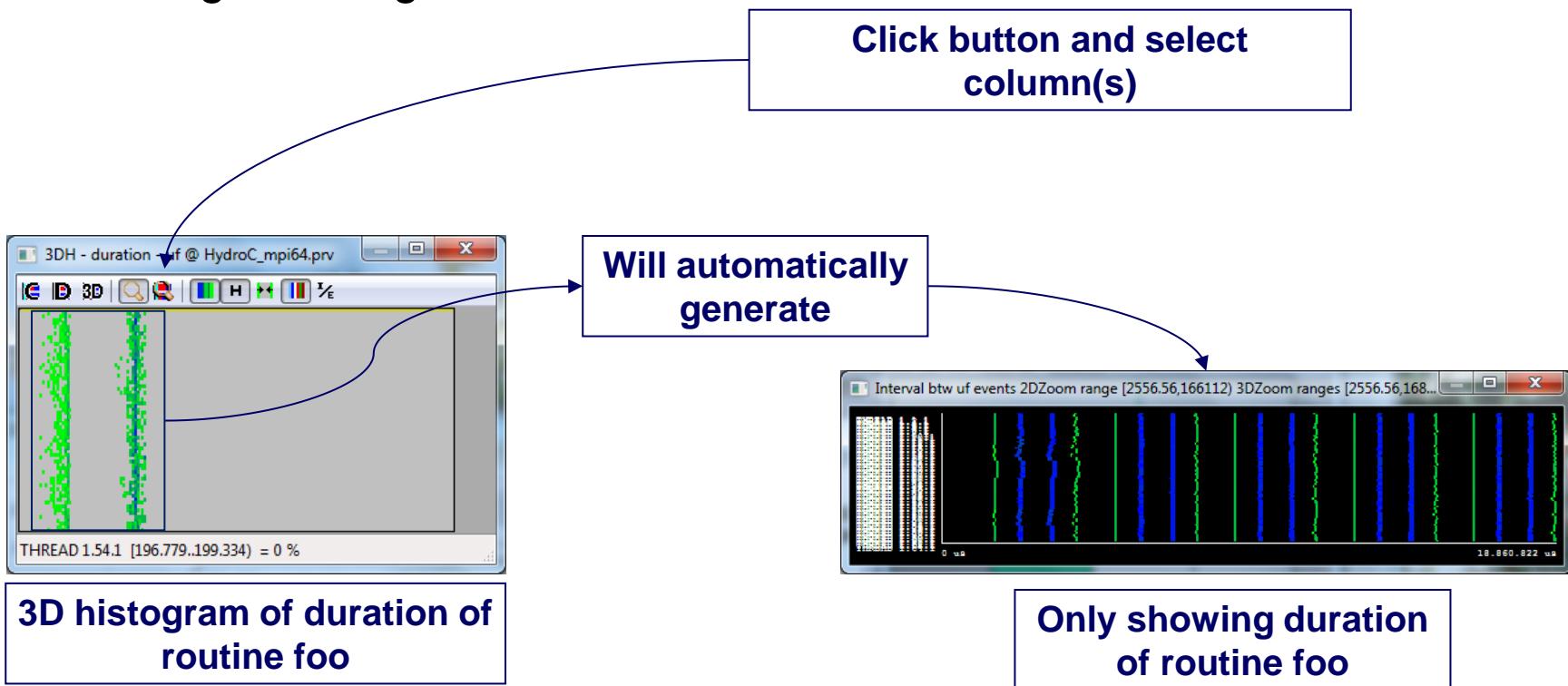
MPI call duration @ 480.chop1.pvt

THREAD 1.3.1 [1845,09..1865,08) = 0 us

- Where in the timeline do the values in certain table columns appear?
  - ie. want to see the time distribution of a given routine?



- Where in the timeline do the values in certain table columns appear?
  - ie. want to see where the timeline happen computation bursts of a given length?



# Trace manipulation

- Paraver data handling utilities
  - If trying to load a very large trace, Paraver will ask if you want to filter it
- Three steps:
  - Filter original trace discarding most of the records only keeping most relevant information (typically computation bursts longer than a given lower bound)
  - Analyze coarse grain structure of trace. Typically useful\_duration.cfg
  - Cut original trace to obtain a fully detailed trace for the time interval considered representative or of interest

Guided hands-on available in

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

Trace to which it will be applied

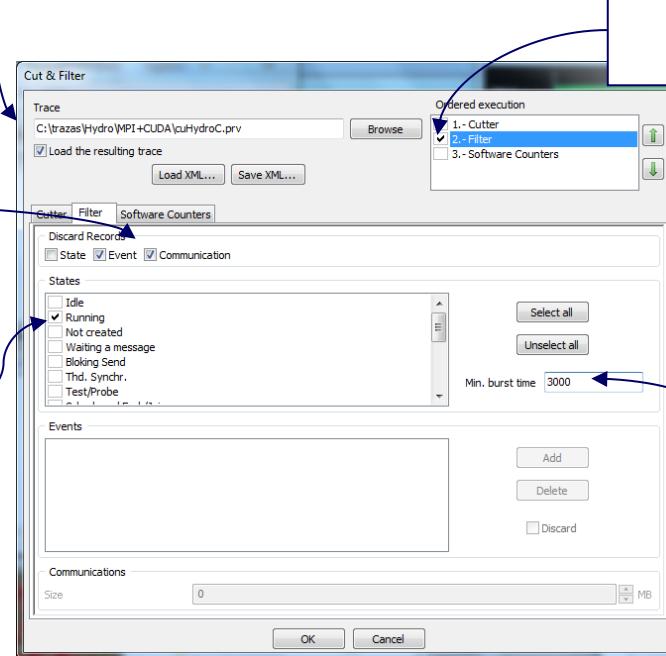
A trace with  
basename.filter1.prv will be  
generated

Discard events and  
communications

Keep only Running bursts  
....

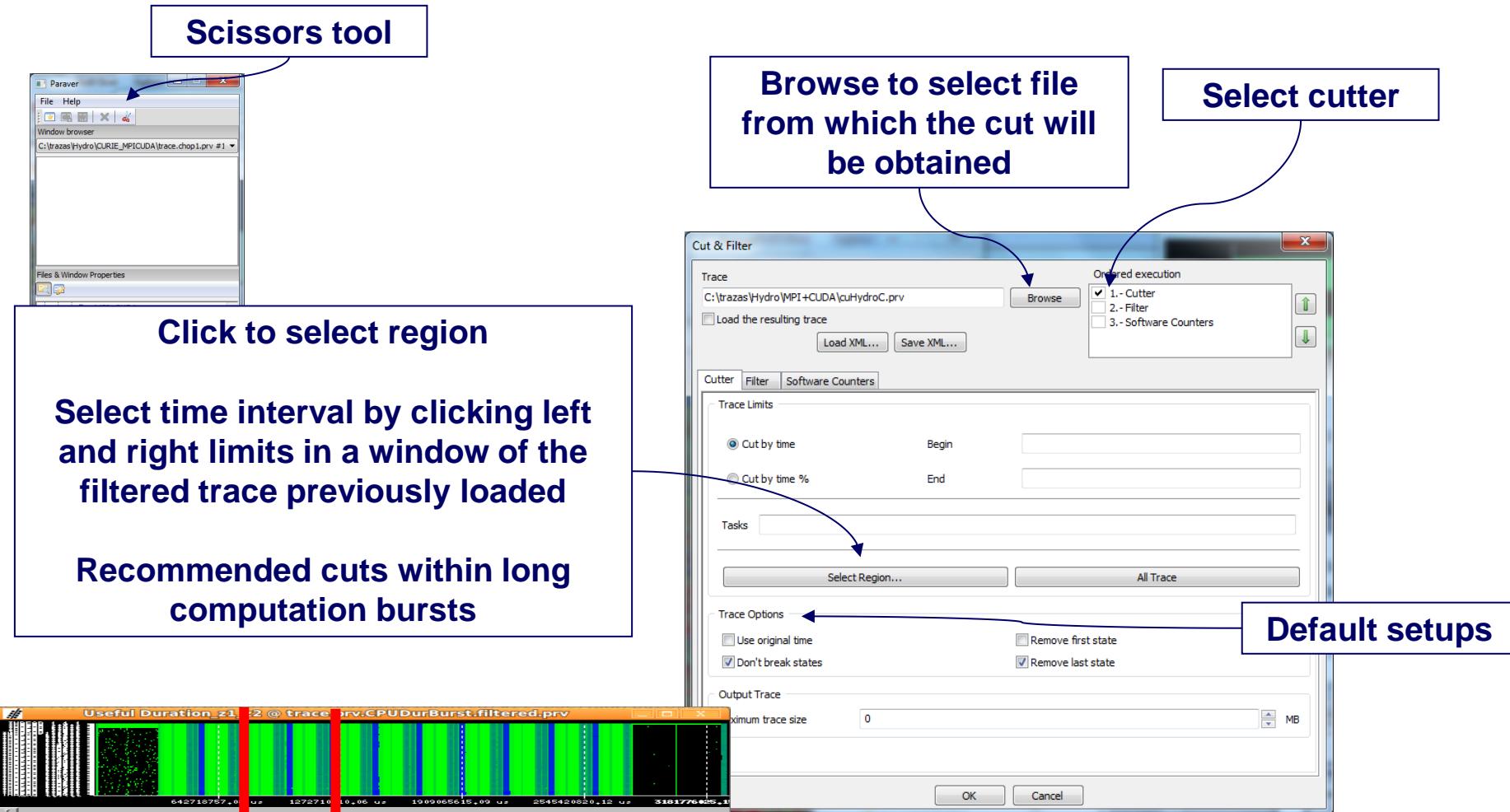
Select filtering  
option

--- longer than  
3000 ns



# Cutting very large traces

- Load a filtered trace and use the scissors tool



# Extrae

# Adapt job submission script

```
#!/bin/bash

export NP=8
export INPUT=$1

cleo-submit -np $NP          ./HydroC -i $INPUT
```

**appl.job**

# Adapt job submission script

```
#!/bin/bash

export NP=8
export INPUT=$1

cleo-submit -np $NP ./trace.sh ./HydroC -i $INPUT
```

**appl.job**

```
#!/bin/bash

export EXTRAE_HOME=/export/hopsta/BSCTools/tools/extrae-2.3
export EXTRAE_CONFIG_FILE=extrae/extrae.xml

export LD_PRELOAD=$EXTRAE_HOME/lib/libmpitrace.so

export EXE=$1
export TRACENAME=${EXE}_$3.prv

$@
```

**trace.sh**

```
<?xml version='1.0'?>                                extrae.xml  
  
<trace enabled="yes"  
    home="/home/judit/tools/extrae-2.3"  
    initial-mode="detail"  
    type="paraver"  
    xml-parser-id="Id: xml-parse.c 799 2011-10-20 16:02:03Z harald $"  
>  
  
<mpi enabled="yes">  
    <counters enabled="yes" />          ← Activate MPI tracing and emit hardware counters at MPI calls  
</mpi>  
  
<openmp enabled="no">           ← Do not activate OpenMP tracing  
    <locks enabled="no" />  
    <counters enabled="yes" />  
</openmp>  
  
<callers enabled="yes">           ← Emit call stack information (number of levels) at acquisition points  
    <mpi enabled="yes">1-3</mpi>  
    <sampling enabled="no">1-5</sampling>  
</callers>  
  
...
```

Details in \$EXTRAE\_HOME/share/example/MPI/extrae\_explained.xml

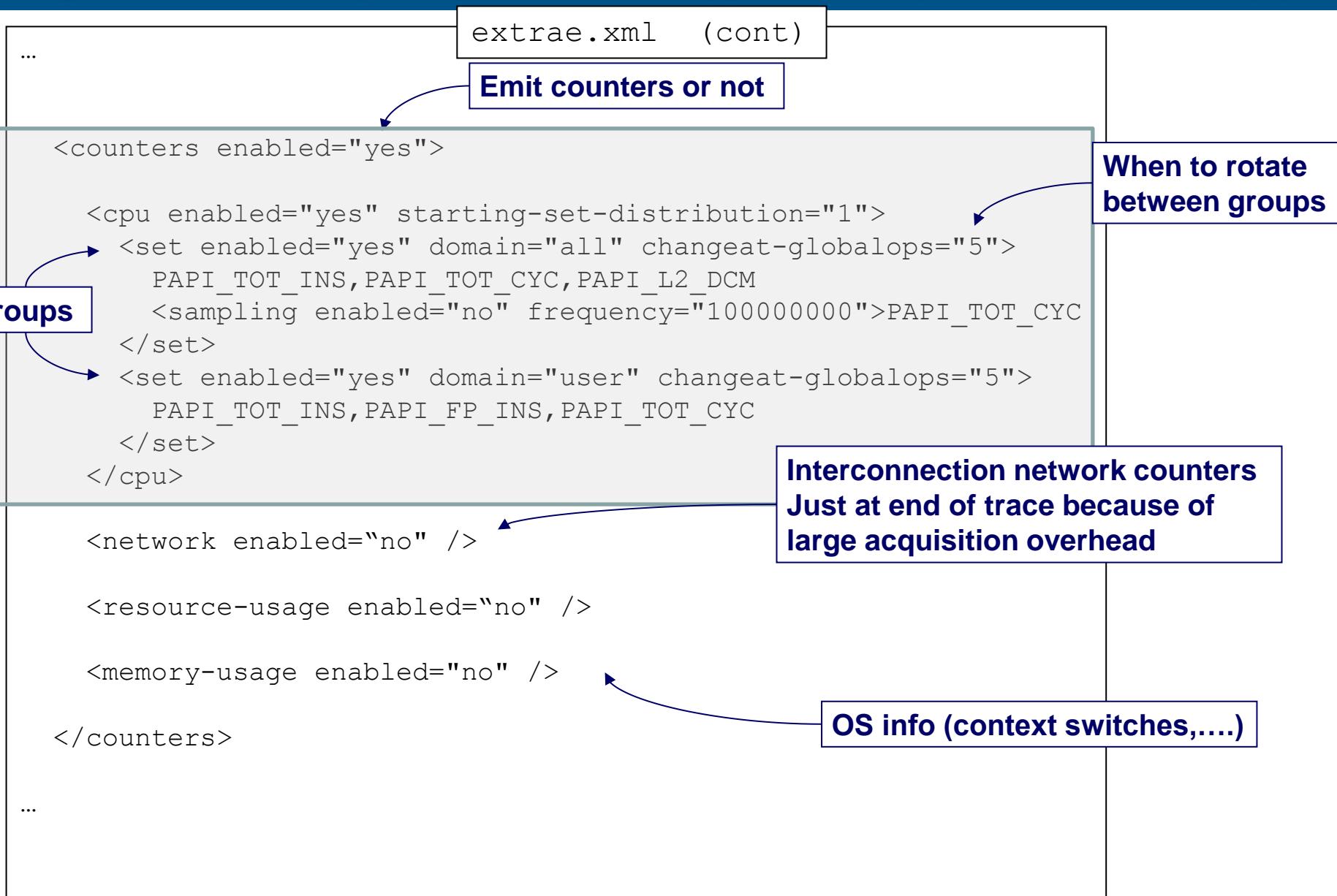
## extrae.xml (cont)

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

**Add instrumentation at specified user functions  
Requires Dyninst based mpitrace**

...

# Trace control .xml (cont)



# Trace control .xml (cont)

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</final-directory>
  <gather-mpits enabled="no" />
</storage>

<buffer enabled="yes">
  <size enabled="yes">500000</size>
  <circular enabled="no" />
</buffer>
```

**Control of emitted trace ...**

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

...

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

**Size of in core buffer (#events)**

## Trace control .xml (cont)

extrae.xml (cont)

...

```
<trace-control enabled="yes">
  <file enabled="no" frequency="5m"/>/gpfs/scratch/bsc41/bsc41273/control</file>
  <global-ops enabled="no"></global-ops>
  <remote-control enabled="no">
    <signal enabled="no" which="USR1"/>
  </remote-control>
</trace-control>
```

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

```
<others enabled="no">
  <minimum-time enabled="no">10M</minimum-time>
  <terminate-on-signal enabled="no">USR2</terminate-on-signal>
</others>
```

**Stop tracing after elapsed time ...**

...

**... or when signal received**

## Trace control .xml (cont)

extrae.xml (cont)

```
<bursts enabled="no">  
  <threshold enabled="yes">500u</threshold>  
  <counters enabled="yes" />  
  <mpi-statistics enabled="yes" />  
</bursts>
```

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

**... plus summarized MPI events**

```
<sampling enabled="no" type="default" period="5m" />
```

**Activate/not time based sampling and how often**

extrae.xml (cont)

...

```
<merge enabled="yes" ←  
    synchronization="default"  
    binary="$EXE$"  
    tree-fan-out="16"  
    max-memory="512"  
    joint-states="yes"  
    keep-mpits="yes"  
    sort-addresses="yes"  
>  
    $TRACENAME$ ←  
</merge>  
  
</trace>
```

Merge individual traces into global  
application trace at end of run ...

... into this trace name

- Library depends on programming model

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

<sup>1</sup> for Fortran codes

# Using Dimemas

- Paraver → Dimemas trace Generation
  - `prv2dim original.prv dimemas.dim`
  - Default: duration of each computation region taken from .prv computation duration

Usage:

```
prv2dim -i <iprobe_miss_threshold> -b <hw_counter_type>,<factor>
<paraver_trace> <dimemas_trace>

-h This help
-n No generate initial idle states
-i <iprobe_miss_threshold> Maximum MPI_Iprobe misses to discard Iprobe
area CPU burst
-b <hw_counter_type>,<factor> Hardware counter type and factor used to
generate burst durations
...

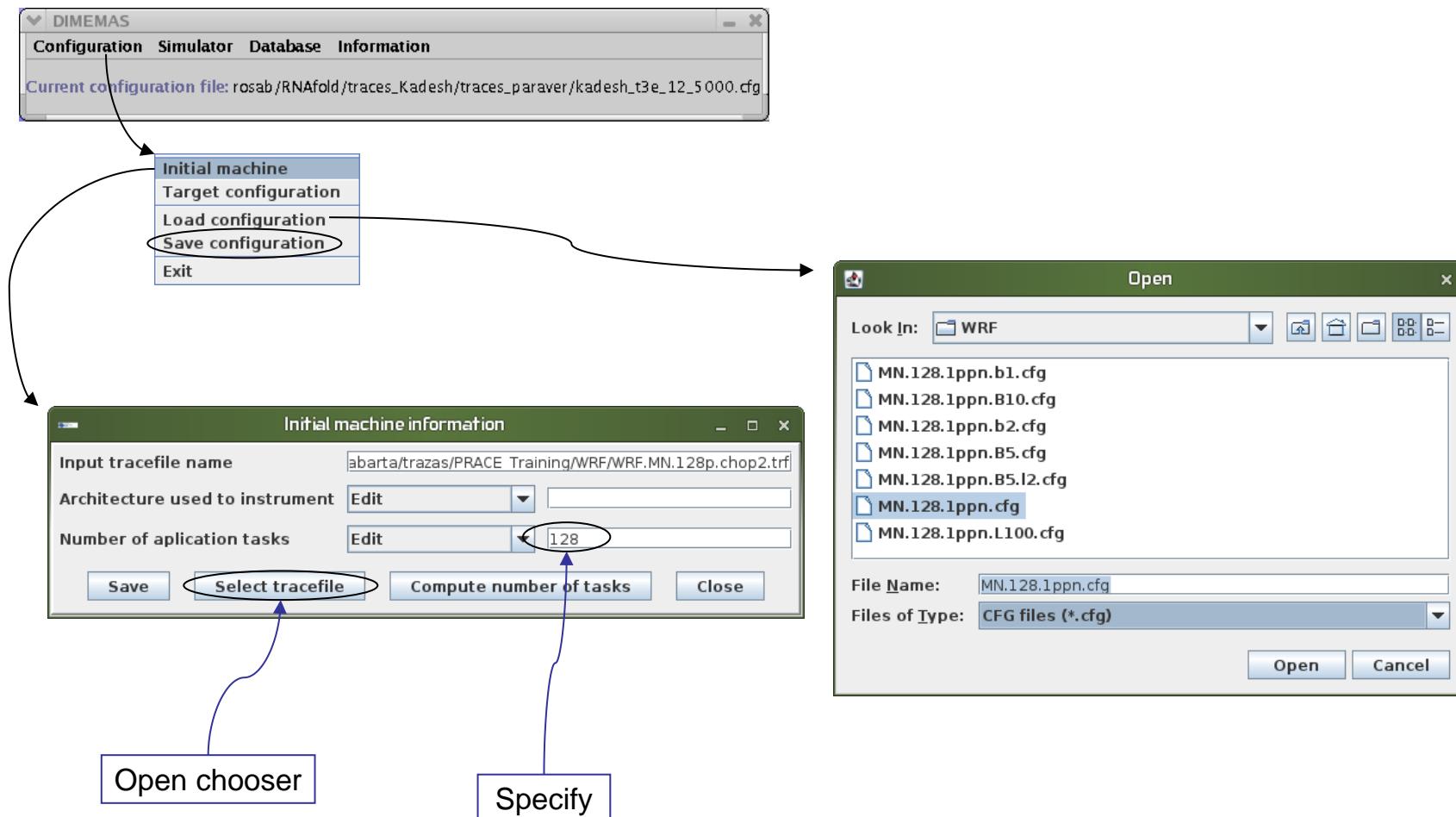
```

Force synchronized start of all threads

Computation region duration derived from hardware counters assuming/modeling a given performance (<factor>) i.e. estimate impact if we get rid of IPC imbalances

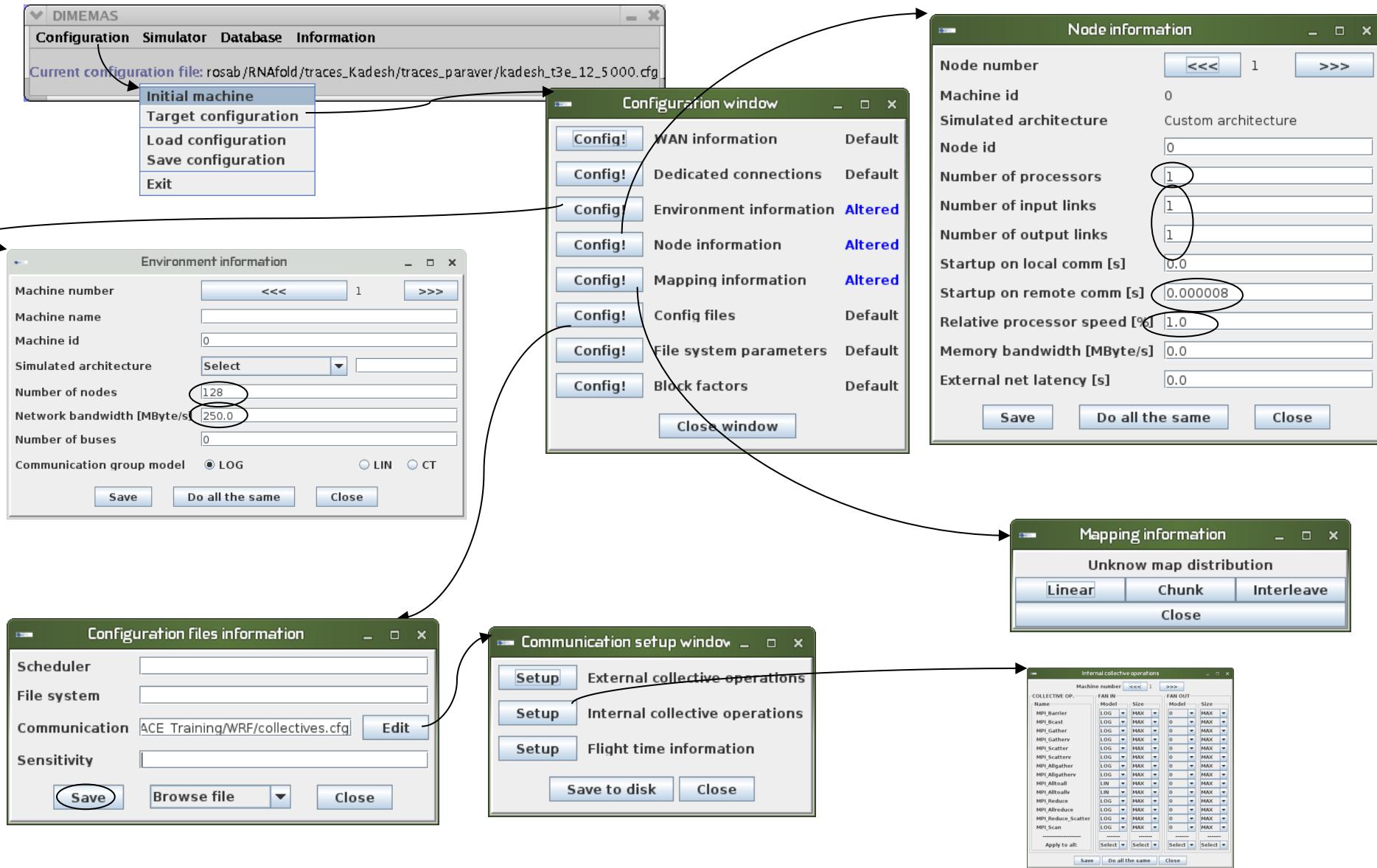
# Dimemas GUI – Specify trace to simulate

VI-HPS



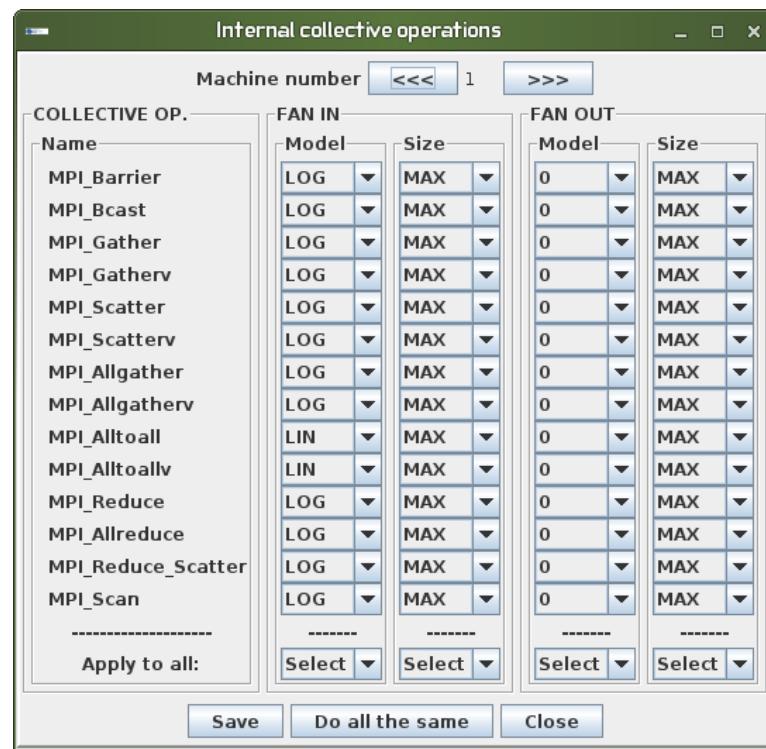
# Dimemas GUI – Specify target machine

VI-HPS



# Collective Communication Model

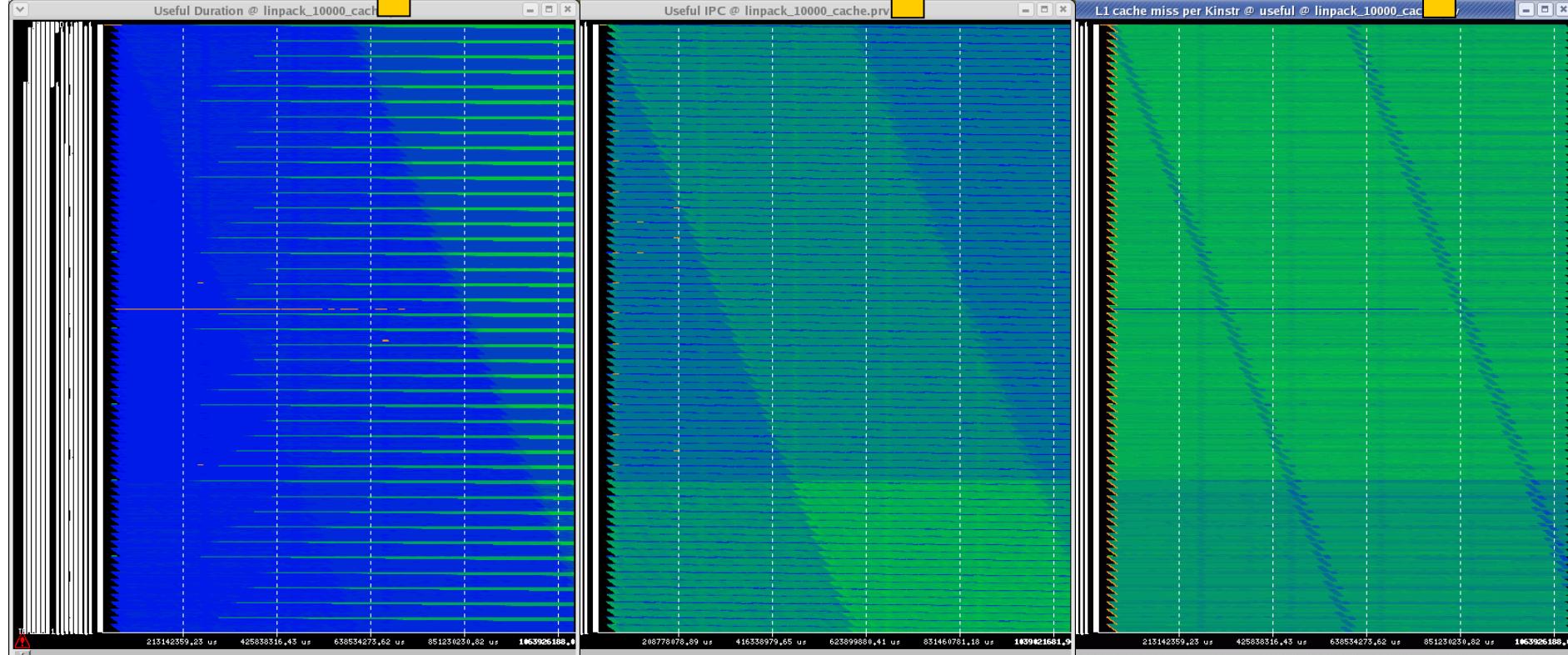
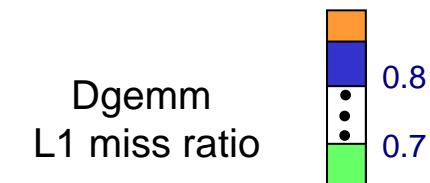
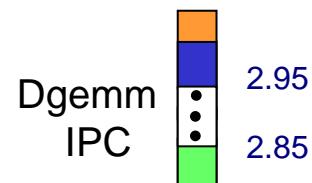
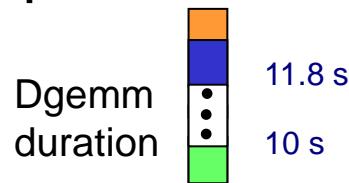
- Per call model
  - Model factor
    - Lin
    - Log
    - Const
  - Size of message
    - Min over all processes
    - Mean over all processes
    - Max over all processes
  - Specified in input file



# Scalability

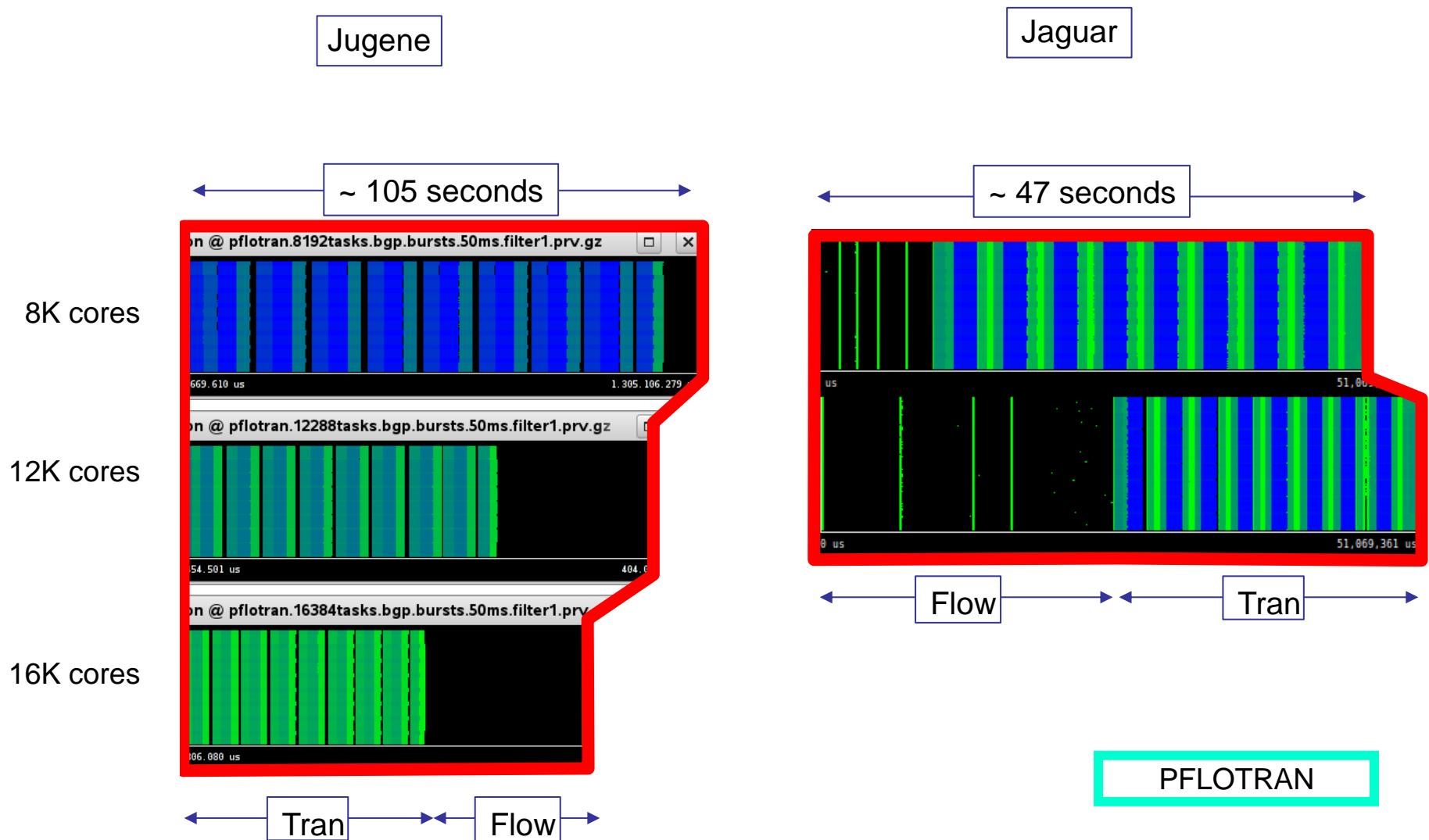
# Scalability of Presentation

- Linpack @ Marenostrum: 10k cores x 1700 s



# Scalability of analysis

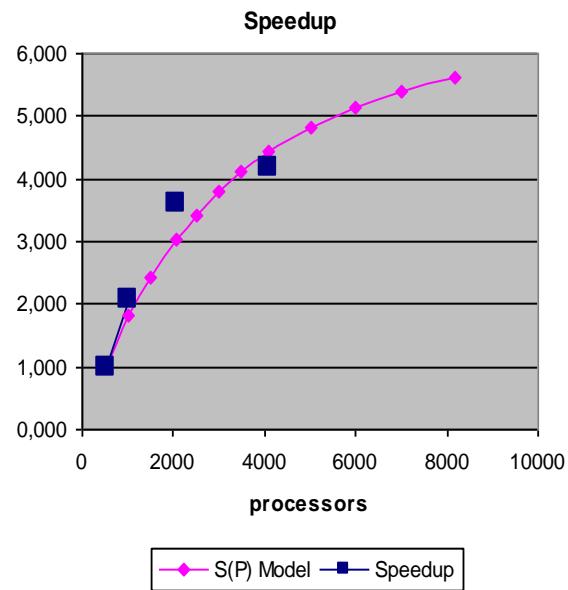
VI-HPS



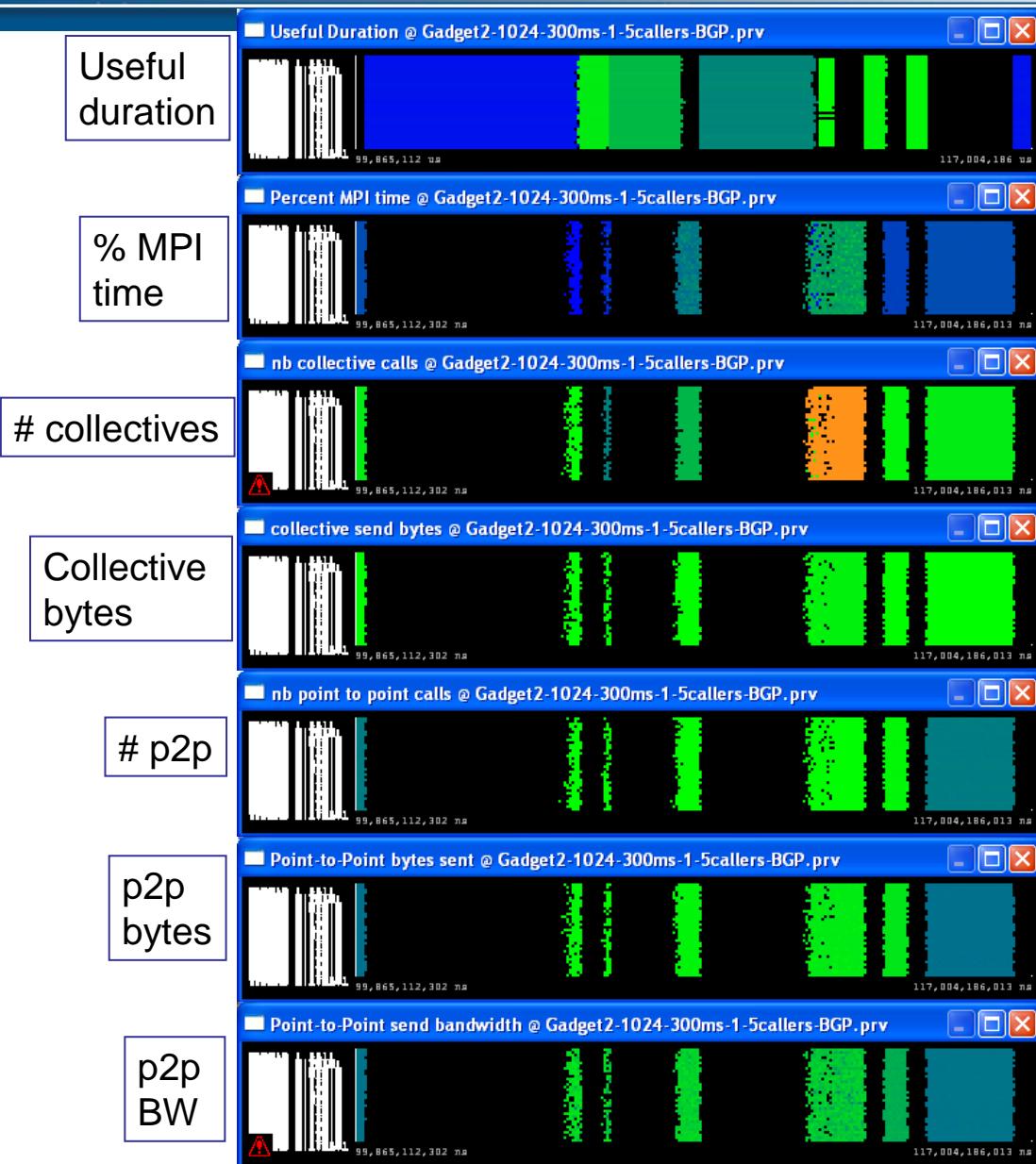
- Software counters
  - Summarize information of some event types (ie. MPI calls) by emitting aggregate counts
  - Emit counts at structurally relevant points (i.e. begin and end of long computation phases)
- Representative cuts
  - Emit full detail only on selected intervals, representative of full program execution
- On and off line combinations
  - By instrumentation
  - By paraver filtering

# Software counters

VI-HPS

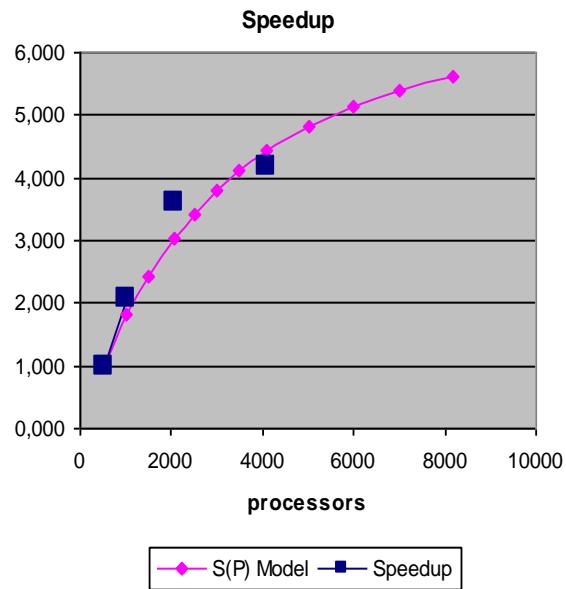


GADGET, PRACE Case A, 1024 procs

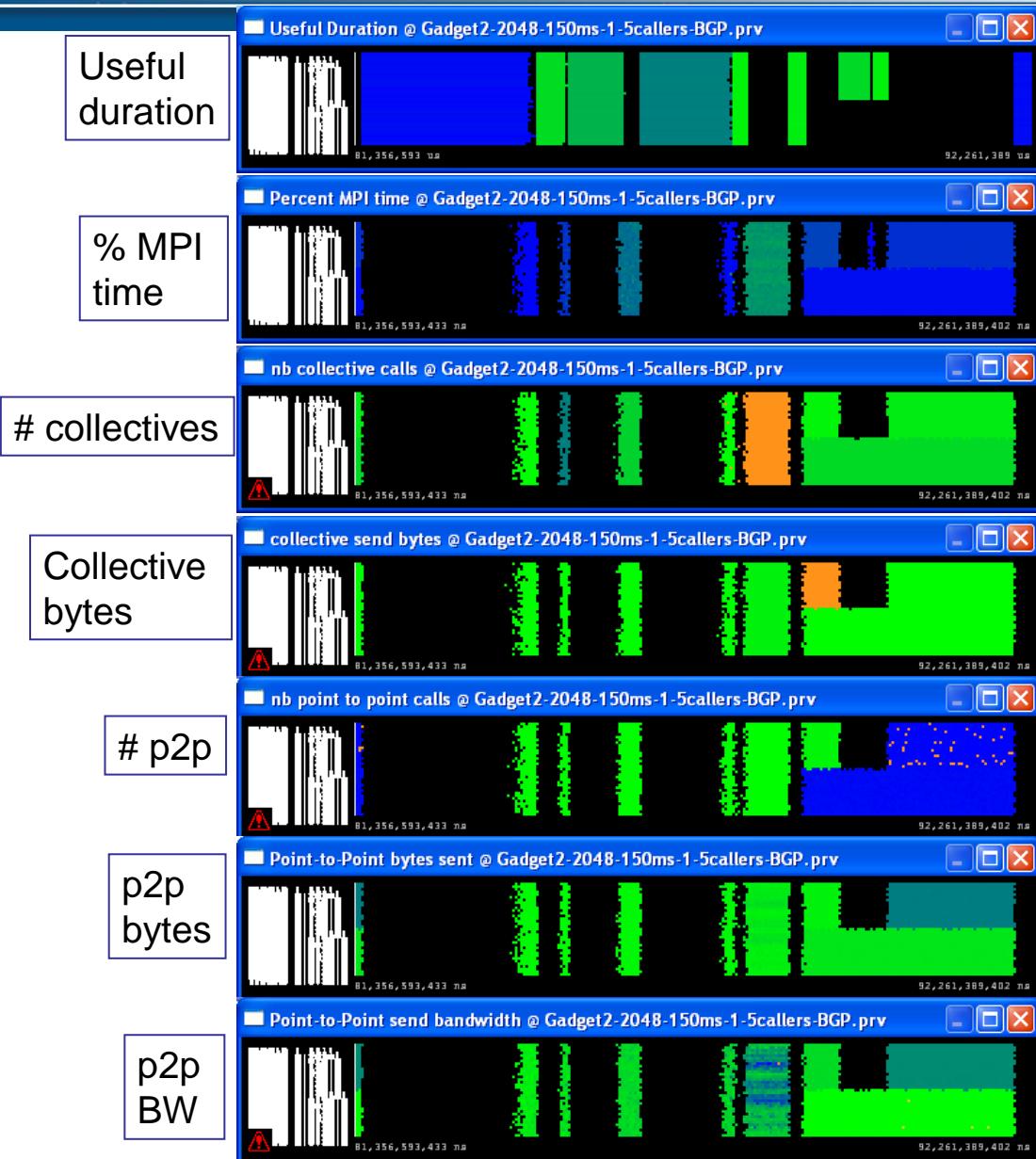


# Software counters

VI-HPS

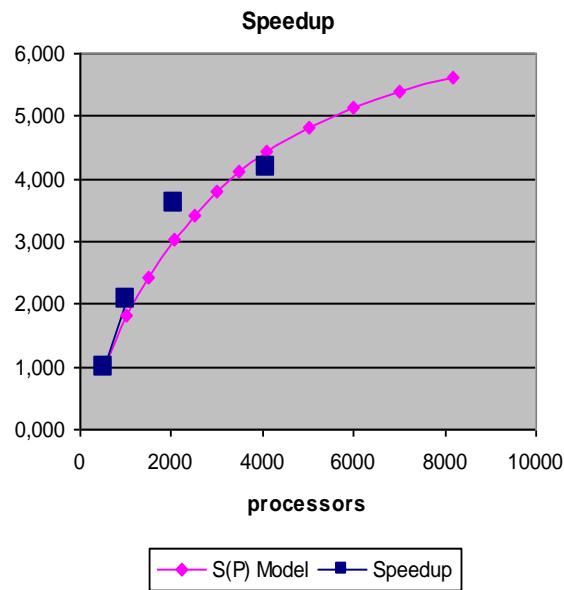


GADGET, PRACE Case A, 2048 procs



# Software counters

VI-HPS

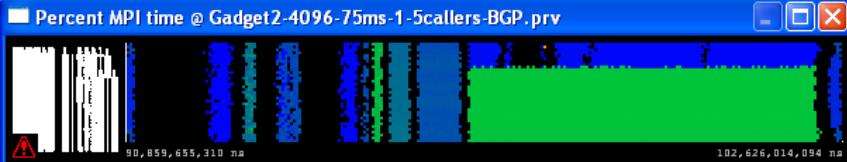


GADGET, PRACE Case A, 4096 procs

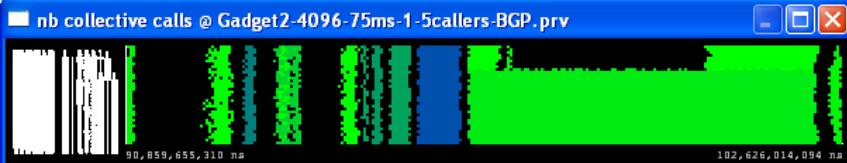
Useful duration



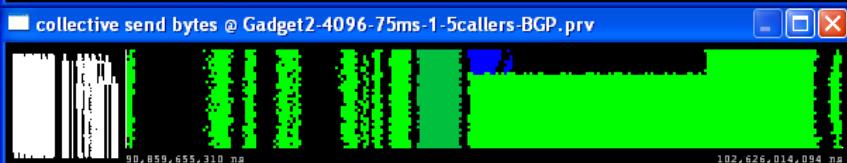
% MPI time



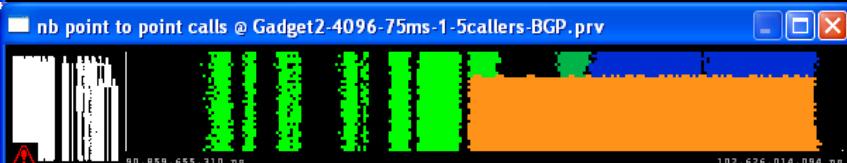
# collectives



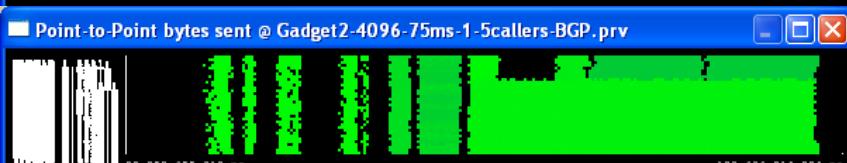
Collective bytes



# p2p



p2p bytes



p2p BW

