

www.bsc.es



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Understanding applications performance with Paraver

Judit Gimenez
BSC

PATC Parallel Programming Workshop
Oct 14-18 2013

Our Tools

« Since 1991

« Based on traces

« Open Source

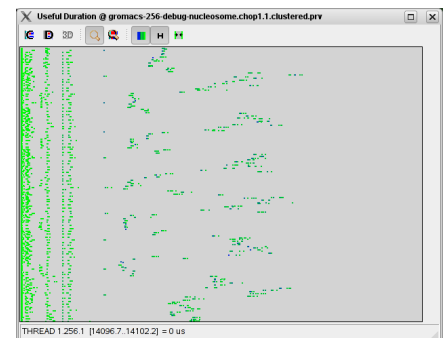
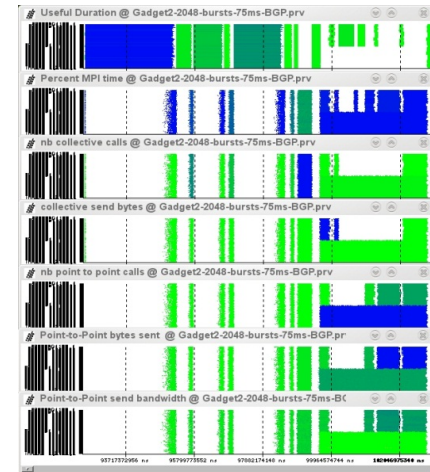
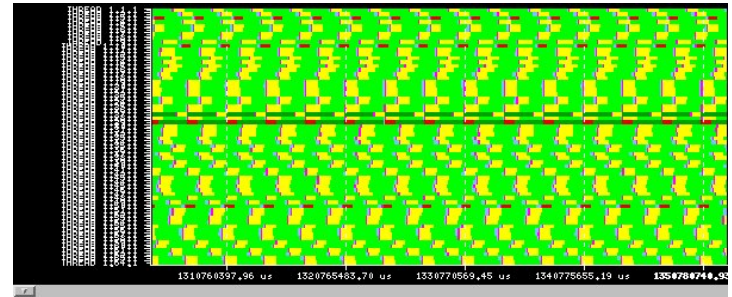
- <http://www.bsc.es/paraver>

« Core tools:

- Paraver (paramedir) – offline trace analysis
- Dimemas – message passing simulator
- Extrae – instrumentation

« Focus

- Detail, flexibility, intelligence



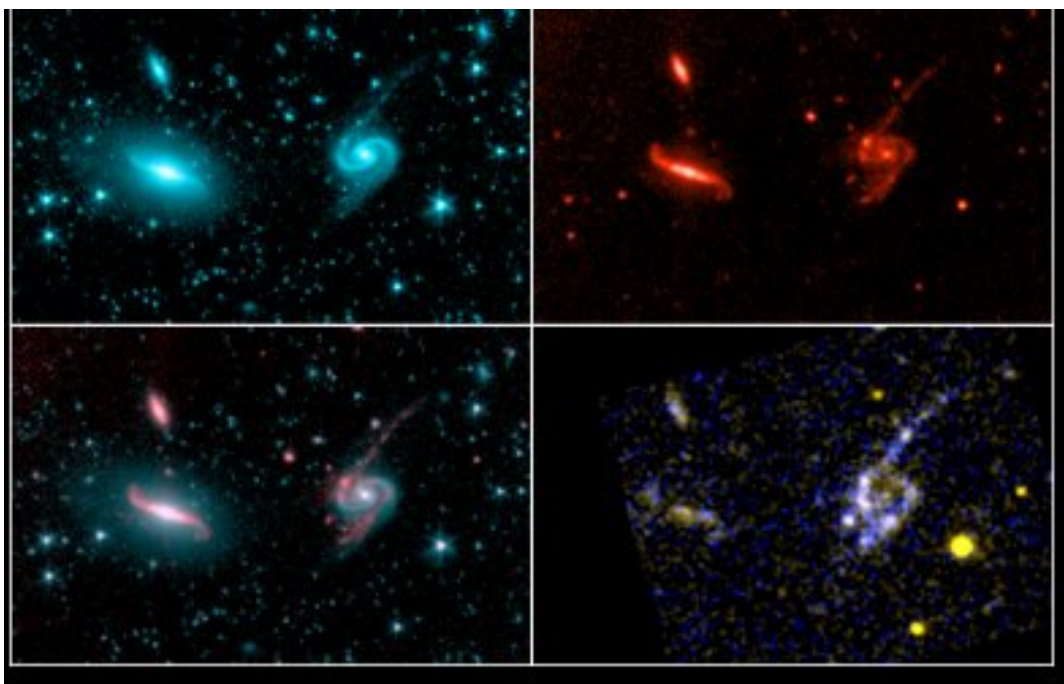
Multispectral imaging

« Different looks at one reality

- Different spectral bands (light sources and filters)

« Highlight different aspects

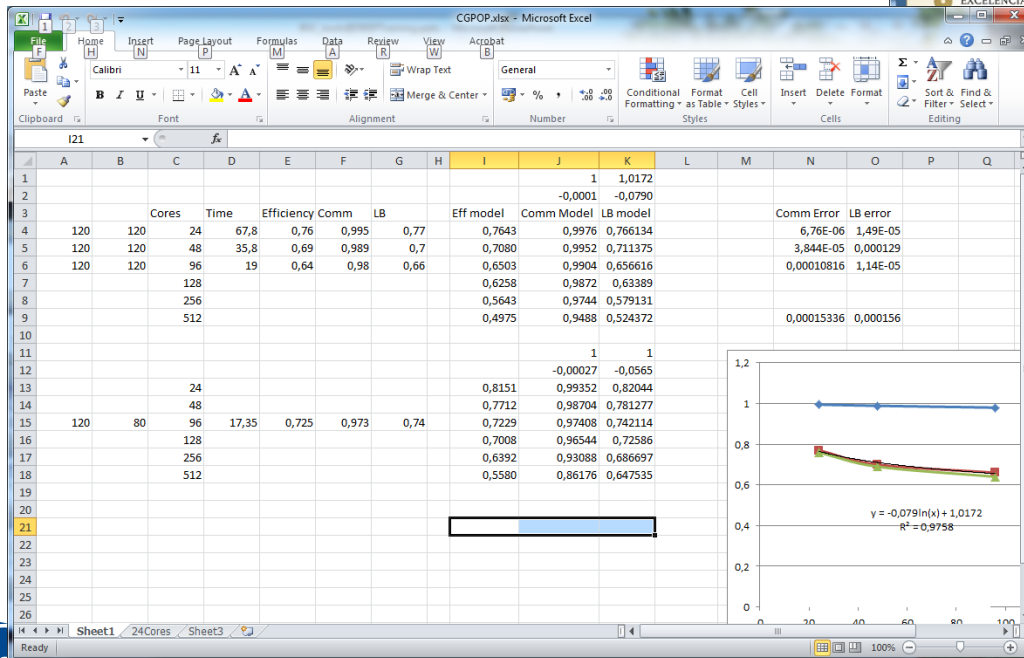
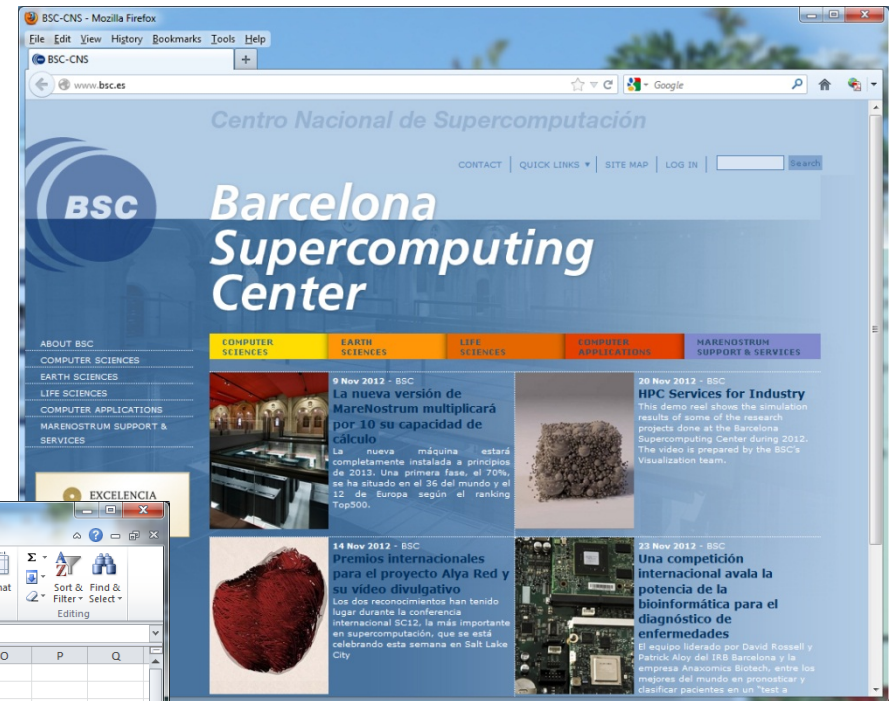
- Can combine into false colored but highly informative images



Spreadsheets and browsers

Display, manipulate data

- Dynamic content
- User defined operations



A “different” view on performance analysis and tools

⌘ 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.

“That what is simple is rarely understood”
my iPads Shanghai cookies

www.bsc.es



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Paraver

PATC Parallel Programming Workshop
Oct 14-18 2013

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**

⌋ Sequence of time stamped records

- Punctual events
 - Something happened: when and where (object/entity: ... thread)
 - One record per specific information (encoded as a type and a value)
 - About the event
 - About the interval from the previous event till this one (i.e. hardware counts,...)
- Relations between objects (... communications)
 - Source and sink
 - Attributes (... size, tag)
- Separate numeric (.prv) and symbolic (.pcf) files

⌋ Only information derived from captured events and data can be reported.

- Trivial but ... often forgotten

« Major BSC instrumentation package

« When / where

- Parallel programming model runtime
 - MPI, OpenMP, pthreads, OmpSs, CUDA, OpenCL, MIC...
 - API entry/exit, OpenMP outlined routines
- Selected user functions
- Periodic samples
- User events

« Additional information

- Counters
 - PAPI
 - Network counters
 - OS counters
- Link to source code
 - Callstack

How does Extrae intercept your app?

« LD_PRELOAD

- Works on production binaries
- Specific library for each combination of runtimes
- Does not require knowledge on the application

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

¹ for Fortran codes

« Dynamic instrumentation

- Works on production binaries
- Just specify functions to be instrumented.

Based on DynInst
U.Wisconsin/U.Maryland

« Other possibilities

- Link instrumentation library statically (i.e., PMPI @ BG/Q, ...)
- OmpSs (instrumentation calls injected by compiler + linked to library)

How to use Extrae?

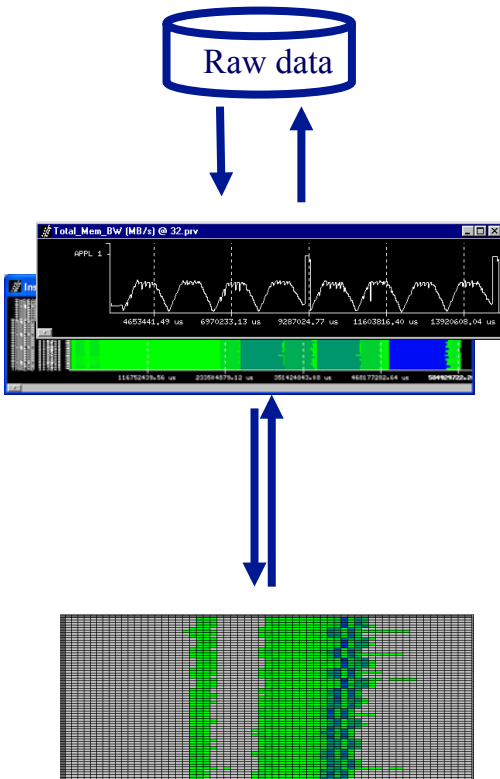
- « Adapt job submission script
 - Specify LD_PRELOAD library and xml instrumentation control file
- « Specify the data to be captured in the .xml instrumentation control file
- « Run and get the trace ...

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

Paraver – Performance data browser



Trace visualization/analysis

+ trace manipulation

Timelines

Goal = Flexibility

No semantics

Programmable

**2/3D tables
(Statistics)**

Comparative analyses

Multiple traces

Synchronize scales

Paraver mathematical foundation

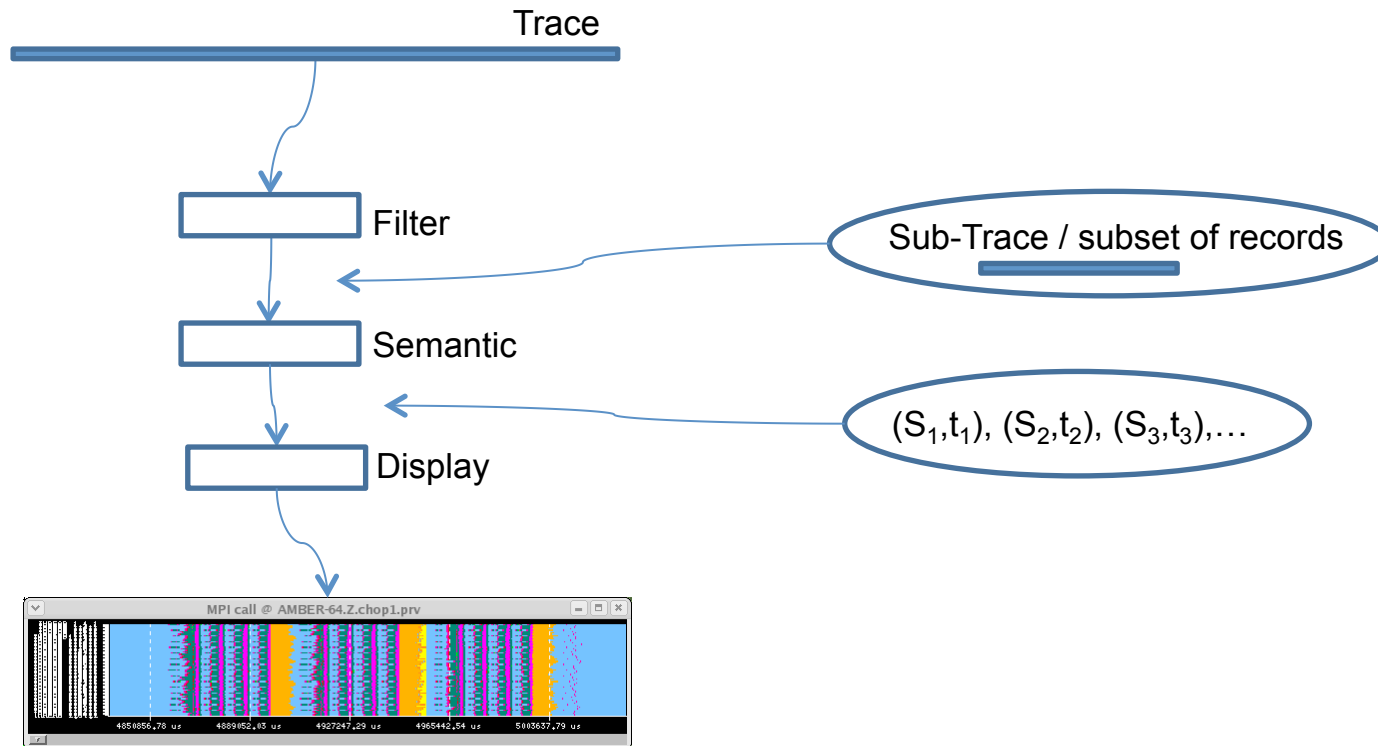
- ⌘ Every behavioral aspect/metric described as a function of time
 - Possibly aggregated along
 - the process model dimension (**thread**, process, application, workload)
 - The resource model dimension (core, node, system)
 - Language to describe how to compute such functions of time (GUI)
 - Basic operators (from) trace records
 - Ways of combining them
- ⌘ Those functions of time can be rendered into a 2D image
 - Timeline
- ⌘ Statistics can be computed for each possible value or range of values of that function of time
 - Tables: profiles and histograms

Parallel mathematical foundation

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

Function of time

Series of values



⌘ 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

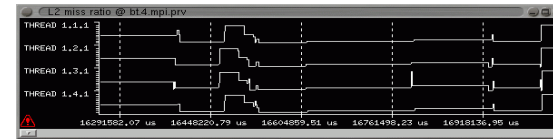
- IPC, L2 miss ratio, Duration of MPI call, duration of computation burst

$$S_i \in R$$

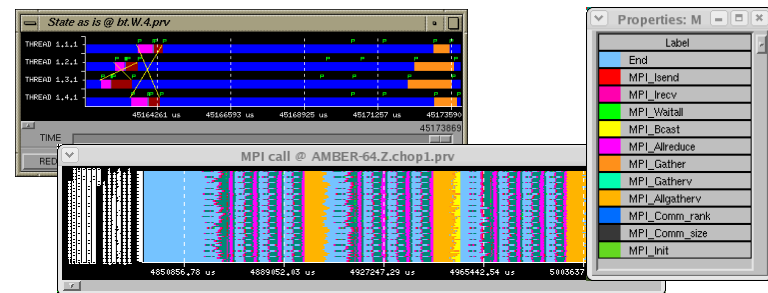
Timelines

Representation

– Function of time

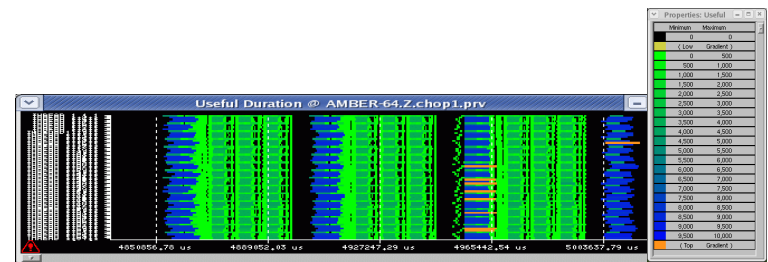


– Colour encoding



– Not null gradient

- Black for zero value
- Light green → Dark blue



Non linear rendering to address scalability

J. Labarta, et al.: “Scalability of tracing and visualization tools”, PARCO 2005

Basic functions of time

Semantic module

From Events to functions of time

- Last event value $s(t) = V(EV_i)$
- Next event value $S(t) = V(EV_{i+1})$
- Average Next Event Value $s(t) = \frac{V(EV_{i+1})}{T(EV_{i+1}) - T(EV_i)}$
- Interval btw. Events $S(t) = T(EV_{i+1}) - T(EV_i)$



Semantic module

From communication records to functions of time

- Send Bytes $s(t) = \sum_j Sz(C_j) \cdot j \mid (T_s(C_j) < t) \wedge (T_r(C_j) > t) \wedge (Dir(C_j) == send)$
- Send Bandwidth $s(t) = \sum_j \frac{Sz(C_j)}{T_s(C_j) - T_r(C_j)} \cdot j \mid (T_s(C_j) < t) \wedge (T_r(C_j) > t) \wedge (Dir(C_j) == send)$
- Msgs in transit $s(t) = \sum_j sign(j) \cdot j \mid (T_s(C_j) < t) \wedge (T_r(C_j) > t) \wedge (Dir(C_j) == send)$
- Recv. Bandwidth $s(t) = \sum_j \frac{Sz(C_j)}{T_s(C_j) - T_r(C_j)} \cdot j \mid (T_s(C_j) < t) \wedge (T_r(C_j) > t) \wedge (Dir(C_j) == recv)$
- Rec. Negative Msgs $s(t) = \sum_j sign(j) \cdot j \mid (T_s(C_j) < t) \wedge (T_r(C_j) > t) \wedge (Dir(C_j) == recv)$
- Comm. Partner $s(t) = Partner(C_j) \cdot j \mid (T_s(C_j) < t) \wedge (T_r(C_j) > t)$
- Bytes btw. Events $S(t) = \sum_j Sz(C_j) \cdot j \mid T_s(C_j) \in [T(EV_i), T(EV_{i+1})) \vee T_r(C_j) \in [T(EV_i), T(EV_{i+1}))]$



Composition

S'(t) = f(S(t))

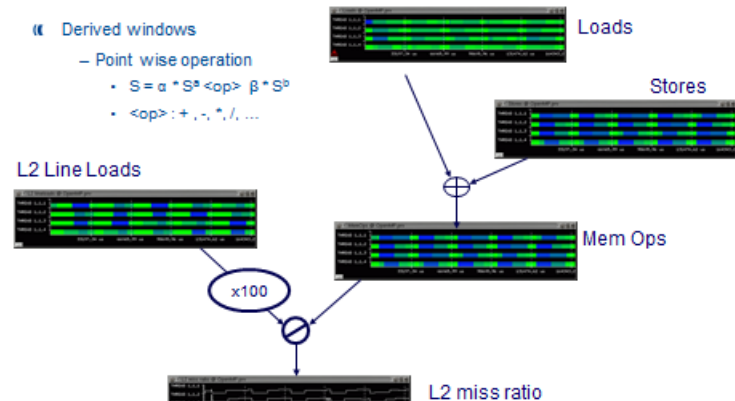
$$S' = f \circ S$$

- 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 * 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
 - $S = \alpha * S^a <op> \beta * S^b$
 - $<op> : +, -, *, /, \dots$

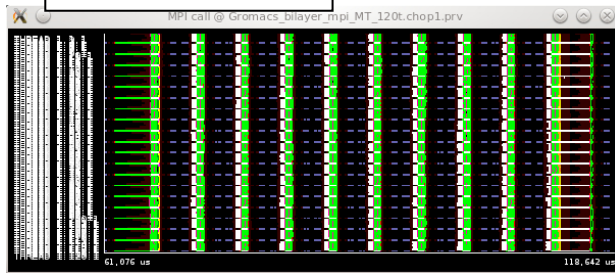


See slides at end of presentation for details

Tables: Profiles, histograms, correlations

« From timelines to tables

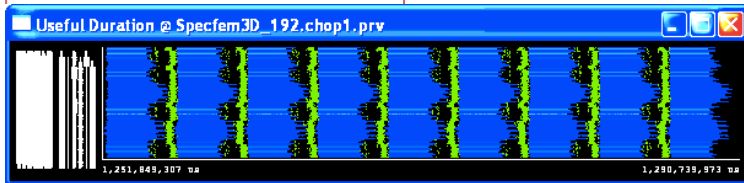
MPI calls



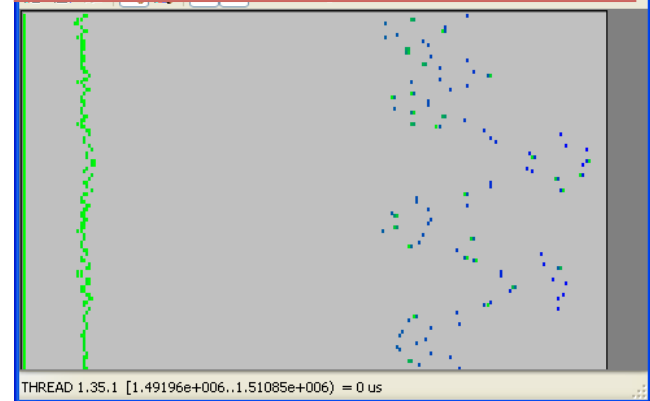
MPI calls profile

	Outside MPI	MPI Send	MPI Recv	MPI Isend	MPI Irecv	MPI Waitall	MPI Bcast	MPI Reduce	MPI Allr
THREAD 1.113.1	67.6081 %	0.0592 %	9.9182 %	2.5777 %	1.7699 %	5.1676 %	0.5934 %	0.1465 %	
THREAD 1.114.1	42.8434 %	-	20.5621 %	1.1947 %	1.0400 %	7.7056 %	-	-	
THREAD 1.115.1	68.6127 %	0.0707 %	9.6223 %	2.2589 %	2.0177 %	5.9825 %	0.5249 %	0.0297 %	
THREAD 1.116.1	74.6039 %	0.0531 %	9.6084 %	2.8813 %	2.5593 %	2.9286 %	0.5095 %	0.0483 %	
THREAD 1.117.1	74.3733 %	0.0691 %	9.7012 %	2.8517 %	2.5240 %	-	-	-	
THREAD 1.118.1	72.7770 %	0.0545 %	9.5489 %	2.8489 %	2.5353 %	-	-	-	
THREAD 1.119.1	66.7994 %	0.0682 %	10.0674 %	2.4206 %	1.9741 %	-	-	-	
THREAD 1.120.1	43.7224 %	-	20.5273 %	1.1912 %	1.0175 %	-	-	-	
Total	8.012.4546 %	7.3174 %	1,370.5276 %	288.6168 %	253.0137 %	54			
Average	66.7705 %	0.0690 %	11.4211 %	2.4051 %	2.1084 %				
Maximum	75.6821 %	0.4390 %	21.2505 %	2.9706 %	2.6369 %				
Minimum	40.5200 %	0.0129 %	8.8583 %	1.1489 %	1.0077 %				
StDev	11.3685 %	0.0474 %	4.0613 %	0.5984 %	0.5406 %				
Avg/Max	0.8822	0.1572	0.5374	0.8096	0.7996				

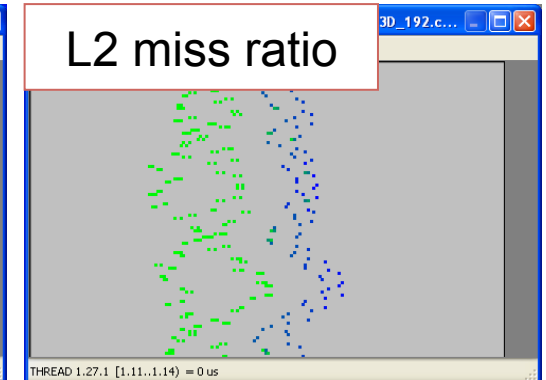
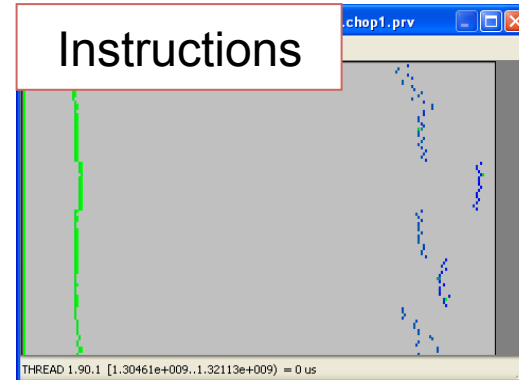
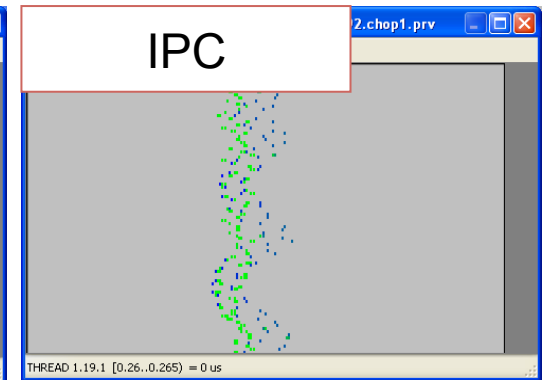
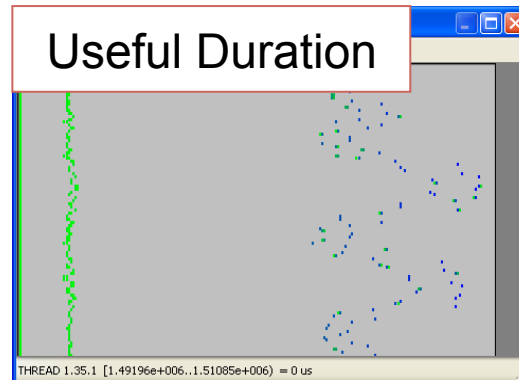
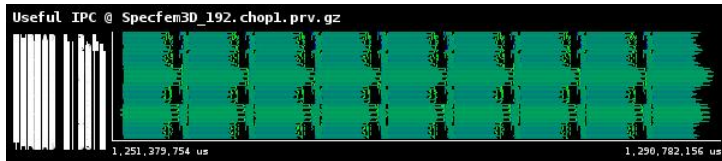
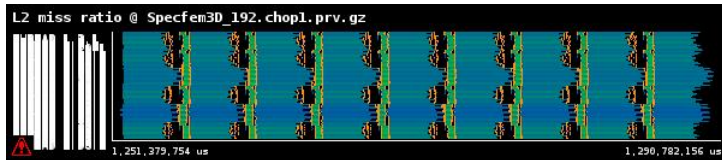
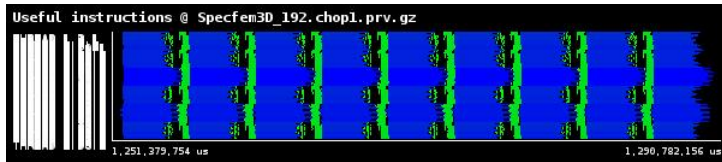
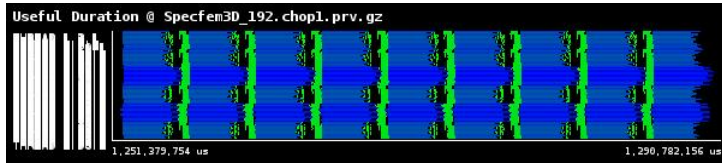
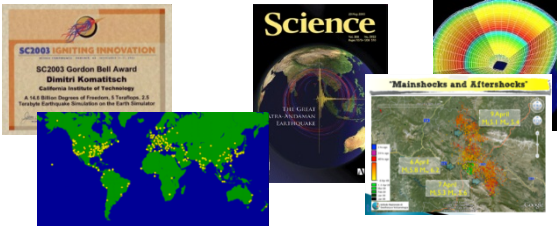
Useful Duration



Histogram Useful Duration

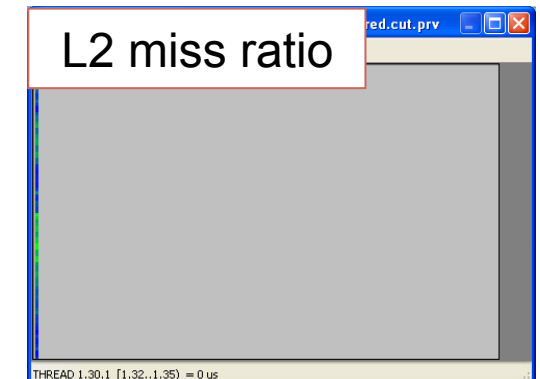
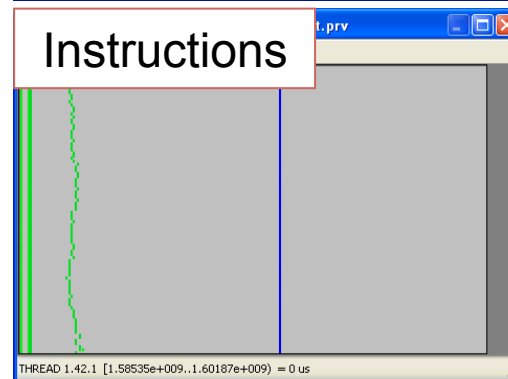
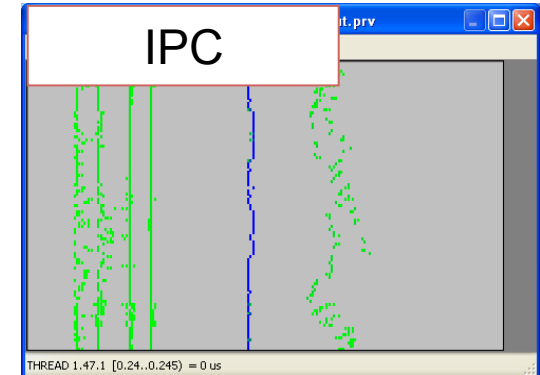
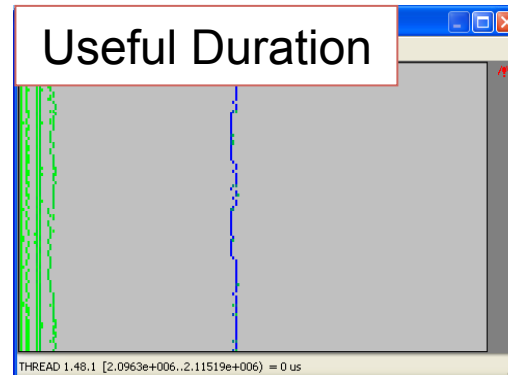


Analyzing variability through histograms and timelines



Analyzing variability through histograms and timelines

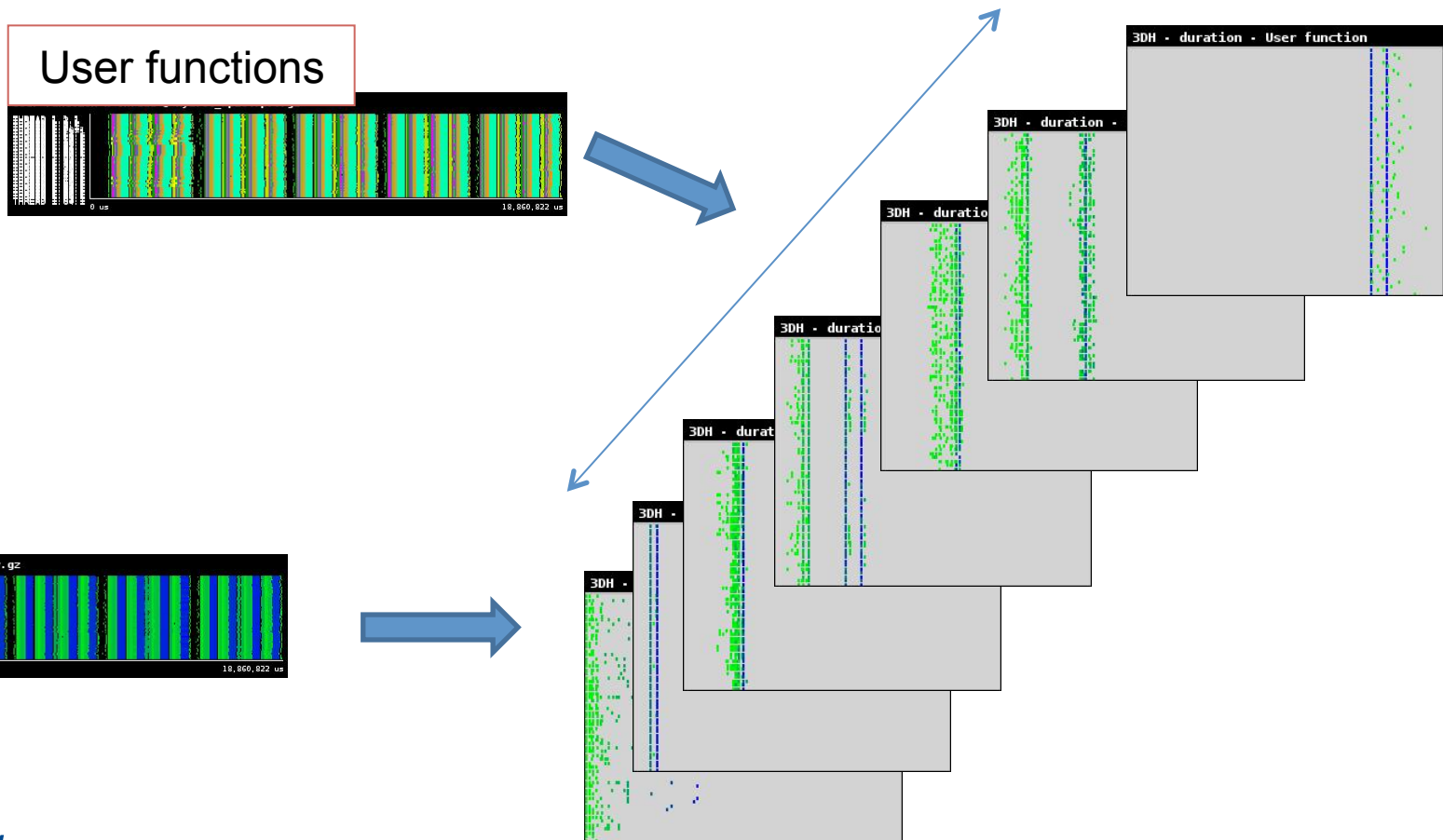
☞ By the way: six months later



3D Tables

« An additional control dimension

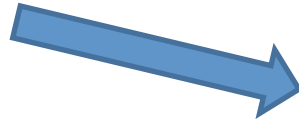
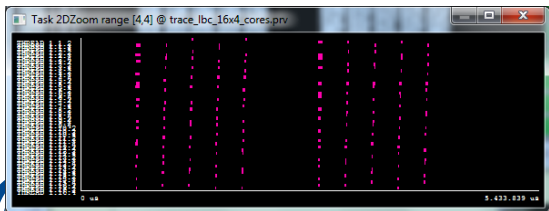
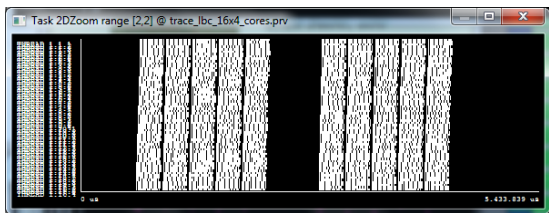
- One table (plane) per value (or range) of 3D window
- i.e. histogram of duration of each function



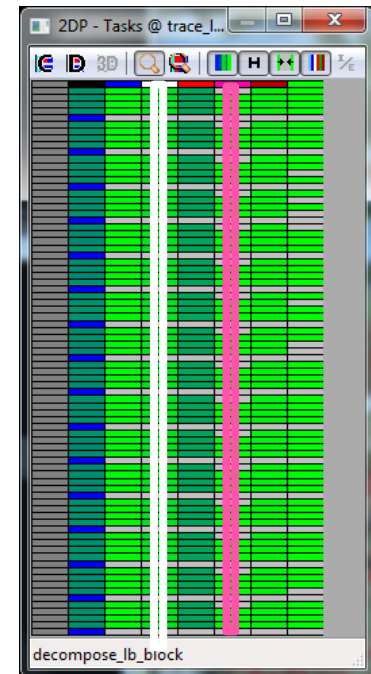
From tables to timelines

« Where in the timeline do the values in certain table columns appear?

– ie. want to see the time distribution of a given routine?

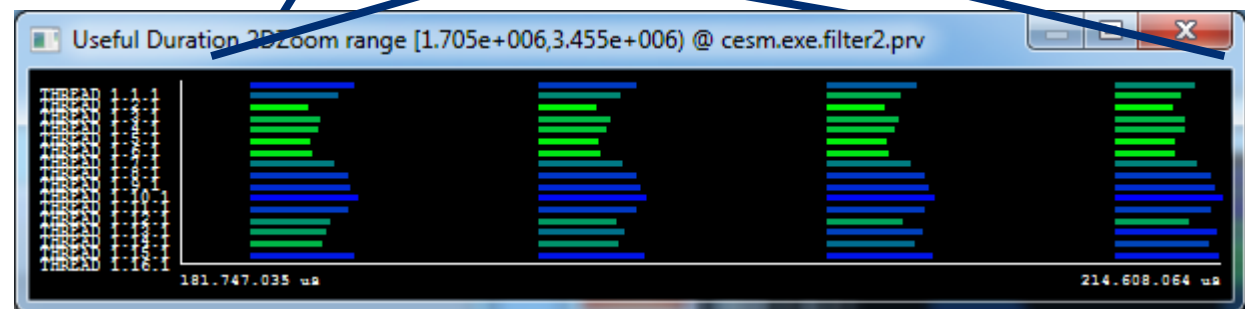
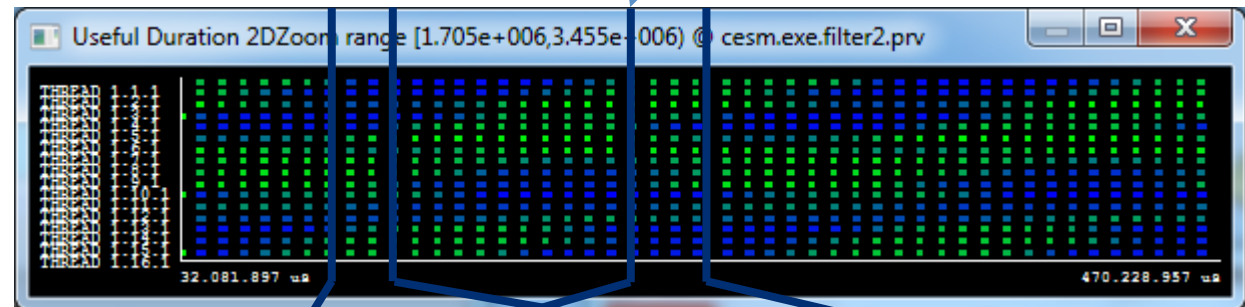
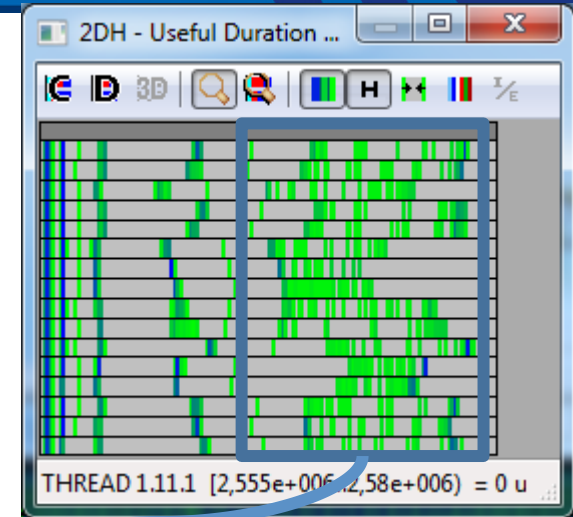


Only showing when a given value happens



Variability ... is everywhere

- ⌘ CISM: 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 ? ☺



Other mechanisms integrated in the GUI

« Trace manipulation

- Cut
- Filter

« Performance analytics

- Clustering
- Folding
- Tracking

« Executing external commands and tools

- BSC Tools
- Scripts
- External tools

Trace manipulation

« Data handling/summarization capability

– Filtering

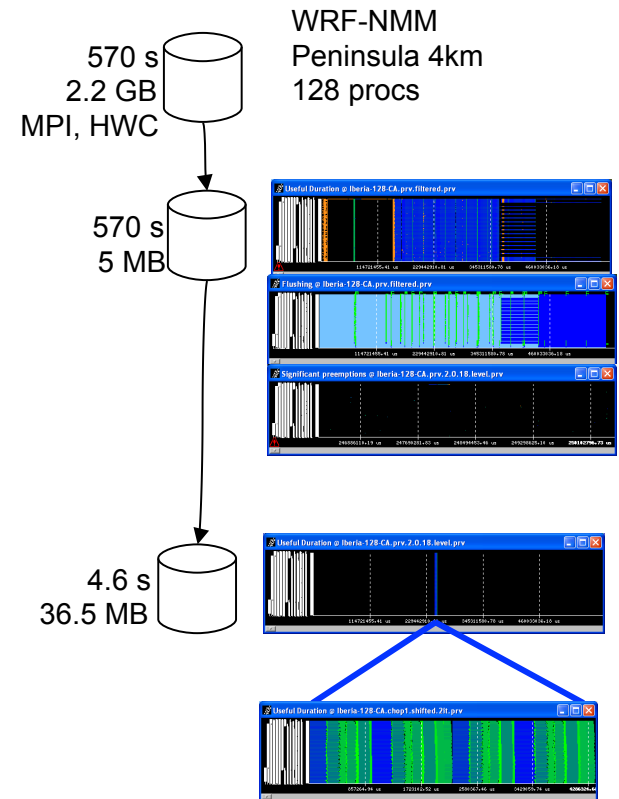
- Subset of records in original trace
- By duration, type, value,...
- Filtered trace IS a paraver trace and can be analysed with the same cfgs (as long as needed data kept)

– 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

External commands and tools

⌘ Execute external commands ...

- Predefined: Dimemas, Stats,...
- User specified binaries or scripts

⌘ ... specifying arguments

- Trace
- Command specific arguments

⌘ Paramedir

- Non graphical version of Paraver
- Reads trace, applies standard cfgs, writes ASCII output (table,...)

⌘ Scripts can use paramedir, Dimemas, clustering ... in parametric sweeps, search/optimization loops,...

Executing external commands

« Example: basic_analysis.py

```
=====
Analysis of XXXXX.prv
=====
Timing:
  Elapsed duration      = 3.887 s
  Ideal time            = 3.223 s
  Compute time         = 2.637952 s
  MPI time              = 1.24930663 s

Parallel Efficiency:
  Total Efficiency      = 0.679
  Load Balance         = 0.950
  Micro Load Balance  = 0.859
  Transfer              = 0.829
  Bweff                = 0.814
  Leff                 = 1.000

Load balance:
  Time Load Balance    = 0.950
  Instructions Load Balance = 0.990
  Cycles load balance  = 0.950
  IPC Load Balance     = 0.960

Computational analysis:
  Total useful instructions = 5.399636e+11
  Average useful instructions per process = 8.436931e+09
  Instructions based microload balance/sync = 0.924

sequential performance:
  Average MIPS         = 3199
  Average IPC          = 1.050
```

Basic communication statistics:

Point to point:

Average number of calls = 554.0
Balance in number of calls = 1.0
Average bytes per process = 728596480.0
Balance in bytes = 0.990

Collectives:

Num collective calls = 6.0
Average bytes per call = 6.670
Max bytes per call = 8.000
Balance in bytes = 1.000

Inter - Intra node communication statistics:

Bytes sent (MB):

Locally = 22229.811
Remotely = 24645.730

Number of sends:

Locally = 3520.000
Remotely = 6592.000

www.bsc.es



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Performance Analytics

PATC Parallel Programming Workshop
Oct 14-18 2013

Performance Analytics

« Dominant practice

- We focus a lot on capturing a lot of data
- but we present either everything or first order statistics
- and require new experiments without squeezing the potential information from the previous one

« Need for performance analytics

- Leveraging techniques from data analytics, mining, signal processing, life sciences,...
- towards insight
- And models

« Some techniques worked on at BSC

- Spectral analysis
- Clustering
- Folding
- Simulation (Dimemas)

Spectral analysis

Repetitive behavior

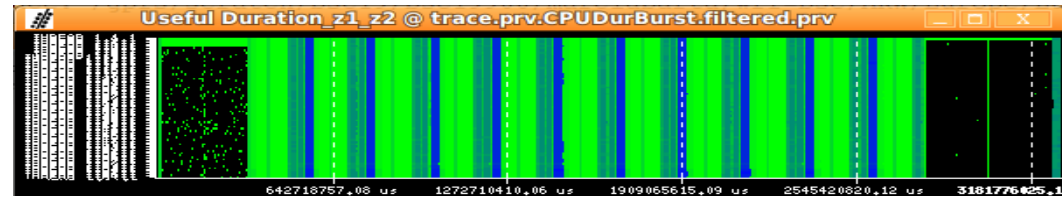
- ⌘ Applications tend to have Iterative behavior
 - Detailed analysis can be applied to a few such iterations

- ⌘ Metrics in Paraver are functions of time
 - Natural target for signal processing techniques to automatically detect such iterative structure
 - Relevant functions of time at global application level
 - # processes in MPI, outside MPI, ...
 - Sum of useful burst duration
 - Semantic: high when many processes are in the middle of very long computation bursts
 - Does capture repetitive structure of application

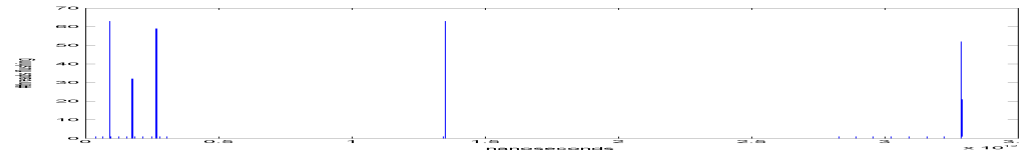
Signal processing applied to performance analysis

Techniques

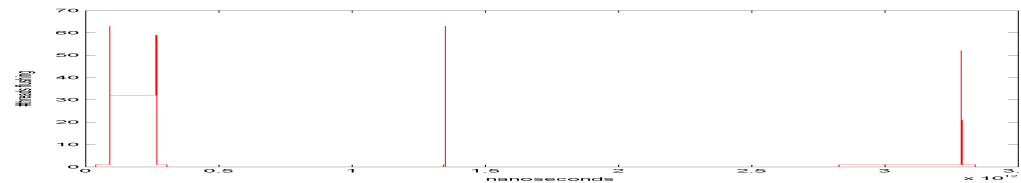
- Mathematical morphology
 - clean up perturbed regions
- Wavelet transform
 - identify coarse regions
- Spectral analysis
 - detailed periodic pattern



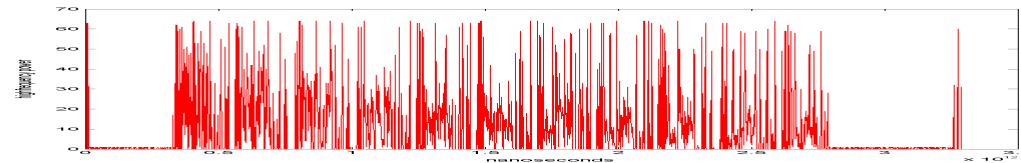
Flushing



Flushing filtered



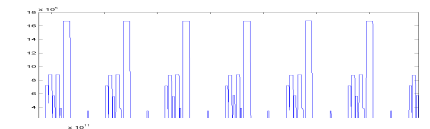
Wavelet High frequency



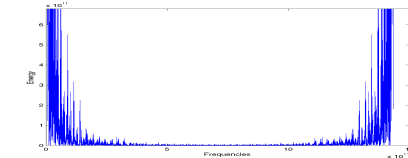
Useful

- Identify structure (periodicity)
- Reduce trace sizes
- Increase precision of profiles (report non perturbed stats)

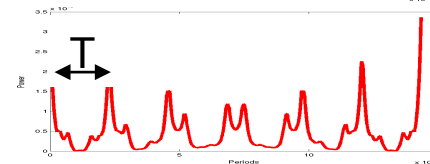
Σ Useful Duration



Spectral density

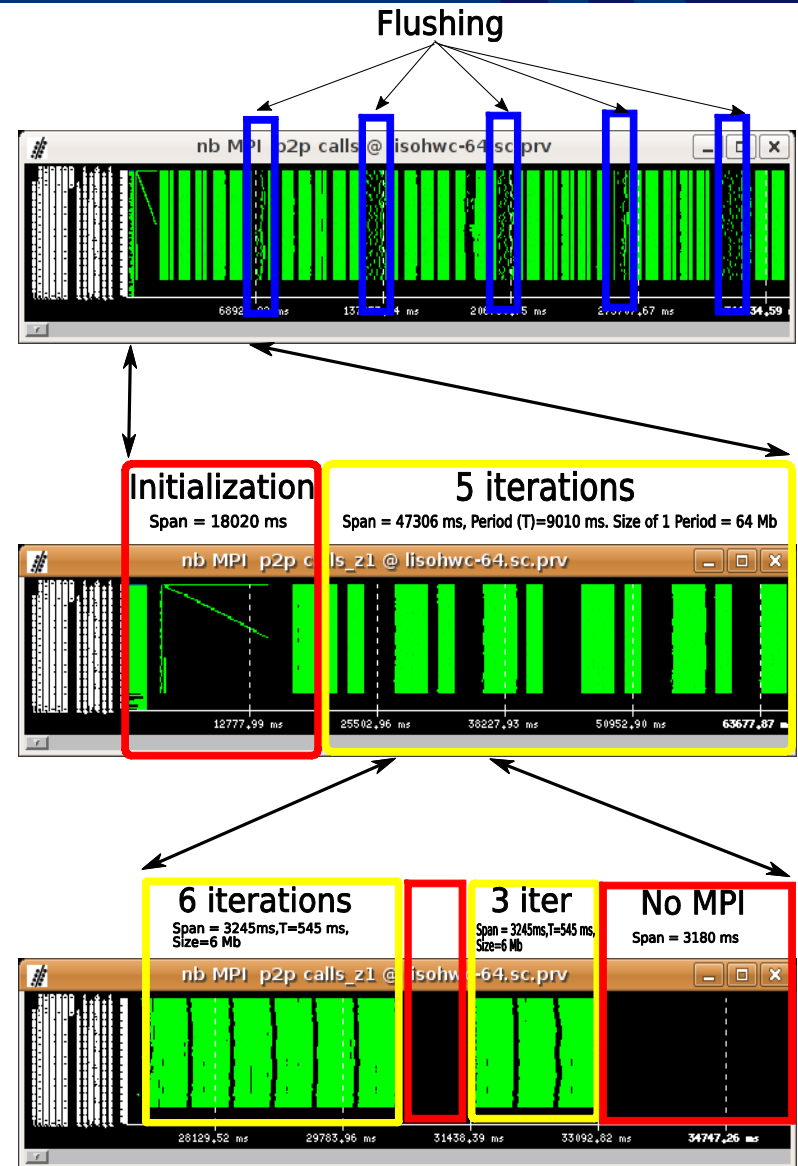


Autocorrelation

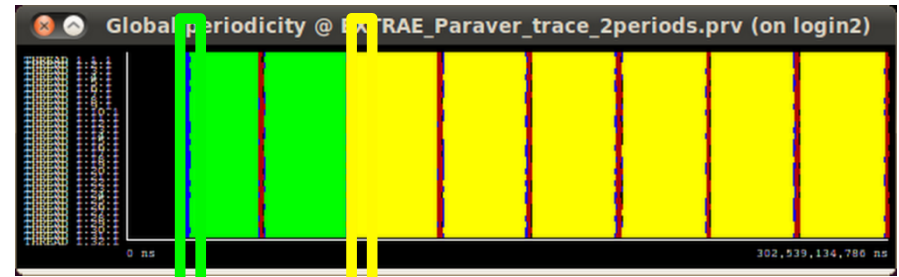
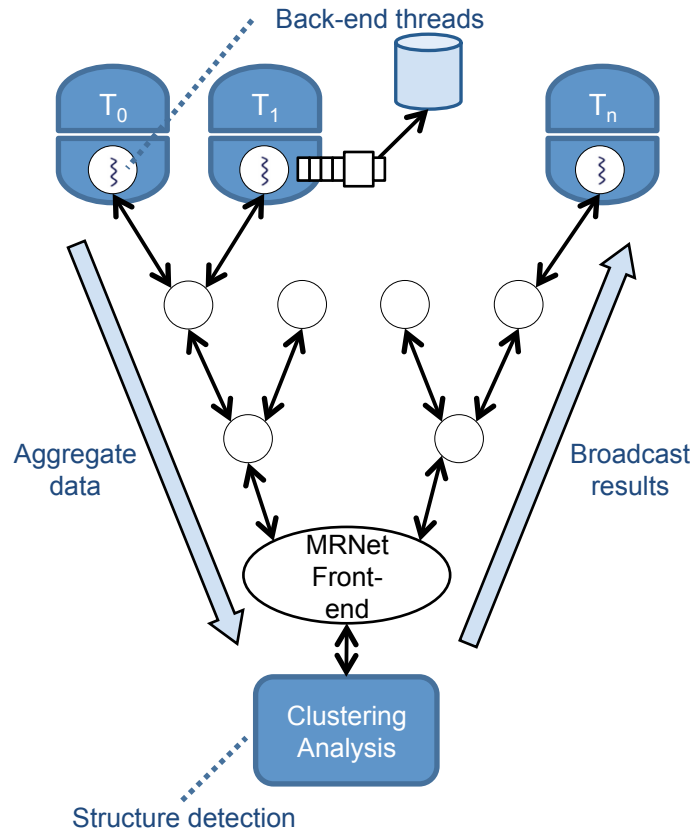


Signal processing applied to performance analysis

« Hierarchical structure identification



Scalability: online automatic interval selection



Detailed trace for only small interval

“ G. Llort et all, “Scalable tracing with dynamic levels of detail” ICPADS 2011

Clustering

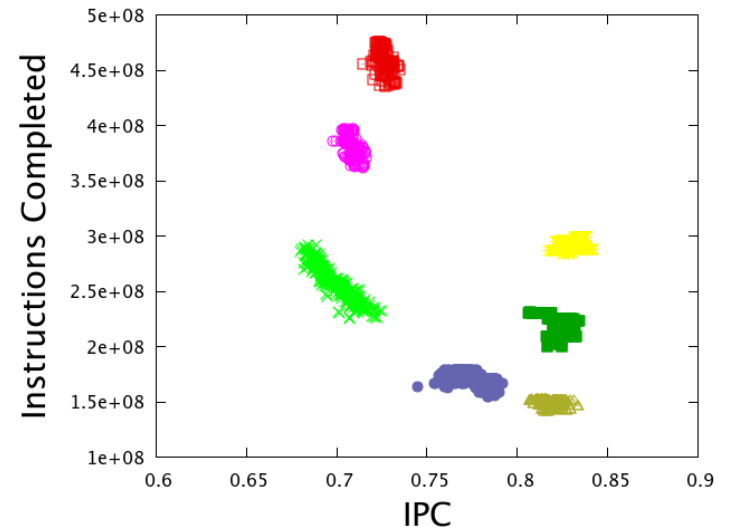
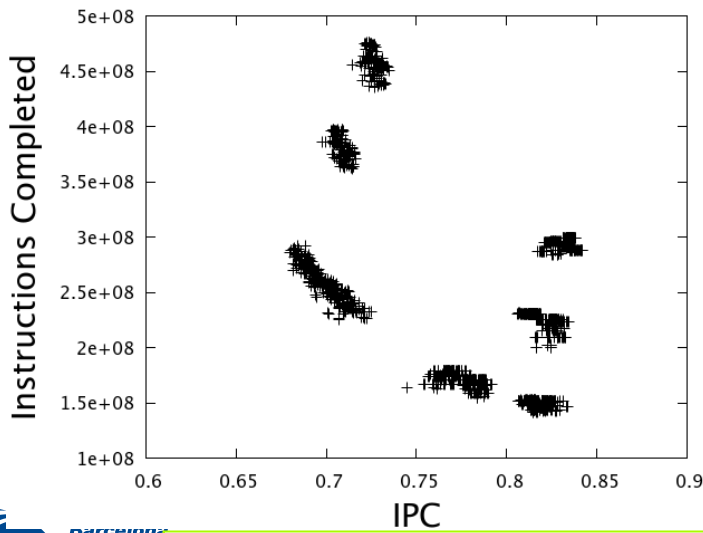
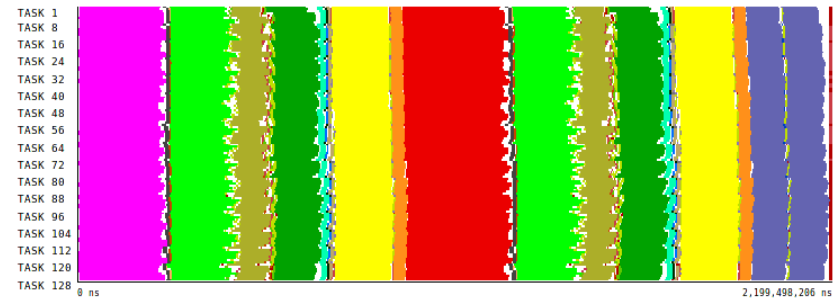
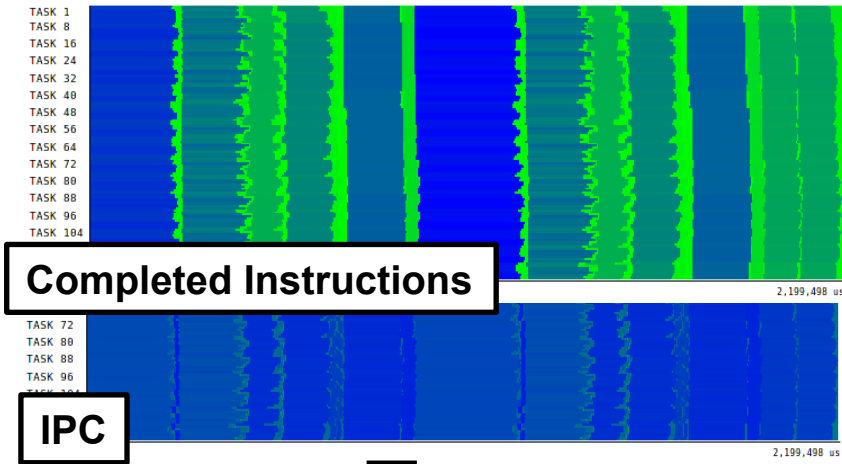
« Identification of computation structure

- CPU burst = region between consecutive runtime (MPI, OpenMP) calls
 - Described with performance hardware counters
 - Associated with call stack data

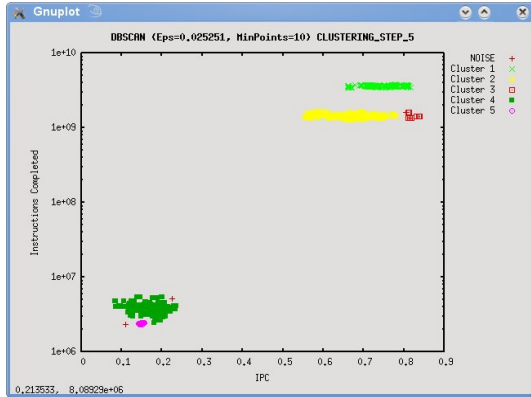
« Scatter plot on some relevant metrics

- Instructions: idea of computational complexity, computational load imbalance,...
- IPC: Idea of absolute performance and performance imbalance
- Automatically Identify clusters

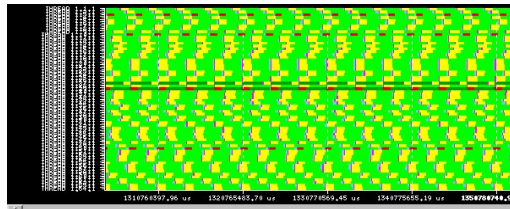
Using Clustering to identify structure



Performance @ serial computation bursts

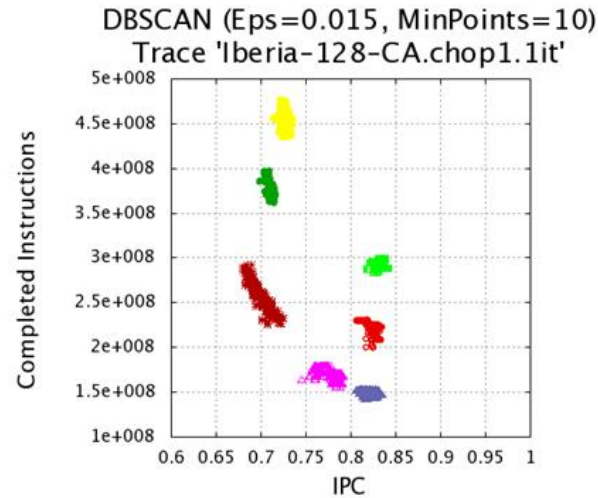


SPECFEM3D

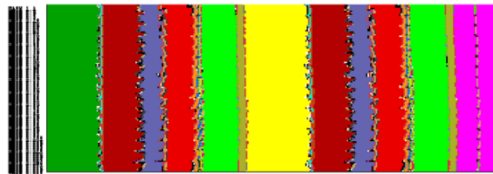


Asynchronous SPMD

Balanced #instr
variability in IPC



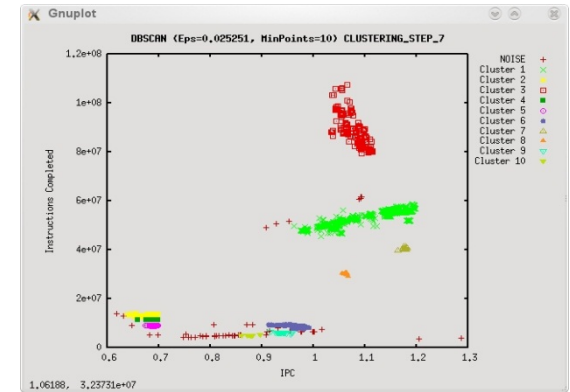
WRF 128 cores



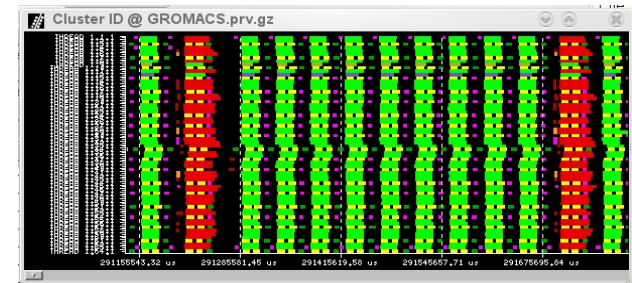
SPMD

Repeated substructure

Coupled imbalance



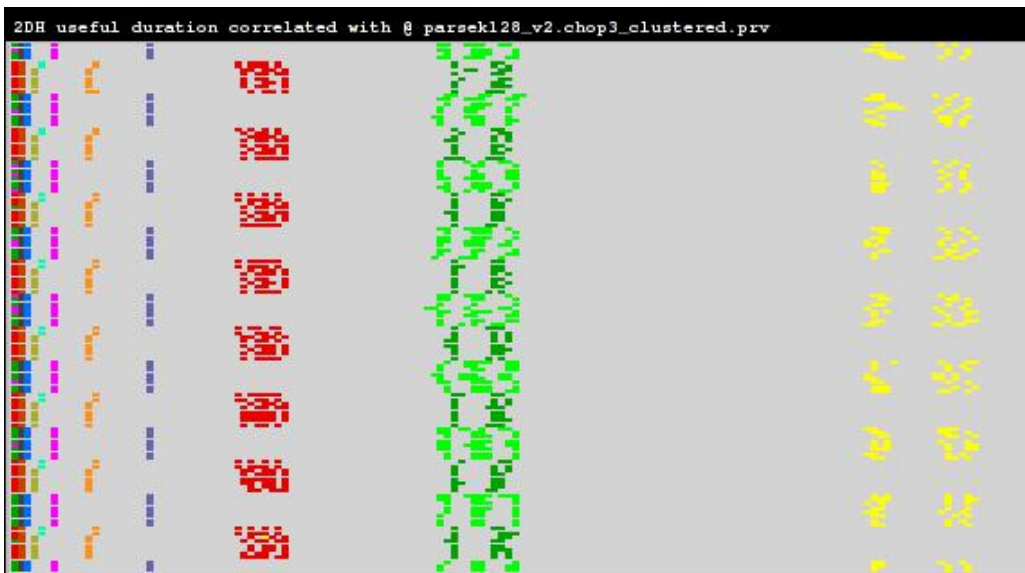
GROMACS



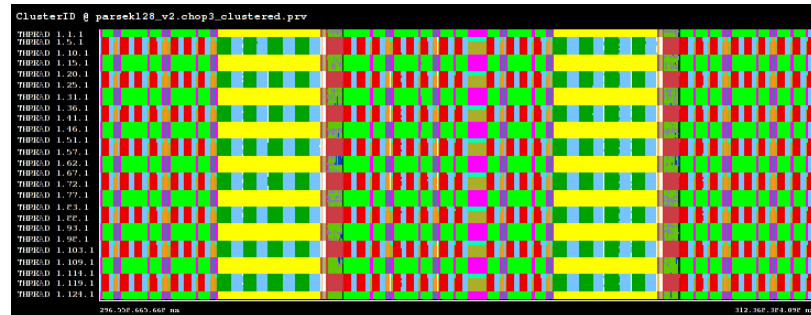
MPMD structure

Different coupled
imbalance trends

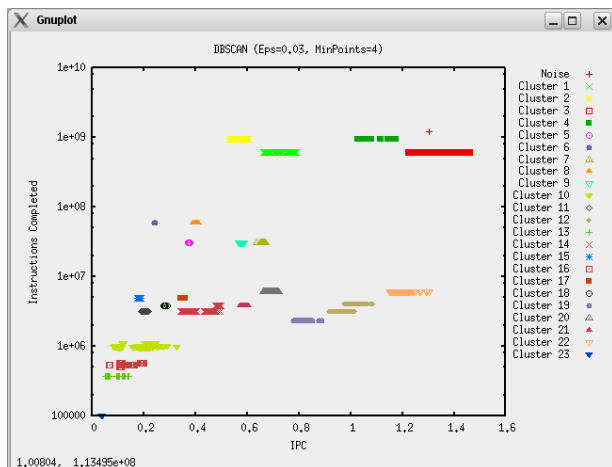
Example PARSEK (DEEP)



duration vs. cluster

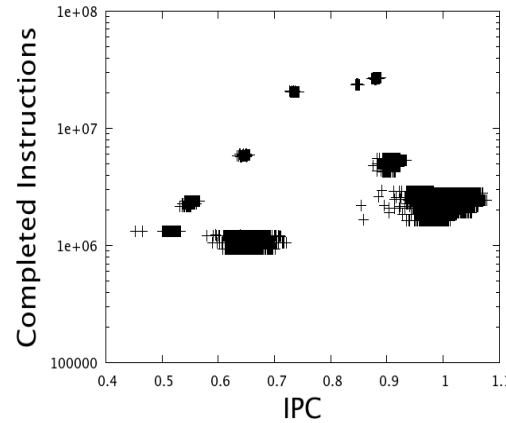
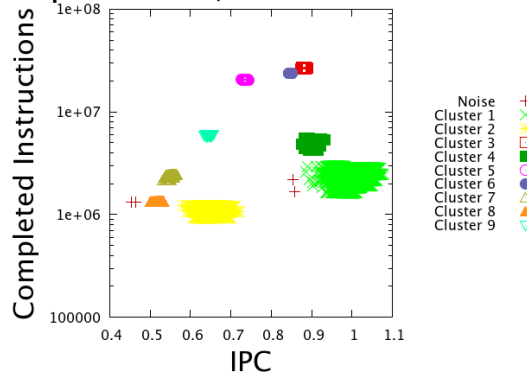


instr. vs. cluster

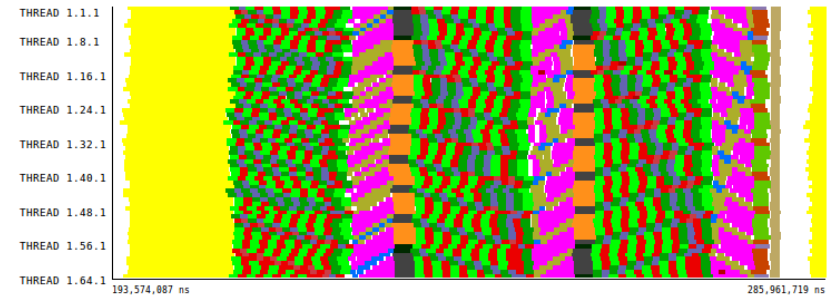
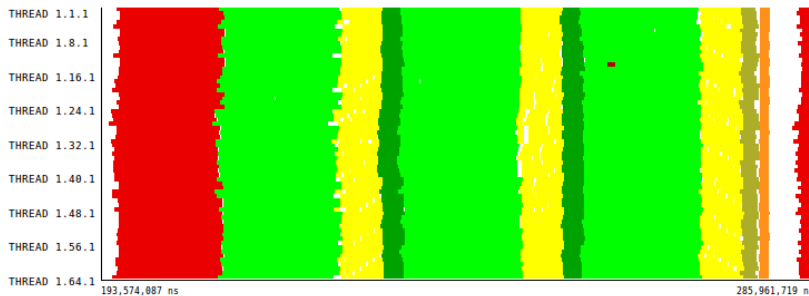
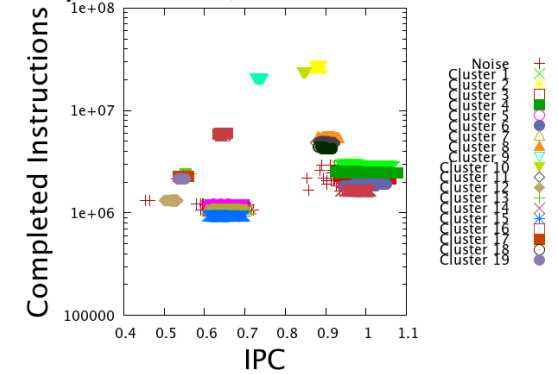


Structure quality?

Eps=0.040, MinPoints=10



Eps=0.014, MinPoints=10



⌋ How many clusters?

⌋ Which is better?

- The two describe interesting structure
- Typically SPMD would be a good first level of description for most apps

« Clustering enables focusing the analysis and opens many different uses

- Analysis
 - Detection of application structure
- Precise instantaneous metrics
 - correlation of sampled data to generate instantaneous metric evolution
- Dimemas:
 - Separate speed factors per cluster on predictive simulations
- Track the evolution of application behaviour effects
- ...

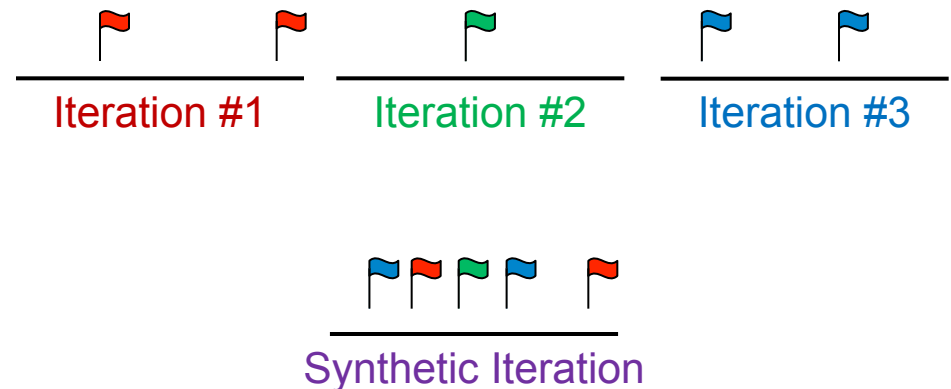
Folding

Mixing instrumentation and sampling ...

« ... to get extreme detail with minimal overhead

« Different roles

- Instrumentation delimits regions
- Sampling report progress within region

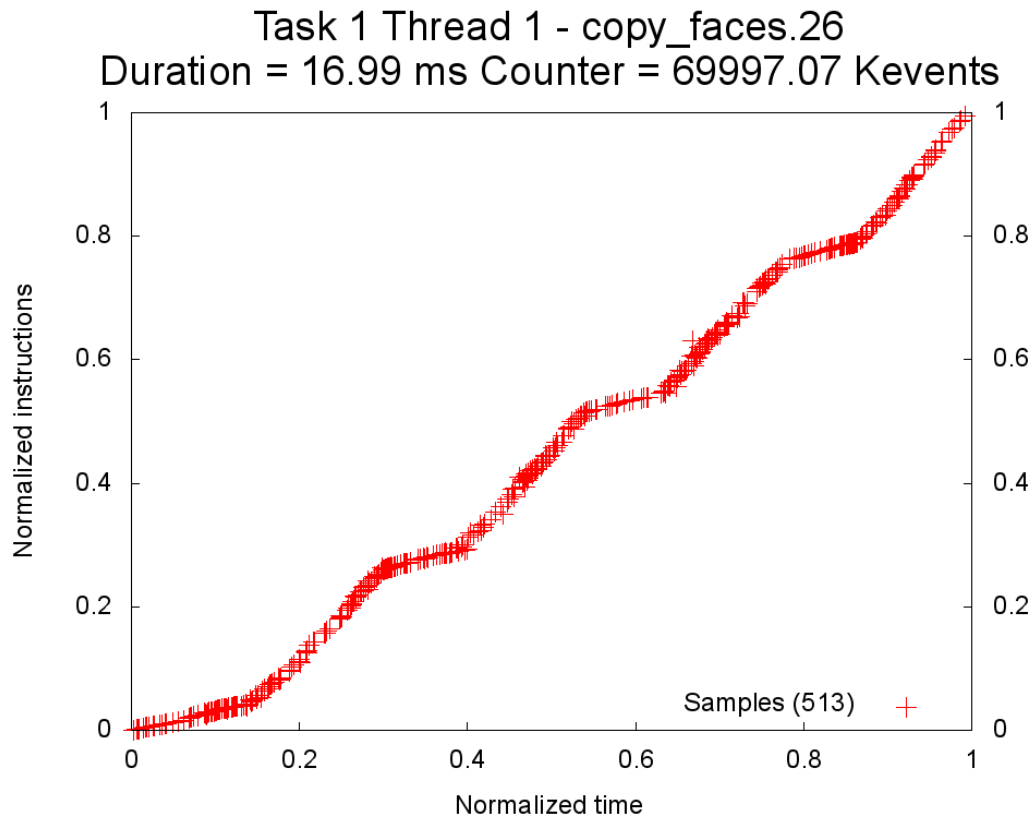


Harald Servat et al. “Detailed performance analysis using coarse grain sampling” PROPER@EUROPAR, 2009

Harald Servat et al. “Unveiling Internal Evolution of Parallel Application Computation Phases” ICPP 2011

Folding hardware counters

- Instructions evolution for routine `copy_faces` of NAS MPI BT.B

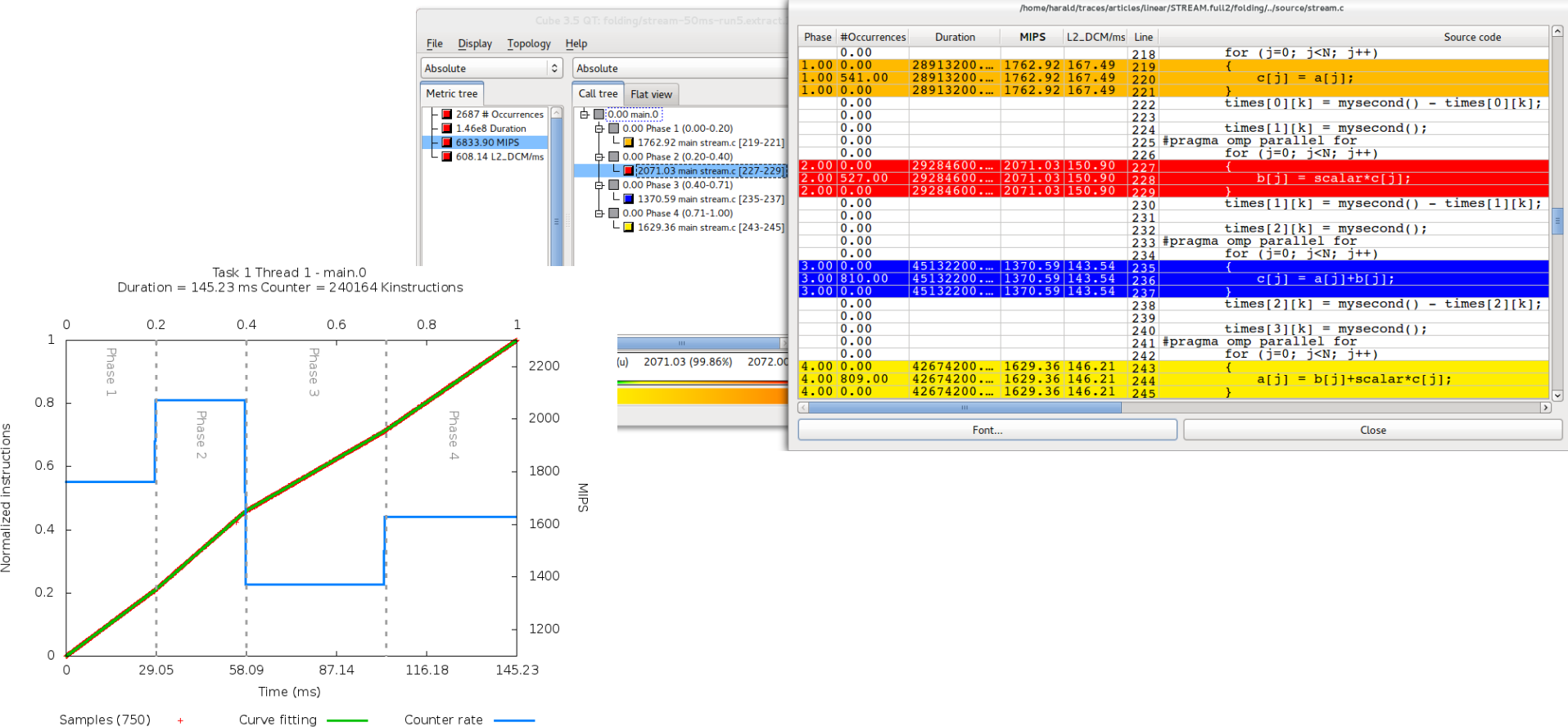


- Red crosses represent the folded samples and show the completed instructions from the start of the routine
- Green line is the curve fitting of the folded samples and is used to reintroduce the values into the tracefile
- Blue line is the derivative of the curve fitting over time (*counter rate*)

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

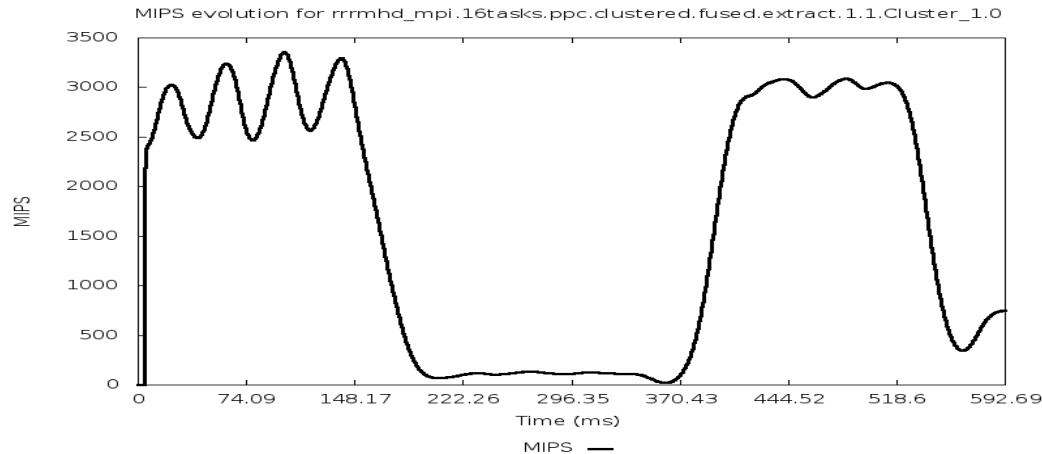


What is a good performance?

Performance of a sequential region = 2000 MIPS

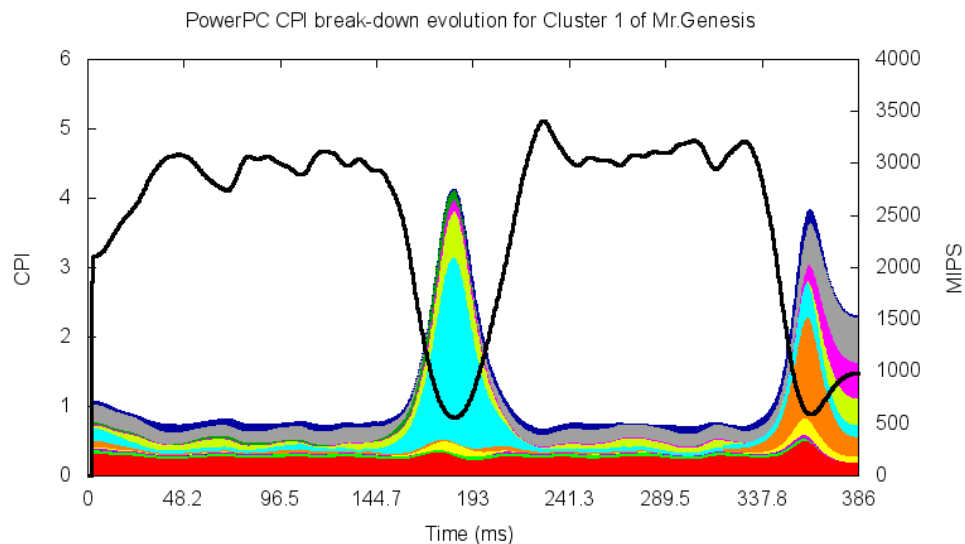
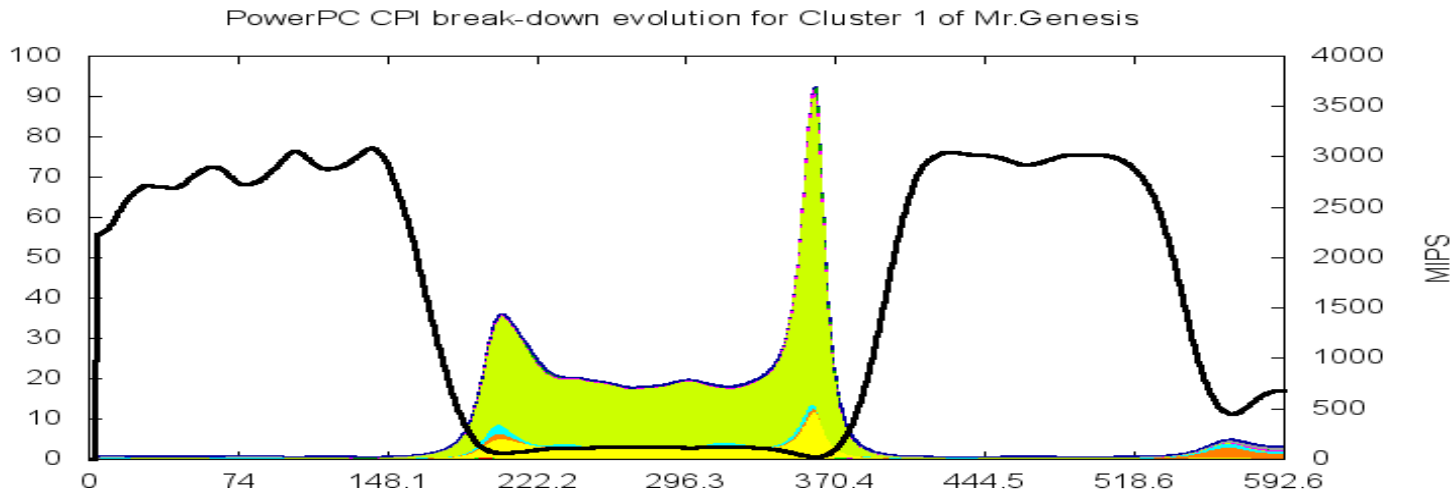
Is it good enough?

Is it easy to improve?



Instantaneous CPI stack

MRGENESIS



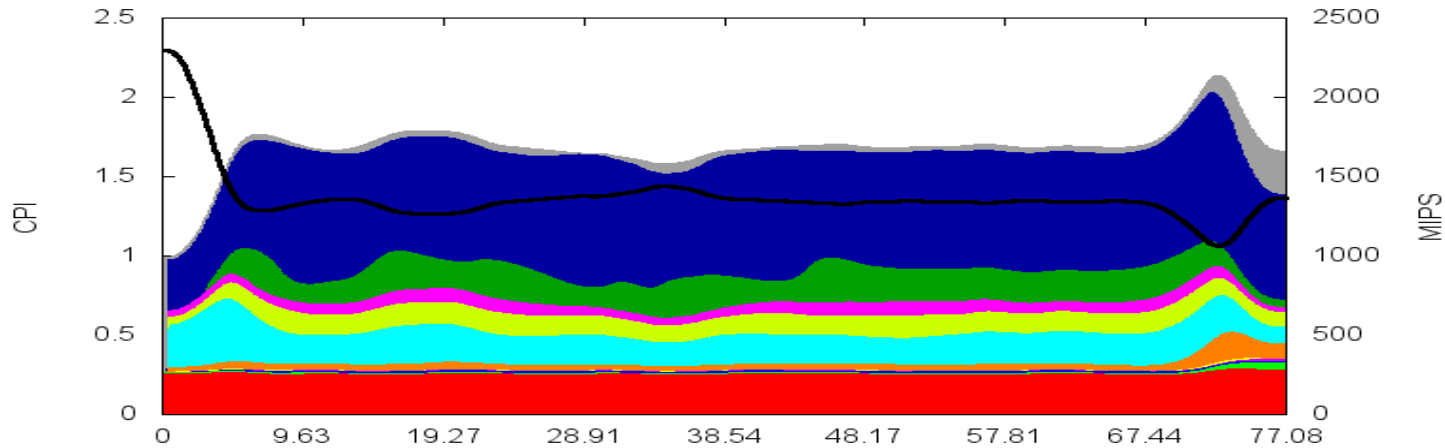
- | | |
|-------------------------|----------------------|
| Useful cycles | LSU: Basic latency |
| I-cache miss | FXU: Div/MSTPR/MSFPR |
| Branch mispredict | FXU: Basic latency |
| Flush penalties, etc | FPU: FDIV/FSQRT |
| LSU: Translation lookup | FPU: Basic latency |
| LSU: Other reject | Other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

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



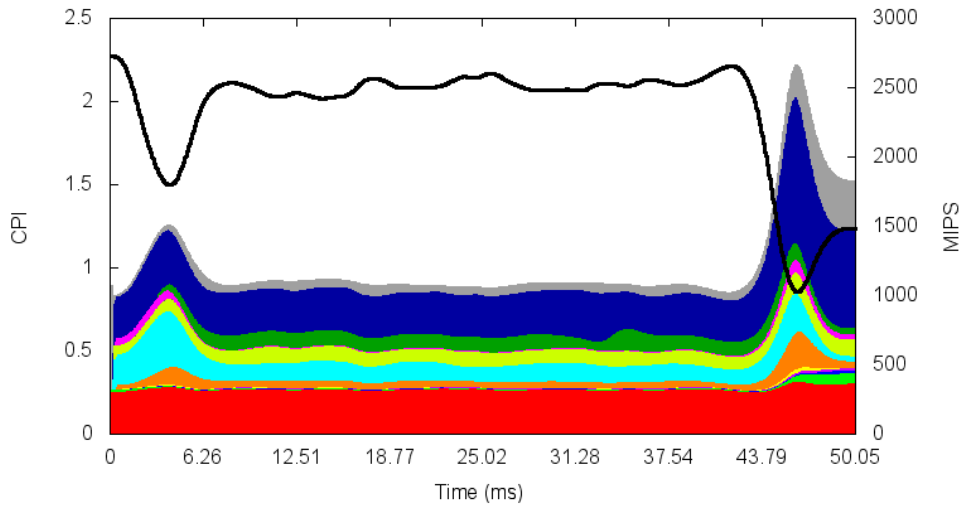
Instantaneous CPI stack

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



PMEMD

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



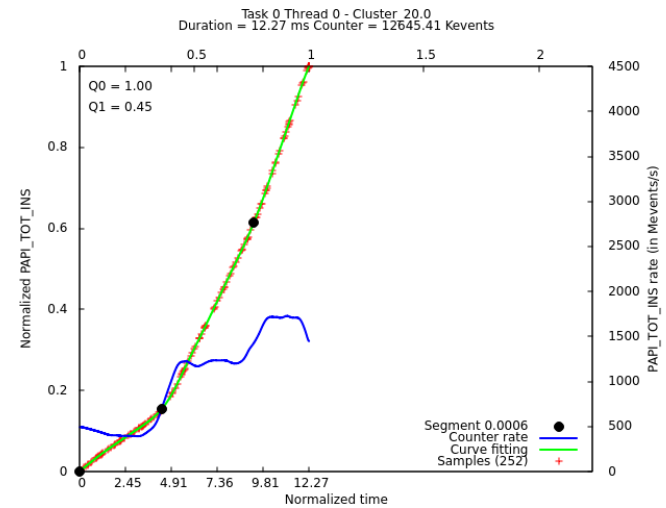
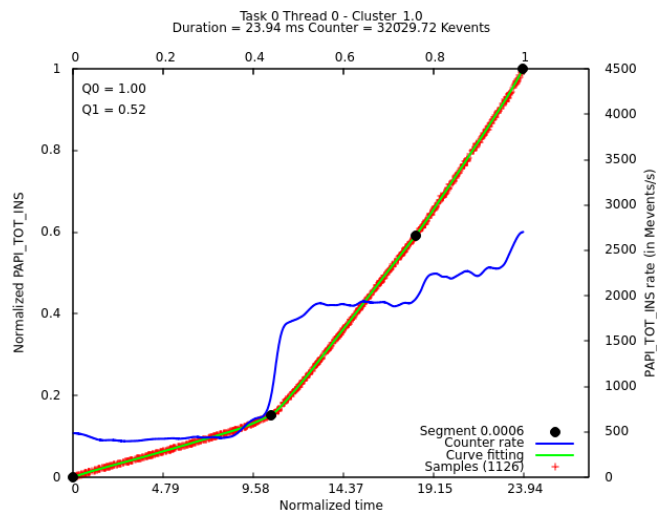
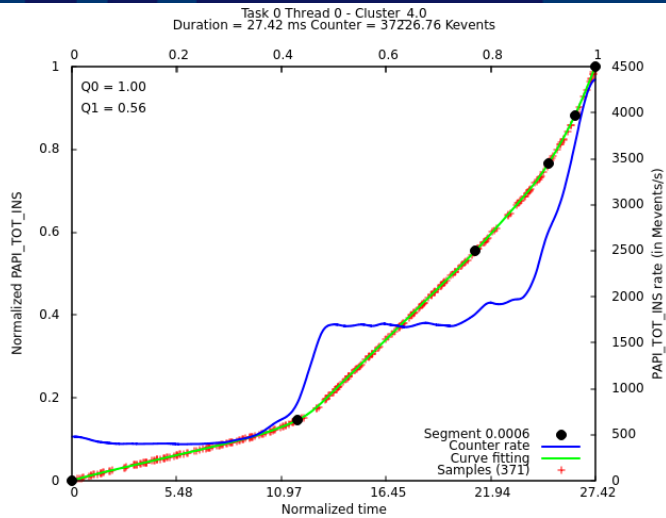
- Fix: Precompute transcendental functions. Use stored value.
- Serialization of computation
- vs memory access
 - Interesting tradeoff



Correlating counters

CG-POP

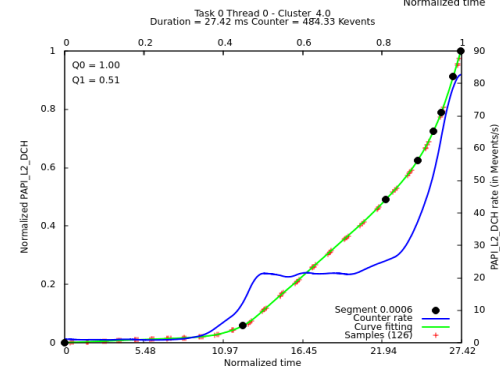
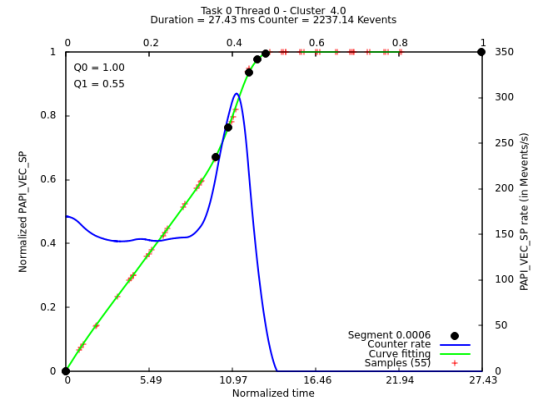
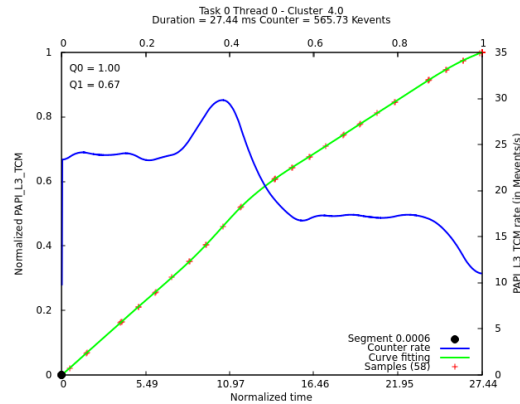
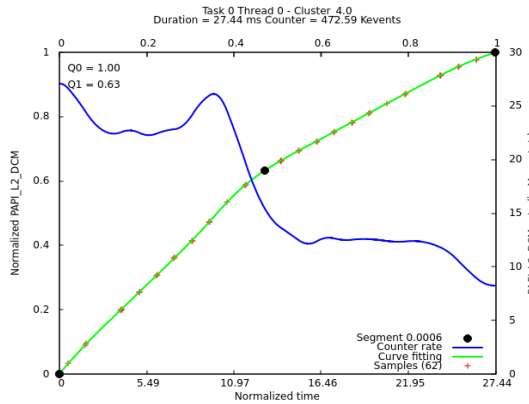
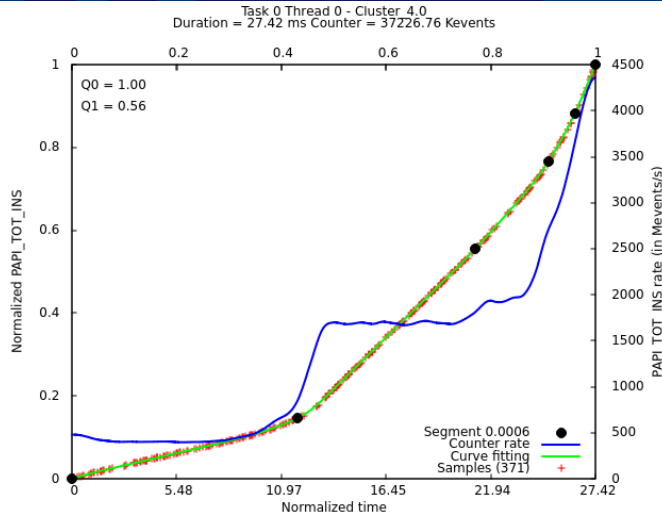
- Between processes
- 3 Algorithmic phases
- Impact of multicore sharing



Correlating counters

CG-POP

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



www.bsc.es



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Methodology

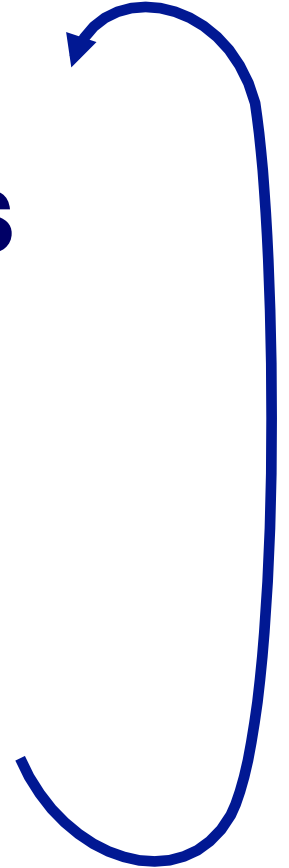
PATC Parallel Programming Workshop
Oct 14-18 2013

Help generate hypotheses

Help validate hypotheses

Qualitatively

Quantitatively



Tools: mechanisms and navigation

« The tools are instruments to address the questions

« Need to know how to use

– First learn to navigate with the tool

« How to load configurations, zoom, fit coloring scales

« How to read

« How to generate timelines form tables

« Second

« Develop a basic understanding of the process of generation of the timelines and histograms.

**Paraver Tutorial:
Introduction to Analysis with Paraver (MPI)**

First steps

- ⌘ Parallel efficiency – percentage of time invested on computation
 - Identify sources for “inefficiency”:
 - load balance
 - Communication /synchronization

- ⌘ Serial efficiency – how far from peak performance?
 - IPC

- ⌘ Scalability – code replication?
 - Total #instructions

- ⌘ Behavioral structure? Variability?

Paraver Tutorial: Introduction to Paraver and Dimemas methodology

Presenting application performance

Factors modeling parallel efficiency

- Load balance (LB)
- Communication
 - Micro load balance (μ LB) or serialization
 - Transfer

CommEff

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

Factors describing serial behavior

- Computational complexity: **#instr**
- Performance: **IPC**

$$\eta_{instr} = \frac{\#instr_0}{\#instr_P}$$

$$\eta_{IPC} = \frac{IPC_P}{IPC_0}$$

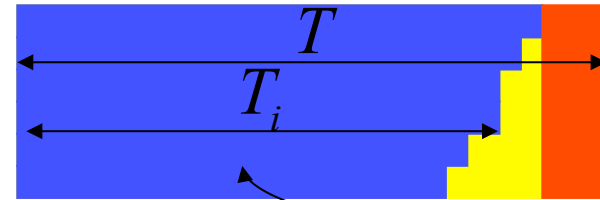
Overall Scaling model

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

Scaling model

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

Directly from real execution metrics



$$eff_i = \frac{T_i}{T}$$

$$CommEff = \max(eff_i)$$

$$LB = \frac{\sum_{i=1}^P eff_i}{P * \max(eff_i)}$$

IPC
#instr

	Outside MPI	MPI_Recv	MPI_Isend	MPI_Irecv
THREAD 1.130.1	87,53 %	9,31 %	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 %
Total	9.309,74 %	306,53 %	1.411,18 %	3,83 %
Average	69,00 %	2,30 %	10,69 %	0,03 %
Maximum	88,18 %	9,62 %	54,97 %	0,04 %
Minimum	30,67 %	0,00 %	0,00 %	0,02 %
StDev	15,27 %	6,06 %	21,40 %	0,00 %
Avg/Max	0,79	0,03	0,19	0,81

η

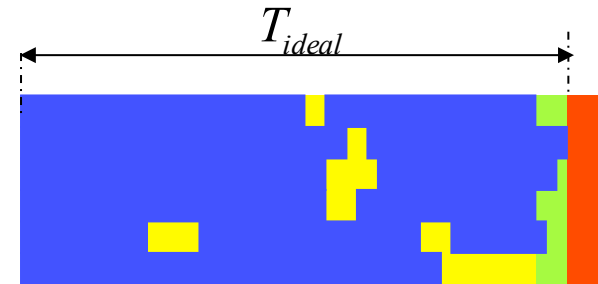
CommEff

LB

Scaling model

- « Dimemas simulation with ideal target
 - Latency = 0; BW = ∞

$$CommEff = \mu LB * Transfer$$



Migrating/local load imbalance
Serialization

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

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

	Outside MPI	MPI_Recv	MPI_Isend	MPI_Irecv
THREAD 1.130.1	87,53 %	9,31 %	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 %
Total	9.309,74 %	306,53 %	1.411,18 %	3,83 %
Average	69,00 %	2,30 %	10,69 %	0,03 %
Maximum	88,18 %	57,62 %	54,97 %	0,04 %
Minimum	30,67 %	0,00 %	0,00 %	0,02 %
StDev	15,27 %	6,06 %	21,40 %	0,00 %
Avg/Max	0,79	0,03	0,19	0,81

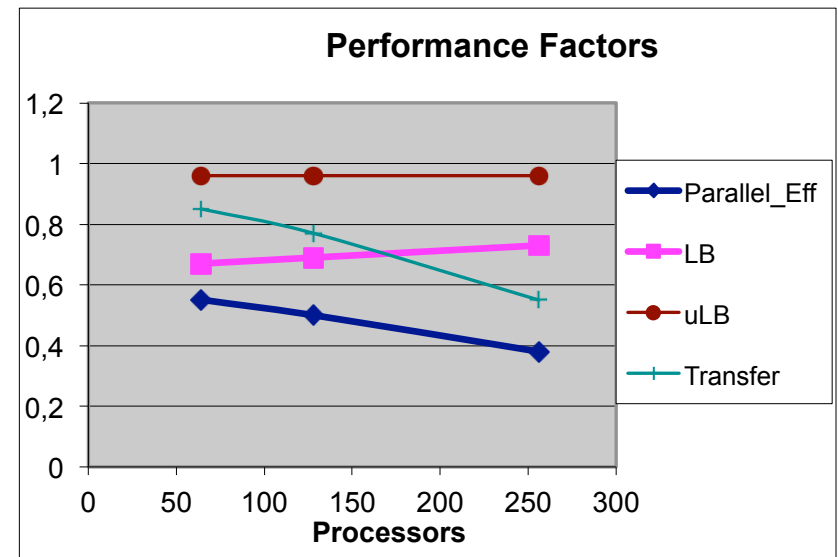
μLB

Scaling model

- « Fundamental behavior
- « Explains bottleneck ...
- « ...how they migrate ...
- « ... and combined effect

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

GROMACsv4.5

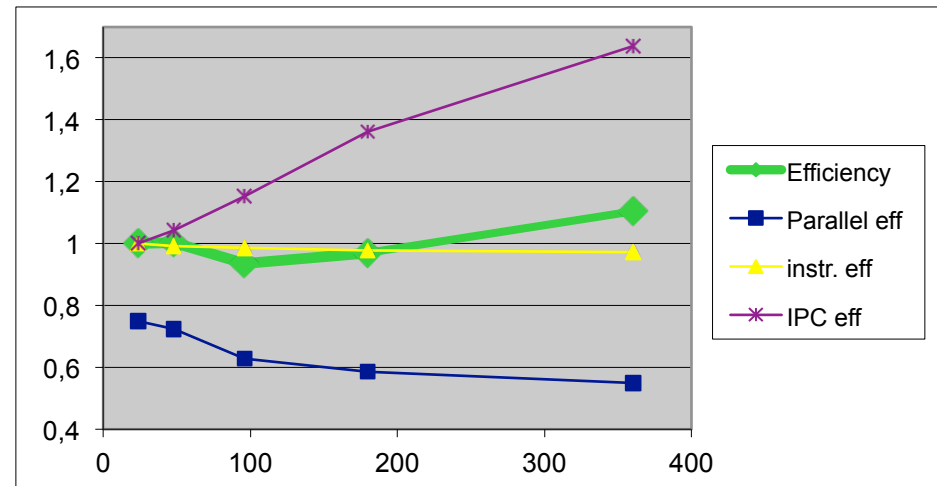
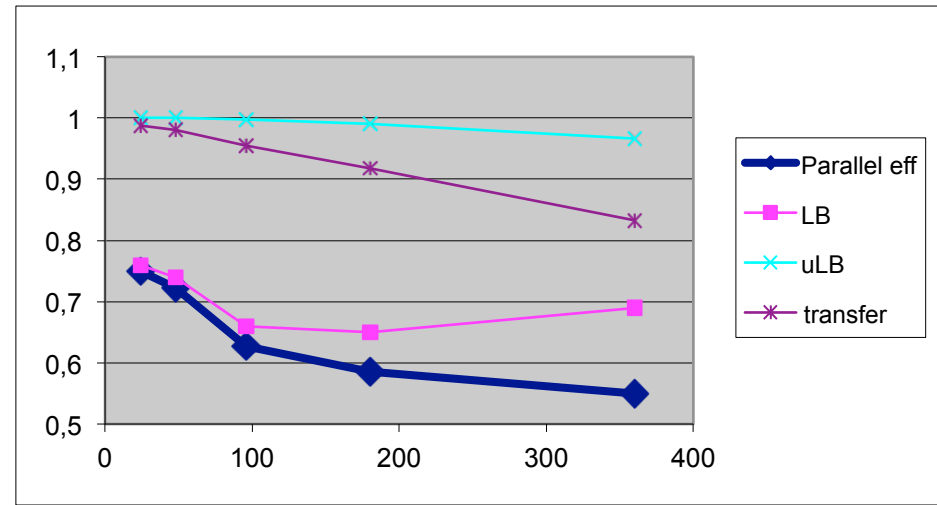
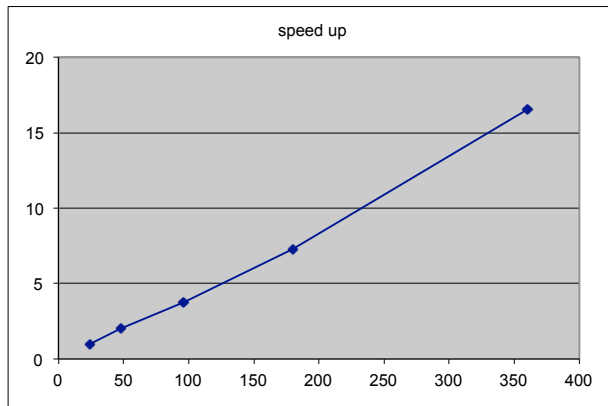


Modelling efficiency

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

CG-POP mpi2s1D - 180x120

Good scalability !!
Should we be happy?



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

« www.bsc.es/paraver

- 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



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación

THANKS

www.bsc.es



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Detailed material

PATC Parallel Programming Workshop
Oct 14-18 2013

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

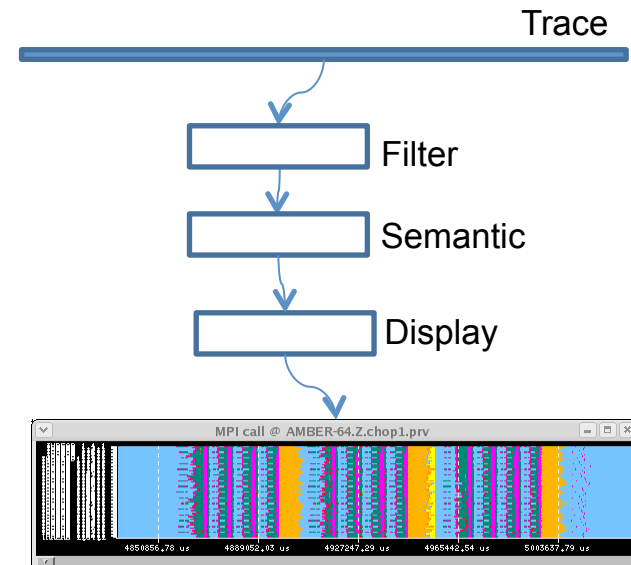
- A sequence of events $Ev_i, i \in N$, states $St_i, i \in N$ and communications $C_i, i \in N$
- For each event let $T(Ev_i)$ be its time and $V(Ev_i)$ its value
- For each state let $T_s(St_i)$ be its start time $T_e(St_i)$ its stop time and $V(St_i)$ its value
- For each Communication let $T_S(C_i)$ 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

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

Function of time

Series of values



Filter module

Communications that pass through the filter

Events that pass through the filter

Property	Value
Logical	<input checked="" type="checkbox"/>
Physical	<input type="checkbox"/>
Comm from	All;
From/To Op	And
Comm to	All;
Comm tag	All;
Tag/Size Op	And
Comm size	All;
Comm bandwic	All;
Events	
Event type	[x,y]; 50000001;50000003
Function	[x,y]
Types	50000001;50000003
Type/Value Op	And
Event value	All;
Function	All
Values	

- All
- !=
- >
- <
- =
- None
- [x,y]

Show list of event types

Semantic module: Control

Paraver

File Help

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

- MPI call profile
- MPI call duration
- MPI call
- MPI call duration

Window properties

Name	MPI call
Begin time	506,27 us
End time	53.300,79 us
Semantic Minimum	0
Semantic Maximum	70
Level	Task
Time Unit	Microseconds

Filter

Semantic

Top Compose 1	As Is
Top Compose 2	As Is
Compose Task	As Is
Task	Thread i
Object	"1"
Compose Thread	As Is
Thread	Last Evt Val

- System
- Node
- CPU
- Workload
- Application
- Task
- Thread

- Adding
- AddingSign
- Average
- Maximum
- Minimum
- Activity
- In Activity
- Mode
- Thread i
- Add Tasks

- As Is
- Sign
- 1-Sign
- Mod
- Mod+1
- Div
- Prod
- Add
- Subs
- Select Range
- Select Range D
- Is In Range
- Is In Range D
- Is Equal
- Is Equal (Sign)
- Floor
- Ceil
- Round
- Abs
- Stacked Val
- In Stacked Val
- Nesting level
- Delta
- Burst Time
- Join Bursts

- State As Is
- Useful
- State Sign
- Given State
- In State
- Not In State
- State Record Dur.

- Last Evt Type
- Last Evt Val
- Last Evt Val w/o Bursts
- Next Evt Type
- Next Evt Val
- Avg Next Evt Val
- Avg Last Evt Val
- Given Evt Val
- In Evt Val
- Int. Between Evt
- Not In Evt Val
- In Evt Range
- Event Bytes
- Event Sent Bytes

- Last Tag
- Comm Size
- Comm Recv. Partner
- Comm Partner
- Last Send Dur.
- Next Recv Dur.
- Send Bytes in Transit
- Send Messages in Transit
- Send BandWidth
- Recv Bytes in Transit
- Recv Messages in Transit
- Recv BandWidth
- Recv Negative Messages
- Recv Negative Bytes
- Number Of Receives
- Number Of Receive Bytes

- State
- Event
- Communication
- Object

- Application ID
- TaskID
- Thread ID
- Node ID
- Cpu ID
- In Appl ID
- In Task ID
- In Thread ID
- In Node ID
- In Cpu ID

Semantic module

« 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
$$s(t) = \sum_j Sz(C_j), j | (T_S(C_j) < t) \wedge (T_R(C_j) > t) \wedge (Dir(C_j) == send)$$
- Send Bandwidth
$$s(t) = \sum_j \frac{Sz(C_j)}{T_R(C_j) - T_S(C_j)}, j | (T_S(C_j) < t) \wedge (T_R(C_j) > t) \wedge (Dir(C_j) == send)$$
- Msgs in transit
$$s(t) = \sum_j sign(j), j | (T_S(C_j) < t) \wedge (T_R(C_j) > t) \wedge (Dir(C_j) == send)$$
- Recv. Bandwidth
$$s(t) = \sum_j \frac{Sz(C_j)}{T_R(C_j) - T_S(C_j)}, j | (T_S(C_j) < t) \wedge (T_R(C_j) > t) \wedge (Dir(C_j) == recv)$$
- Rec. Negative Msgs
$$s(t) = \sum_j sign(j), j | (T_R(C_j) < t) \wedge (T_S(C_j) > t) \wedge (Dir(C_j) == recv)$$
- Comm. Partner
$$s(t) = Partner(C_j), j | (T_S(C_j) < t) \wedge (T_R(C_j) > t)$$
- Bytes btw. Events
$$S(i) = \sum_j Sz(C_j), j | T_S(C_j) \in [T(Ev_i), T(Ev_{i+1})) \vee T_R(C_j) \in [T(Ev_i), T(Ev_{i+1}))$$

Composition

$$\llcorner S'(t) = f(S(t))$$

$$S' = f \circ S$$

– Sign

$$S'(t) = \text{sign}(S(t))$$

– 1-sign

$$S'(t) = 1 - \text{sign}(S(t))$$

– Select range

$$S'(t) = S(t) \in [a, b] ? S(t) : 0$$

– Sign ° Is equal

$$S'(t) = \text{sign}(S(t) = a) ? S(t) : 0$$

– Delta

$$S'(t) = S_{i+1} - S_i$$

– Stacked value

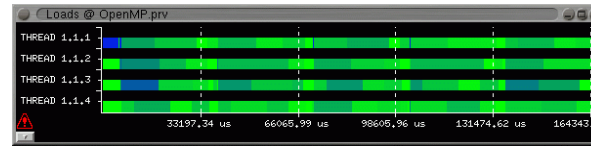
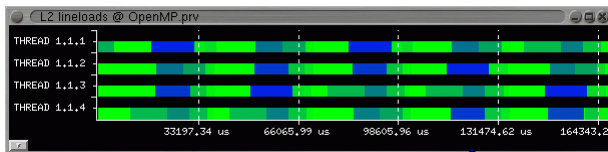
Semantic module

Derived windows

– Point wise operation

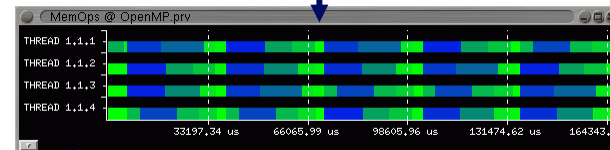
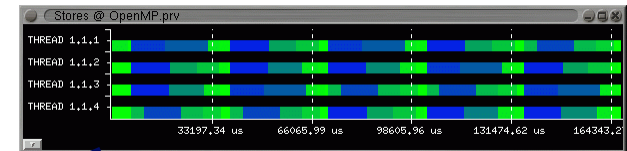
- $S = \alpha * S^a \langle op \rangle \beta * S^b$
- $\langle op \rangle : +, -, *, /, \dots$

L2 Line Loads



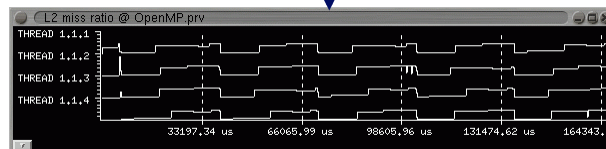
Loads

Stores



Mem Ops

x100



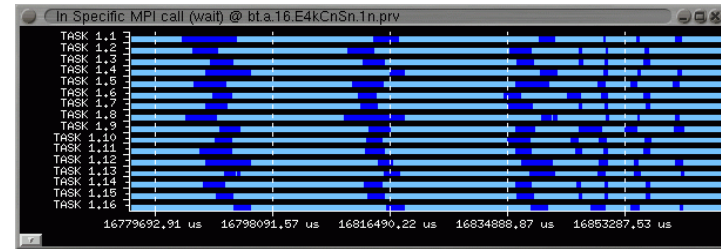
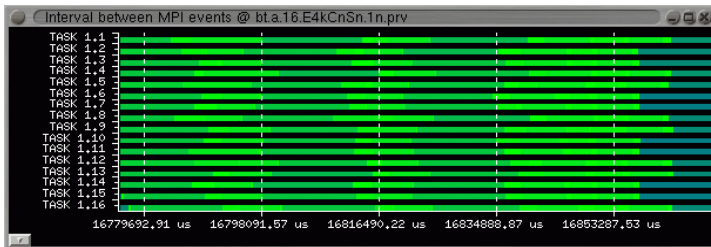
L2 miss ratio

Semantic module

Derived windows

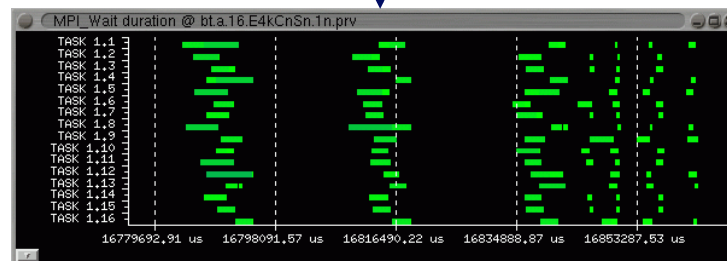
– Point wise operation

- $S = \alpha * S^a \langle op \rangle \beta * S^b$
- $\langle op \rangle : +, -, *, /, \dots$



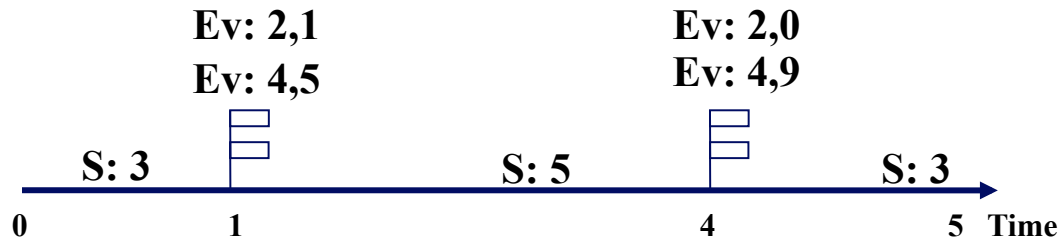
Interval between MPI events

In MPI call

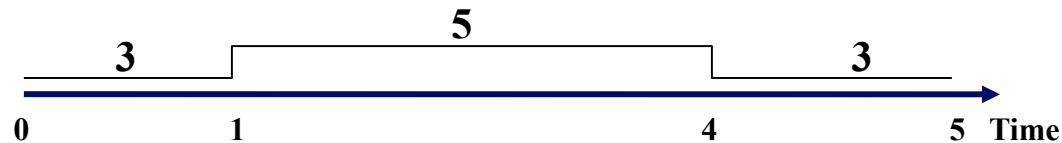


MPI call duration

Semantic module: Examples

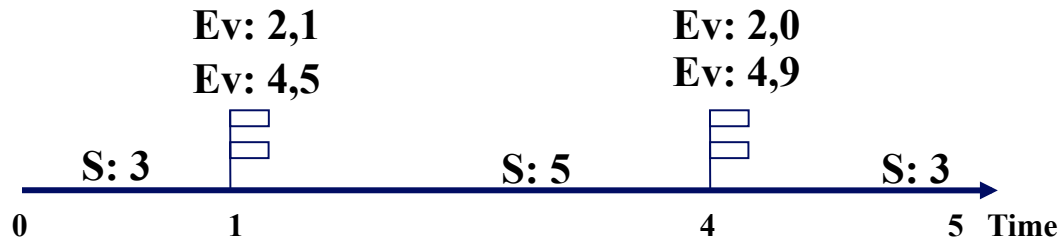


⌋ Thread function: State as is



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

Semantic module: Examples



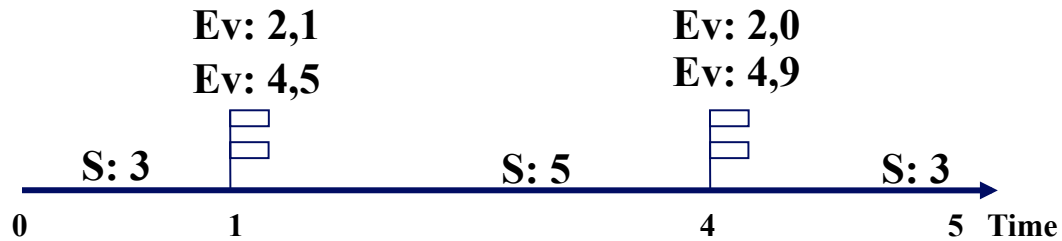
Filter: type == 2

– Thread function: Last event value



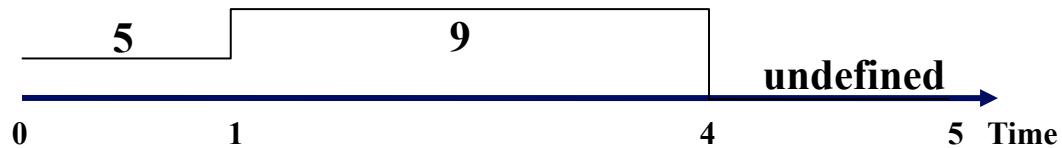
- Useful for
 - In parallel region
 - Mutual exclusion
 - Variable values: iteration,.....

Semantic module: Examples



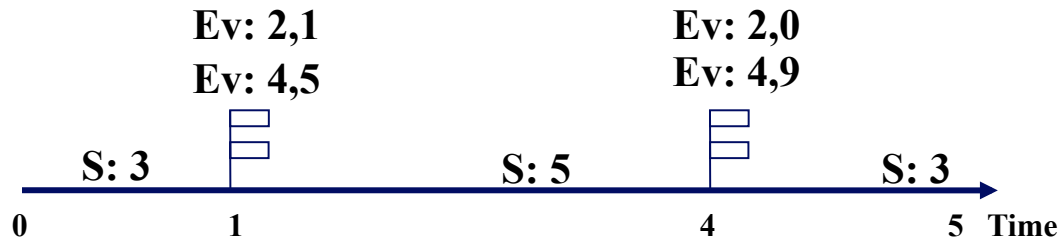
Filter: type == 4

– Thread function: Next event value



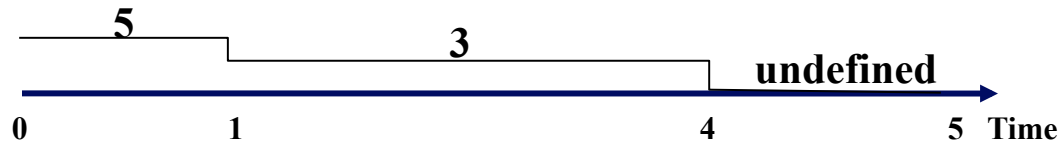
- Useful for
 - Hwc events (TLB, L1 misses,...) within interval

Semantic module: Examples



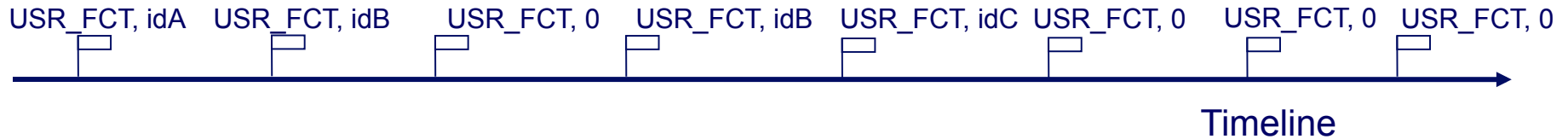
Filter: type == 4

– Thread function: Average next event value

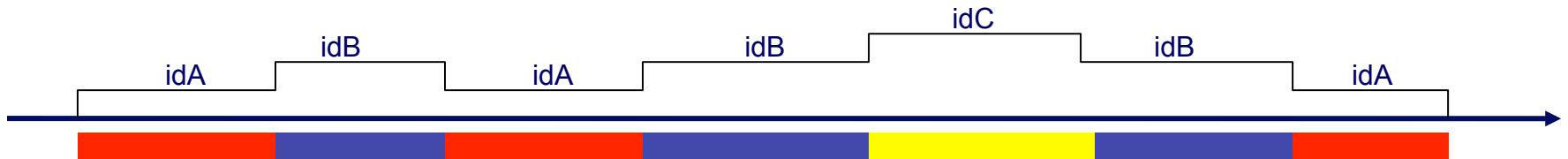


- Useful for
 - Hwc events (TLB, L1 misses,...) per time unit within interval

Semantic module: Examples



- Filter: `type == USR_FCT`
Thread function: Last event value
Compose: Stacked value



- Useful for
 - Routine

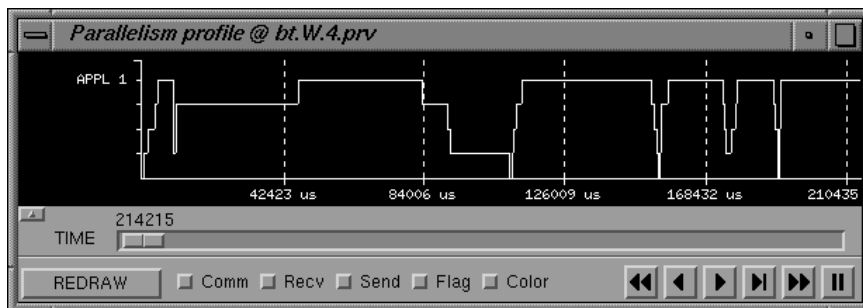
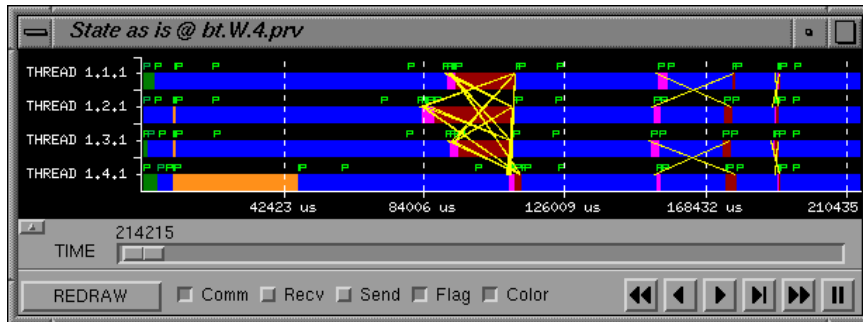
Semantic module perspective

Process model

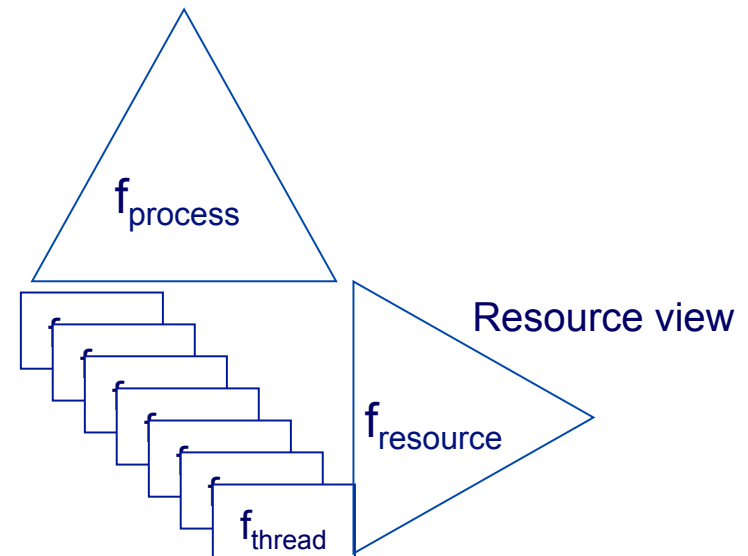
- Thread, task, application, workload

Resource model

- CPU, node, system

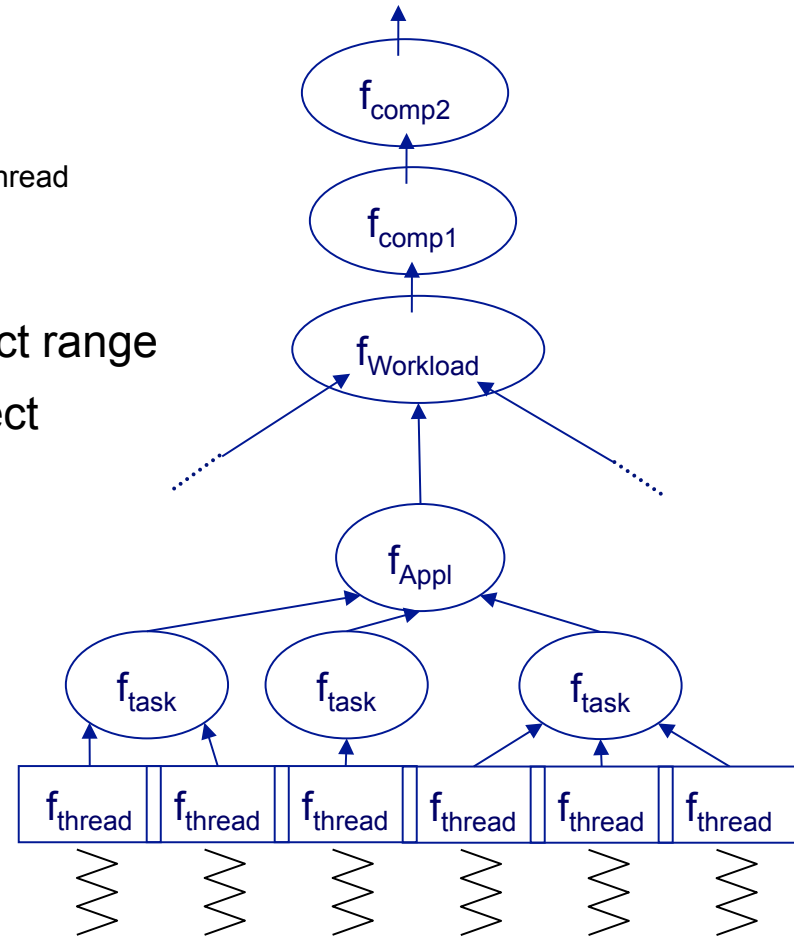


Process view



Process model perspective

- Semantic value: $S(t)$
- $S = f_{\text{comp2}} \circ f_{\text{comp1}} \circ f_{\text{Workload}} \circ f_{\text{Application}} \circ f_{\text{task}} \circ S_{\text{thread}}$
- Semantic functions
 - $f_{\text{comp2}}, f_{\text{comp1}}$: sign, mod, div, in range, select range
 - $f_{\text{Application}}, f_{\text{Workload}}$: 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, ...



Resource model perspective

- $Sf_{\text{resource}} = f_{\text{comp2}} \circ f_{\text{comp1}} \circ f_{\text{System}} \circ f_{\text{Node}} \circ f_{\text{CPU}} \circ S_{\text{thread}}$

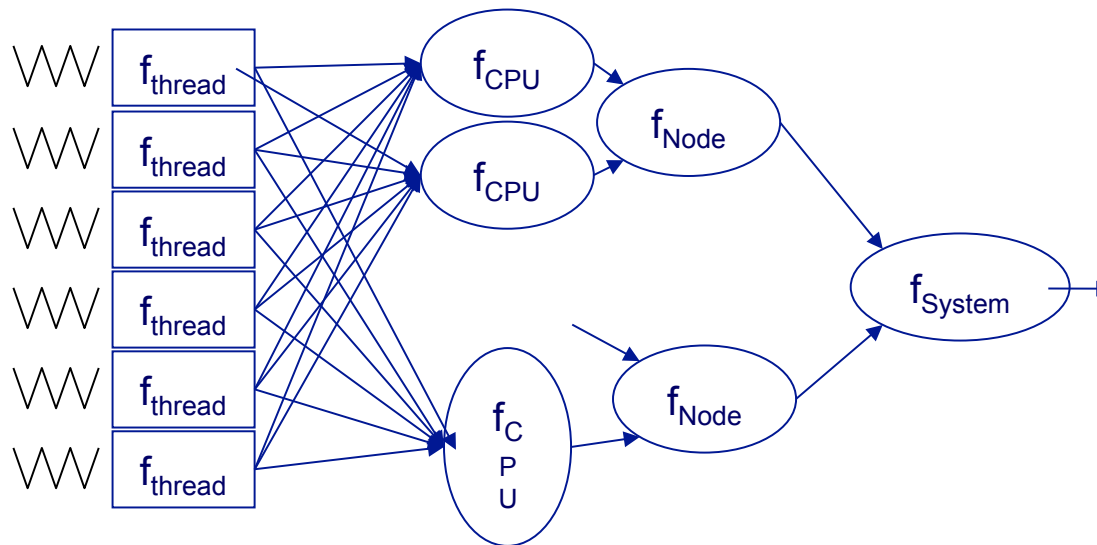
- Semantic functions

- f_{System} : add, average, max, select

- f_{Node} : add, average, max, select

- f_{CPU} : active thread, select

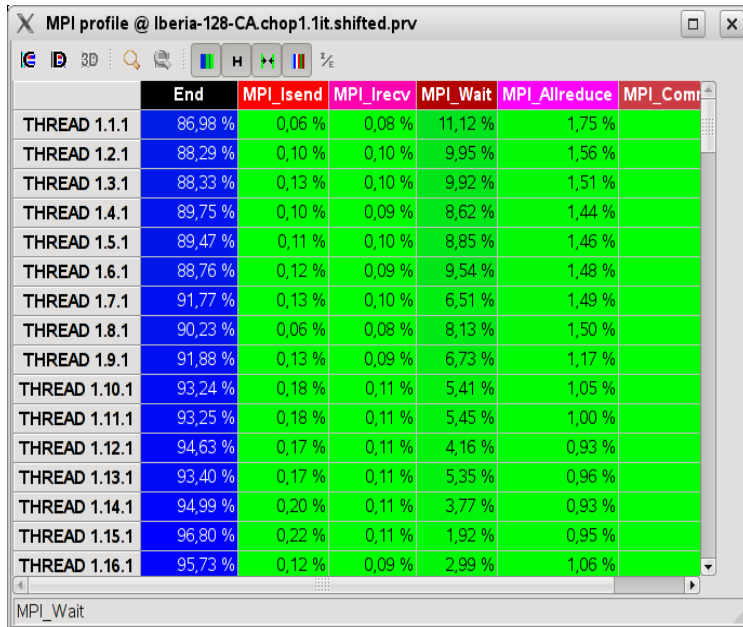
- S_{thread} : in state, useful, given state, next event value, thread_id



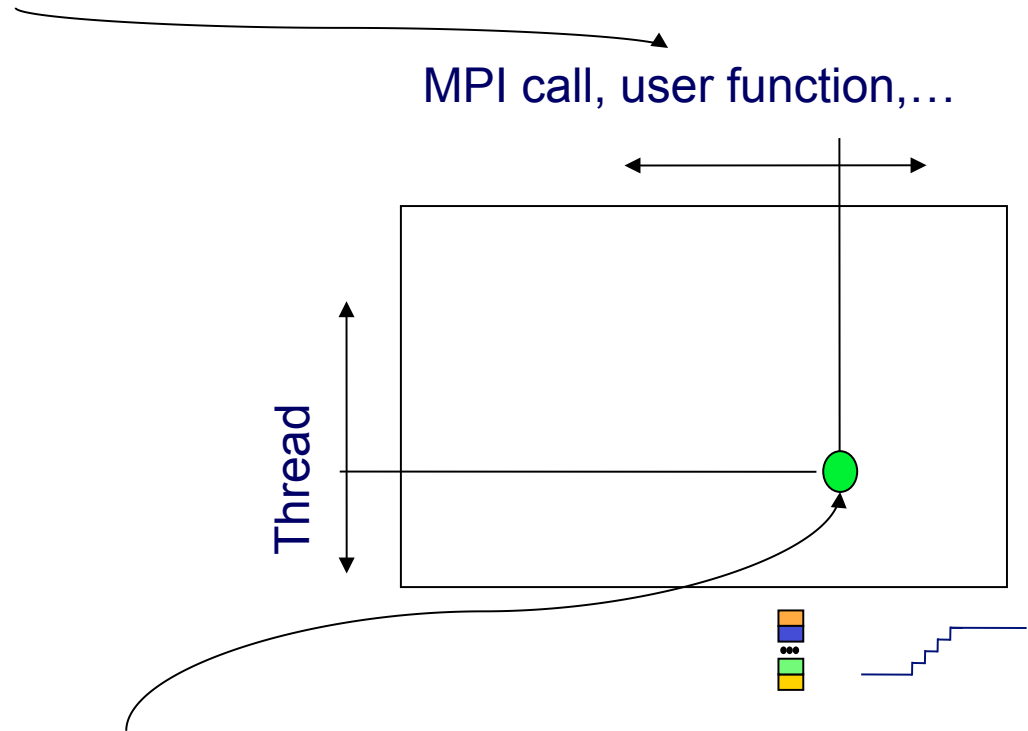
Analysis Module

How to read profiles

One column 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 %	



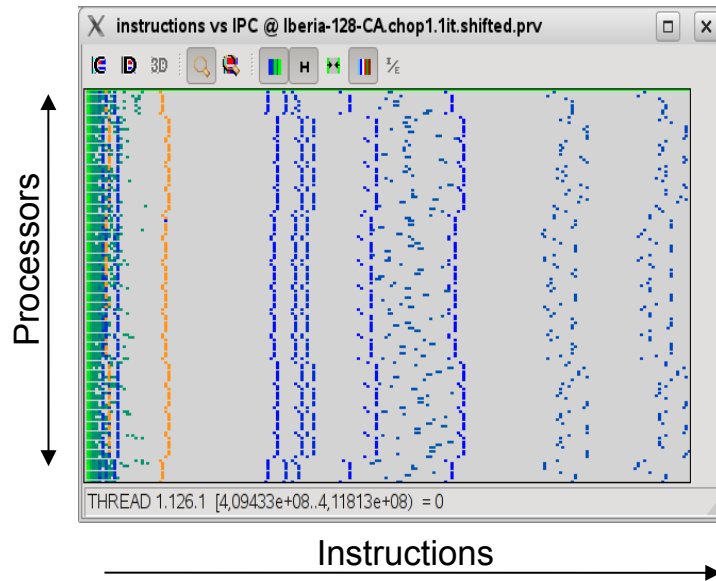
Value/color is a statistic computed for the specific thread when control window had the value corresponding to the column

Relevant statistics:

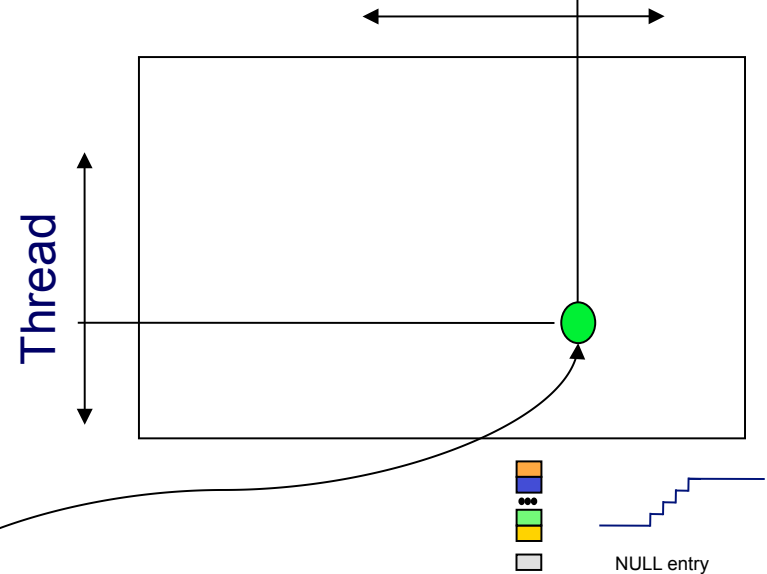
Time, %time, #bursts, Avg. burst time
Average of **Data window**

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**

Tables

⌘ 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

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

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

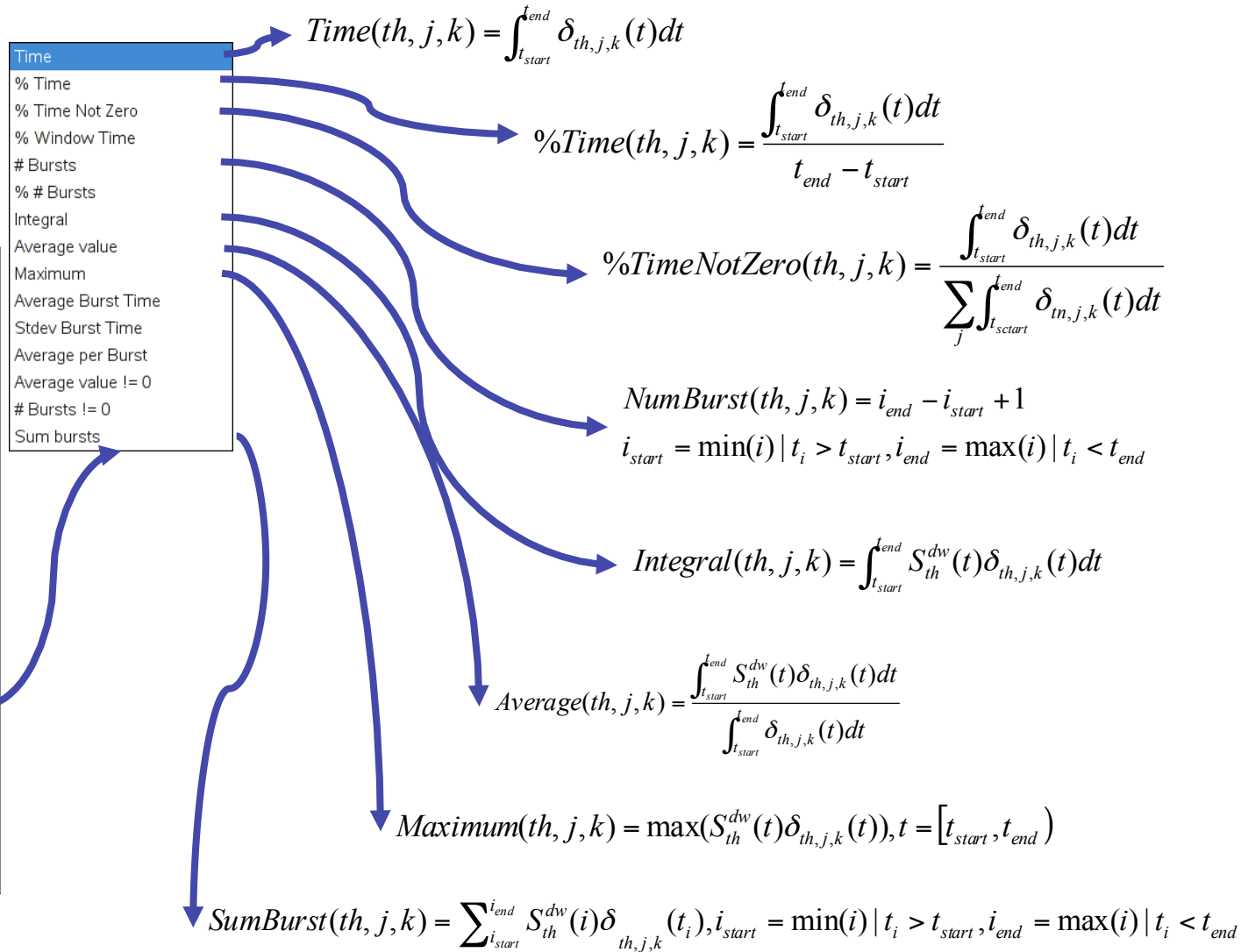
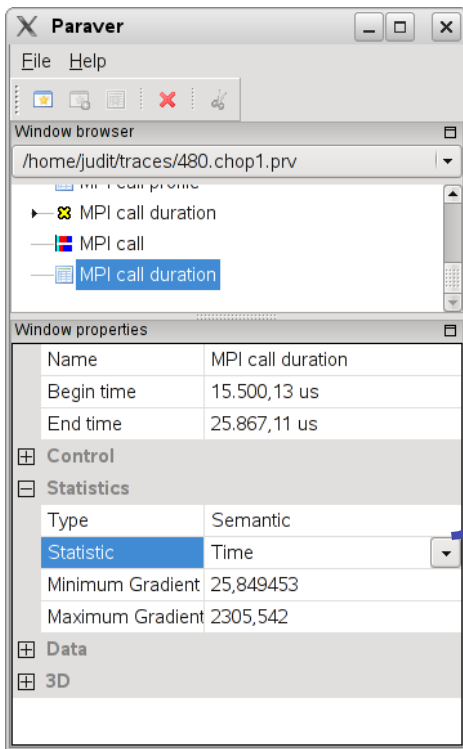
Identify regions with cw 's within the (j,k) bin

⌘ 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



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

The image displays two screenshots of the Paraver application window. The left screenshot shows the 'Window browser' and 'Files & Window Properties' panels. The right screenshot shows the 'Files & Window Properties' panel with a table of active trace characteristics.

Select to browse in lower panel for traces or cfgs

Select to browse characteristics of active view or table

Active trace

Available views and tables
Active view or table highlighted

Property	Value
End time	11.315.554,69 us
Control	
Window	Interval btw uf events
Minimum	1
Maximum	460001
Delta	2555,555556
Statistics	
Type	Semantic
Statistic	% Time
Minimum Gradient	4,782857
Maximum Gradient	50,018864
Data	
Window	Interval btw uf events
3D	
3rd Window	User function
Minimum	0
Maximum	60
Delta	1
Plane	riemann

Load configuration files

The image shows the Paraver application interface. On the left, the 'File' menu is open, with 'Load Configuration...' selected. An arrow points from this menu item to the 'Load Configuration' dialog box on the right. The dialog box shows a directory tree on the left and a list of files on the right. The file 'point2point' is selected in the list. An arrow points from the 'point2point' file to the 'Abrir' button. Below the dialog box, a list of configuration files is shown, with the first file highlighted. An arrow points from the 'Abrir' button to this list.

Select directory

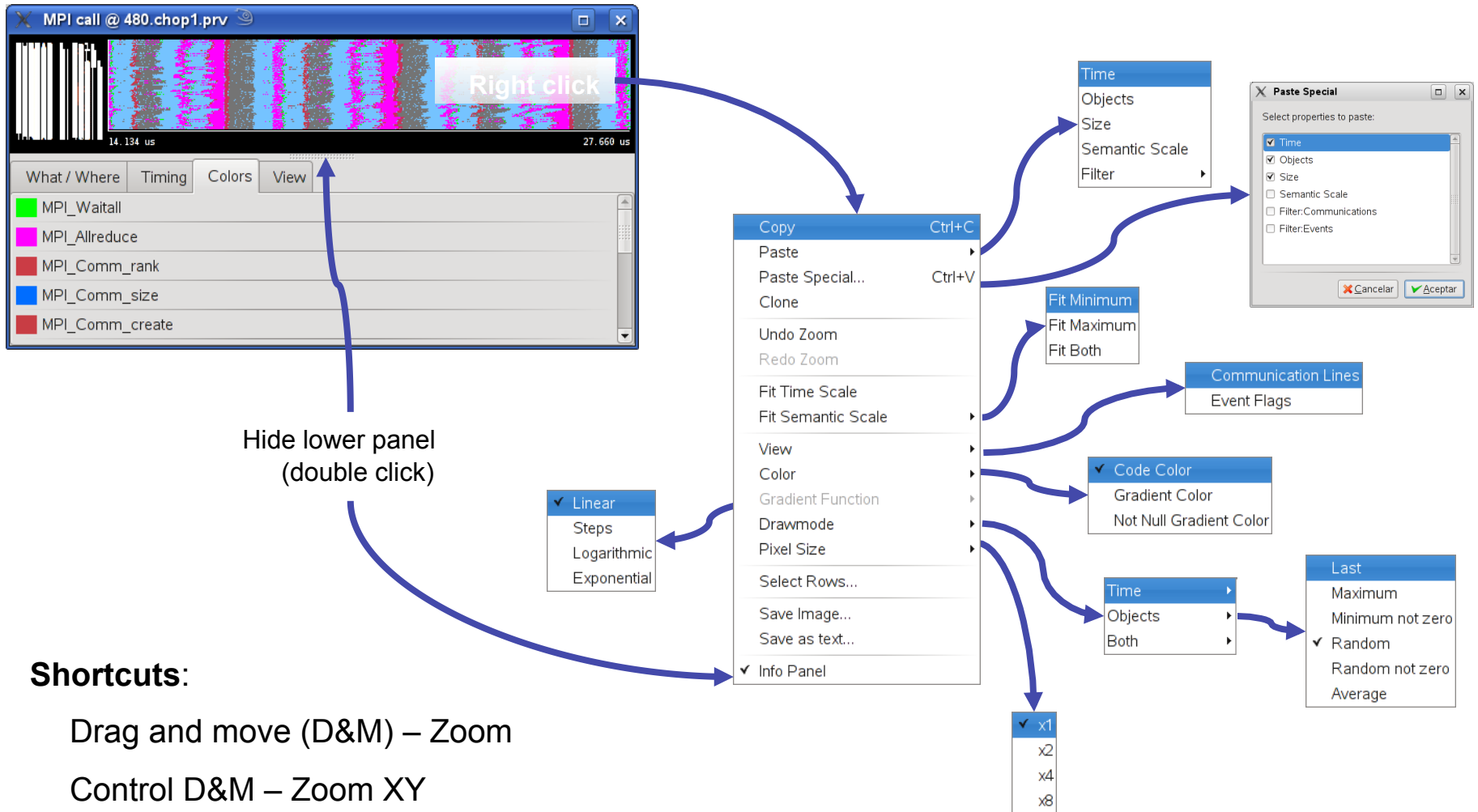
Navigate through directory tree

List of directories and configuration files in current directory

APPLIED TO THE CURRENT TRACEFILE

```
/home/judit/tools/etc/cfgs/counters_PAPI/performance/IPC.cfg
/home/judit/tools/etc/cfgs/mpiViews/MPI_call.cfg
stacked9k.cfg
/home/judit/users/xavit/stacked9k.cfg
/home/judit/users/laurent/3d_count_samples_by_value_at_phases.cfg
/home/judit/users/count_samples.cfg
/home/judit/users/paraver/paraver-cfgs/2dh_function_in_frame_0.cfg
/home/judit/users/paraver/paraver-cfgs/execution_phases.cfg
/home/judit/users/paraver/paraver-cfgs/bug.cfg
/home/judit/traces/jaguar/3d_mpi_duration.cfg
/home/judit/tools/etc/cfgs/mpi/analysis/3dh_duration_MPIcall.cfg
/home/judit/tools/etc/cfgs/counters_PAPI/performance/cycles_per_us.cfg
/home/judit/tools/etc/cfgs/mpiViews/collectives/collective_size.cfg
/home/judit/traces/480.chop.cfg
/home/judit/users/laurent/Transfer Durations.cfg
/home/judit/users/juli/load_cpu_node.cfg
/home/judit/users/juli/NodeMBUsed.cfg
/home/judit/users/juli/MB.cfg
/home/judit/traces/jaguar/large_allreduce_31_1.cfg
/home/judit/users/rhodri/useful_duration.cfg
```

Navigation



Shortcuts:

Drag and move (D&M) – Zoom

Control D&M – Zoom XY

Shift D&M – Timing

How to generate table and change statistic

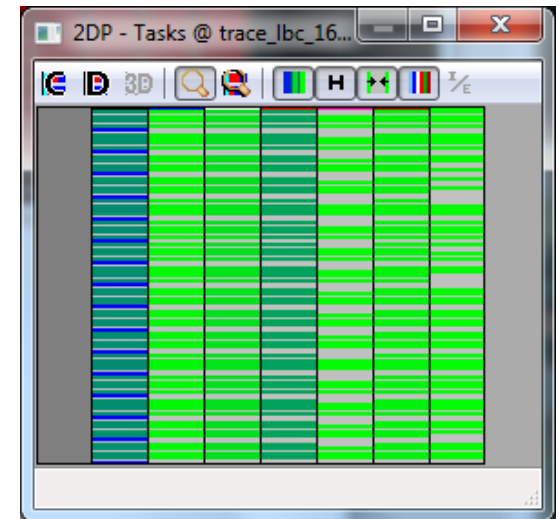
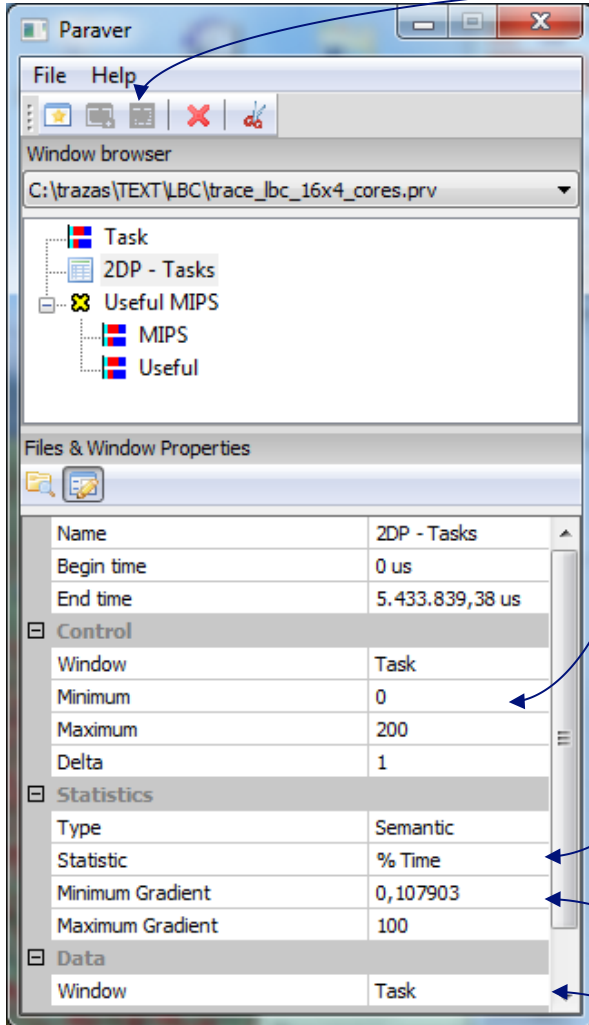
To generate table: click button and select region of the window whose values will determine the columns of the table

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)



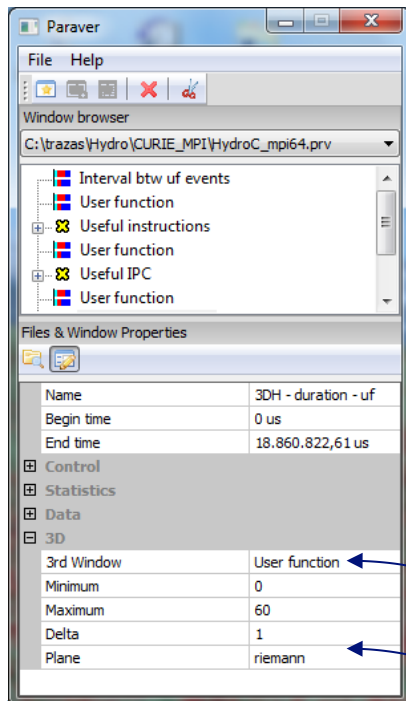
3D tables

☞ 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



The screenshot shows the Paraver application interface. The 'Window browser' pane on the left lists several categories: 'Interval btw uf events', 'User function', 'Useful instructions', 'Useful IPC', and 'User function'. The 'Files & Window Properties' pane on the right shows details for a window named '3DH - duration - uf', including its begin and end times. At the bottom, the '3D' section contains a table with the following data:

3rd Window	User function
Minimum	0
Maximum	60
Delta	1
Plane	riemann

**3D control window:
determines planes**

**Actual Plane on
display**

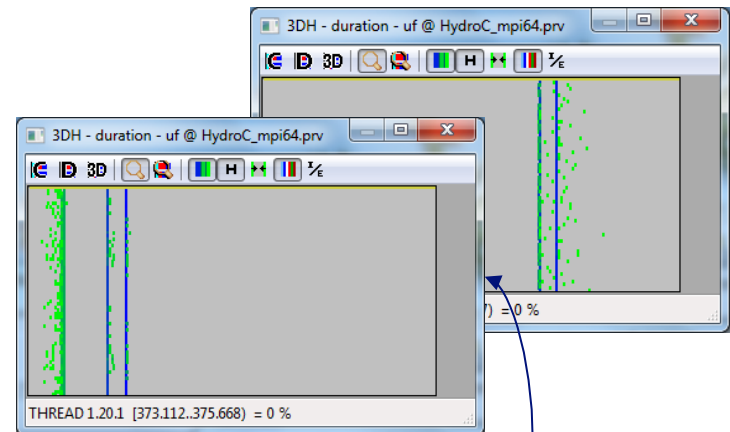


Table information and control

Create a new table

Paraver window properties for 'MPI call profile':

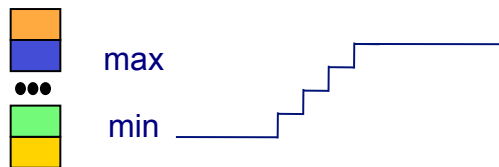
- Name: MPI call profile
- Begin time: 15.500,13 us
- End time: 25.867,11 us
- Control: Window (MPI call), Minimum (0), Maximum (124), Delta (1)
- Statistics: Type (Semantic), Statistic (% Time), Minimum Gradient (1,188944), Maximum Gradient (49,410556)
- Data window: Data (selected), 3D (unchecked)

Region analyzed

Bin definition

Change Data window

Color encoding



Display whole table / cell text

Color/not cells

Transpose

Hide null columns

MPI call profile @ 480.chop1.prv

	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 %

Activate 3D analysis

Table information and control

Open Data window

Open Control window

Open 3D window

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

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

Right click

- Copy (Ctrl+C)
- Paste
- Paste Special... (Ctrl+V)
- 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 Accept

- ✓ Linear
- Steps
- Logarithmic
- Exponential

- Semantic
- Objects
- Both

Generate ASCII file with table data

Selected plane

Shortcuts (zoom mode only):

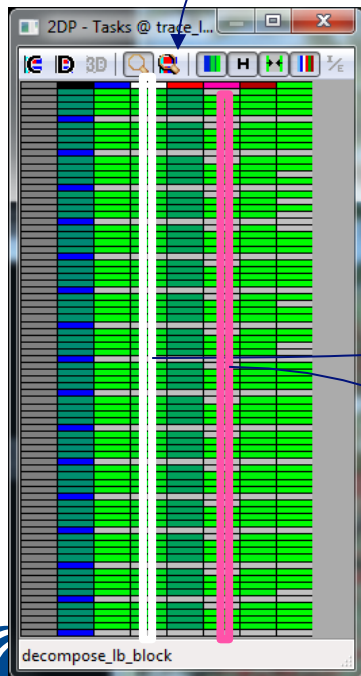
Drag and move (D&M) – Zoom

Control D&M – Zoom XY

From tables to timelines

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

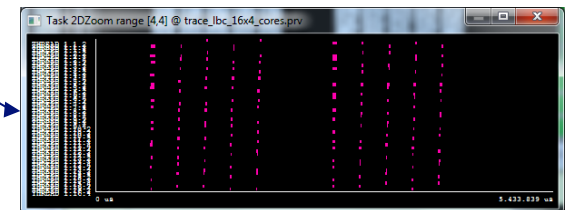
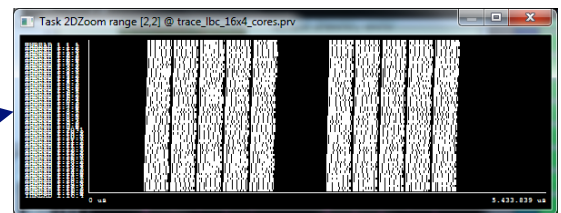
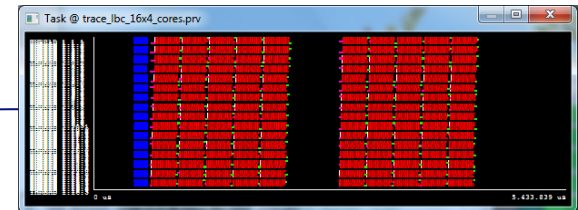
Click button and select column(s)



Will automatically generate derived views from the global view

Only showing when is routine white executing

Only showing when is routine pink executing

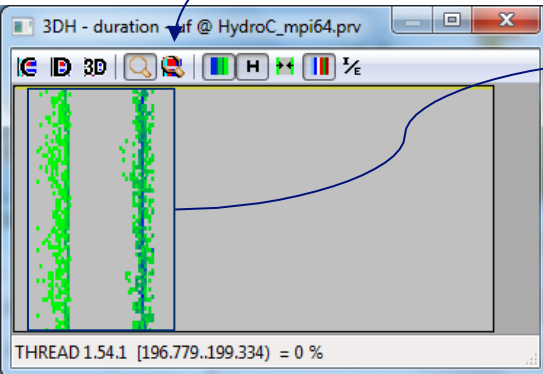


From tables to timelines

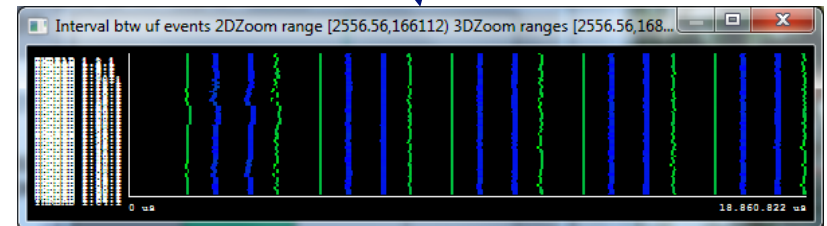
« Where in the timeline do the values in certain table columns appear?

- ie. want to see where the timeline happen computation bursts of a given length?

Click button and select column(s)



Will automatically generate



3D histogram of duration of routine foo

Only showing duration of routine foo

Trace manipulation

Handling very large traces

« 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

Filtering very large traces

Trace to which it will be applied

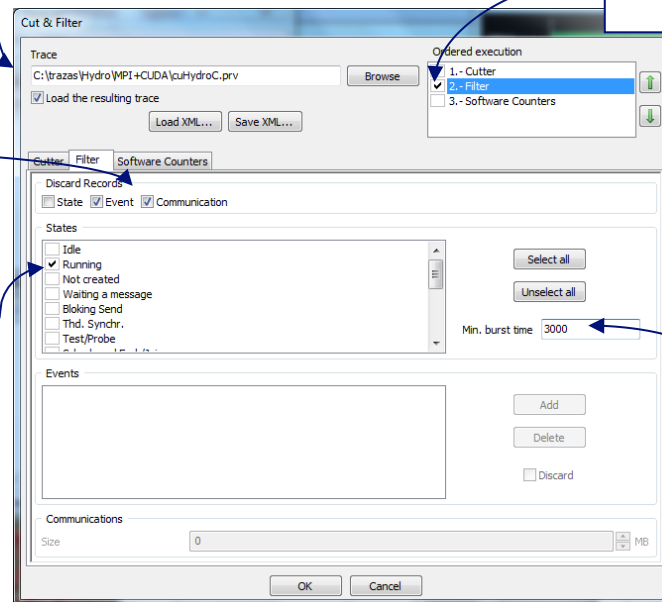
A trace with
basename.filter1.prv will be
generated

Select filtering
option

Discard events and
communications

Keep only Running bursts
....

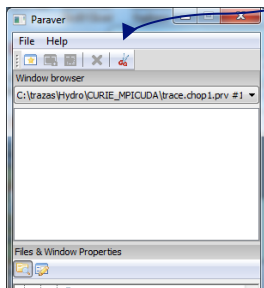
--- longer than
3000 ns



Cutting very large traces

« Load a filtered trace and use the scissors tool

Scissors tool



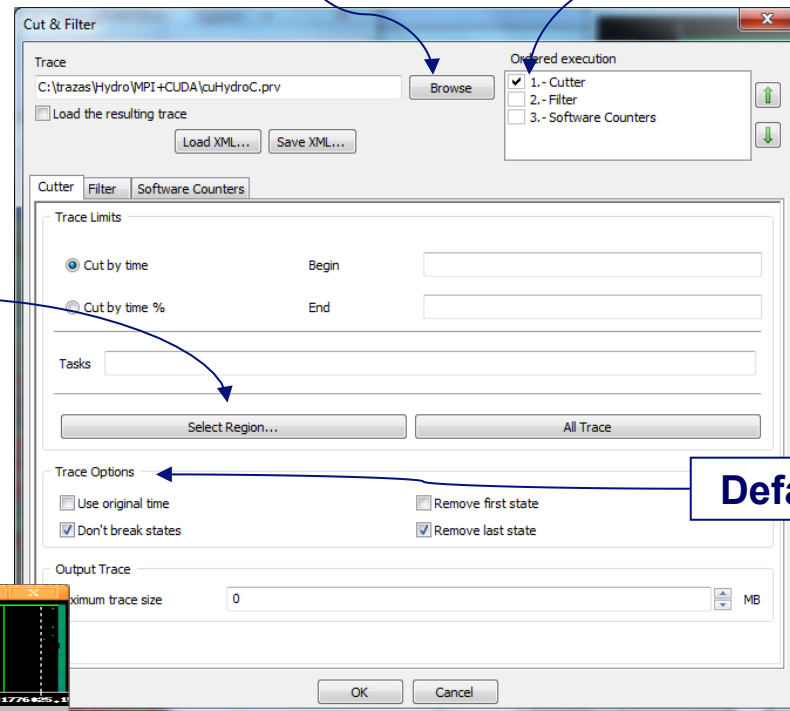
Click to select region

Select time interval by clicking left and right limits in a window of the filtered trace previously loaded

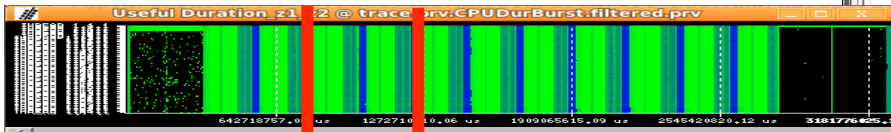
Recommended cuts within long computation bursts

Browse to select file from which the cut will be obtained

Select cutter



Default setups



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

trace.sh

```
#!/bin/bash

export EXTRAE_HOME=/export/hopsa/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 control .xml

```
<?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" />
</mpi>
```

Activate MPI tracing and emit hardware counters at MPI calls

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

Do not activate OpenMP tracing

```
<callers enabled="yes">
  <mpi enabled="yes">1-3</mpi>
  <sampling enabled="no">1-5</sampling>
</callers>
```

Emit call stack information (number of levels) at acquisition points

...

Trace control .xml (cont)

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)

extrae.xml (cont)

Emit counters or not

```
<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_L2_DCM
```

```
      <sampling enabled="no" frequency="100000000">PAPI_TOT_CYC
```

```
    </set>
```

```
    <set enabled="yes" domain="user" changeat-globalops="5">
```

```
      PAPI_TOT_INS,PAPI_FP_INS,PAPI_TOT_CYC
```

```
    </set>
```

```
  </cpu>
```

```
<network enabled="no" />
```

```
<resource-usage enabled="no" />
```

```
<memory-usage enabled="no" />
```

```
</counters>
```

...

When to rotate between groups

Groups

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

OS info (context switches,....)

Trace control .xml (cont)

mpitrace.xml (cont)

Control of emitted trace ...

```
<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>
```

... name, tmp and final dir

...

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

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

Size of in core buffer (#events)

Trace control .xml (cont)

mpitrace.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)

mpitrace.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

Trace control .xml (cont)

mpitrace.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"
```

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

```
>  
  $TRACENAME$  
</merge>
```

... into this trace name

```
</trace>
```

LD_PRELOAD library selection

Library depends on programming model

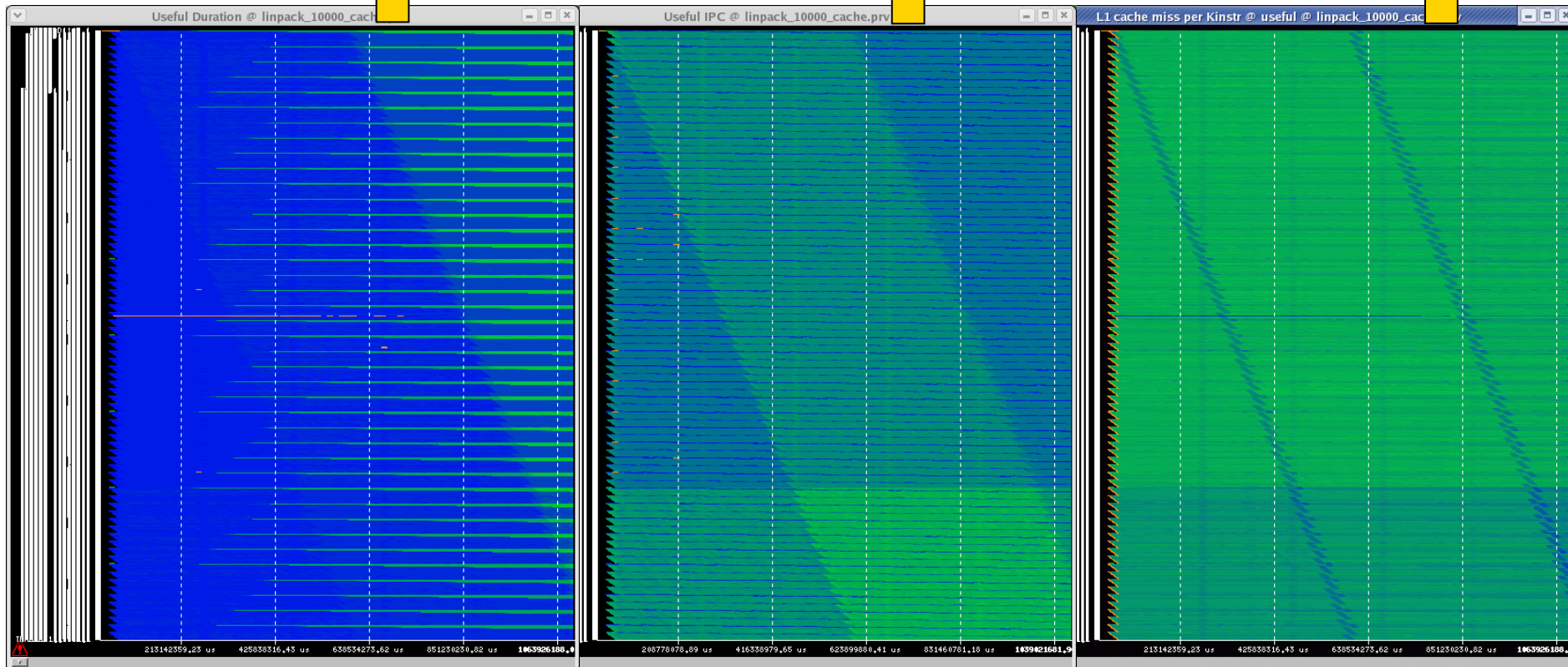
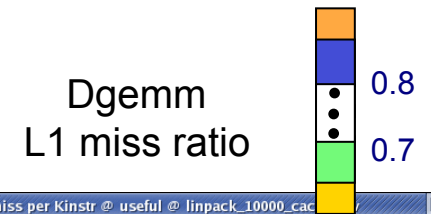
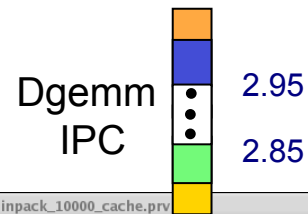
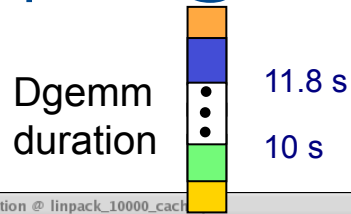
Programming model	Library
Serial	libseqtrace
Pure MPI	libmpitrace[f] ¹
Pure OpenMP	libompitrace
Pure Pthreads	libpttrace
CUDA	libcudatrace
MPI + OpenMP	libompitrace[f] ¹
MPI + Pthreads	libptmpitrace[f] ¹
Mpi + CUDA	libcudampitrace[f] ¹

¹ for Fortran codes

Scalability

Scalability of Presentation

Linpack @ Marenostrom: 10k cores x 1700 s

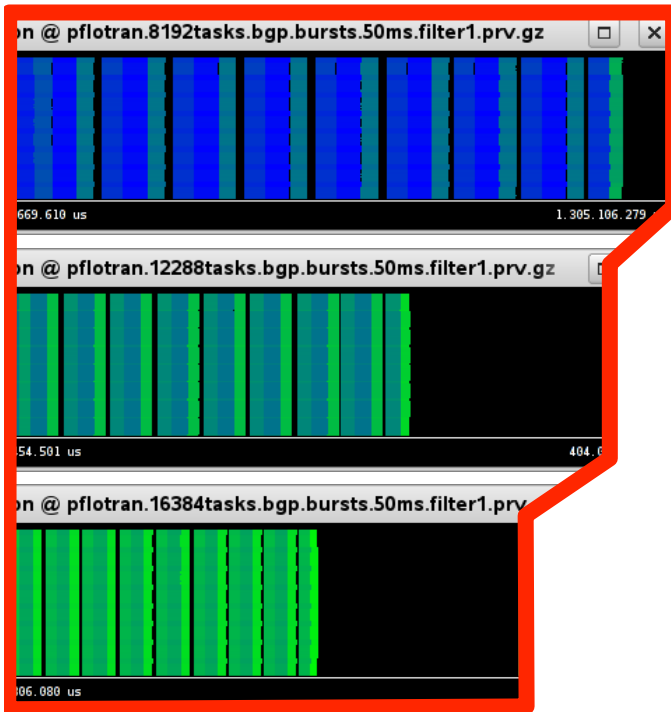


Scalability of analysis

Jugene

~ 105 seconds

8K cores

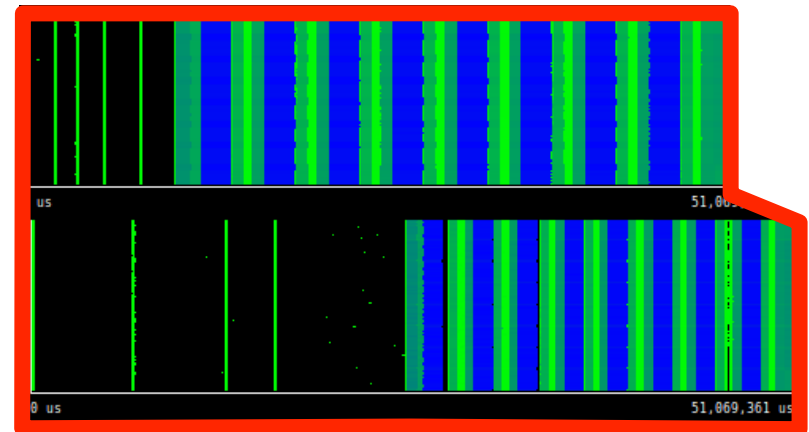


Tran

Flow

Jaguar

~ 47 seconds



Flow

Tran

PFLOTRAN

Data reduction techniques

« 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)

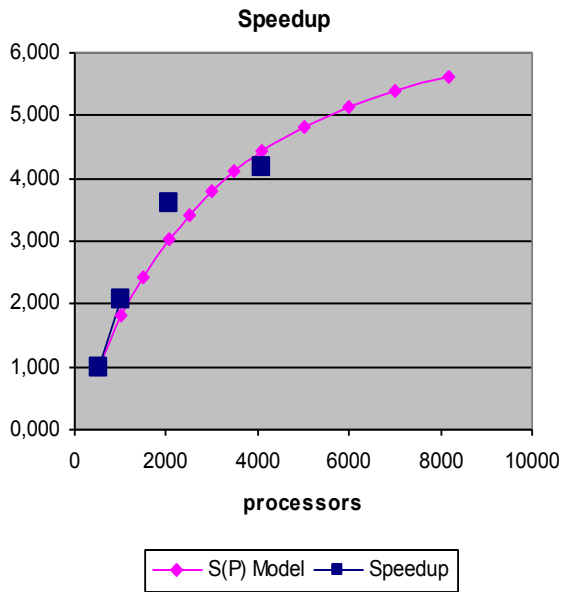
« Representative cuts

- Emit full detail only on selected intervals, representative of full program execution

« On and off line combinations

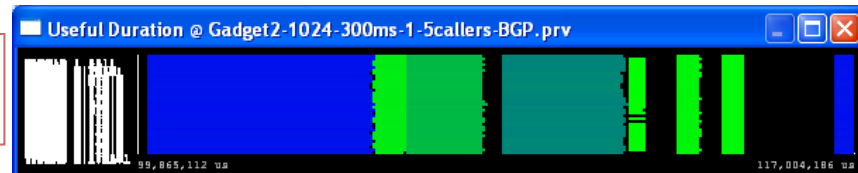
- By instrumentation
- By paraver filtering

Software counters

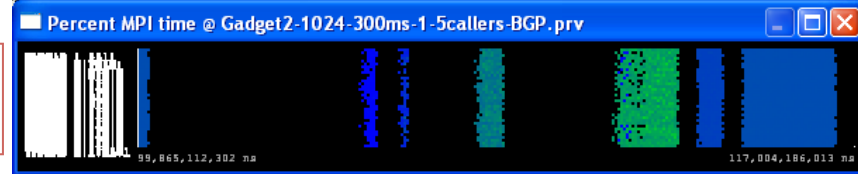


GADGET, PRACE Case A, 1024 procs

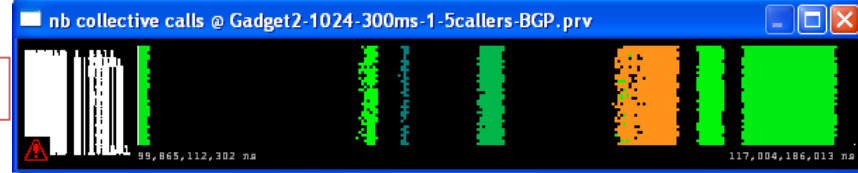
Useful duration



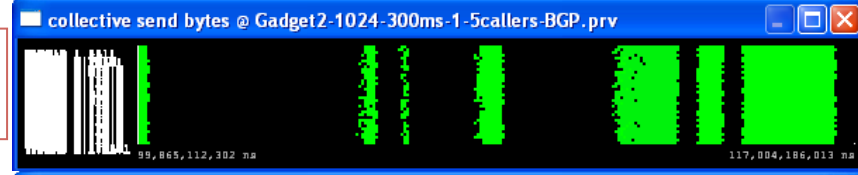
% MPI time



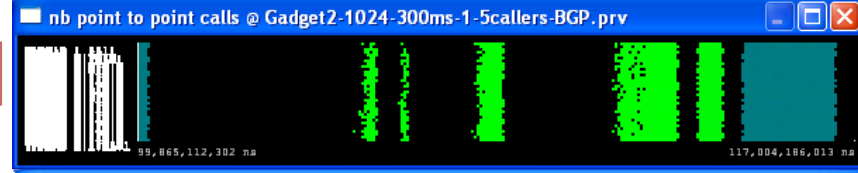
collectives



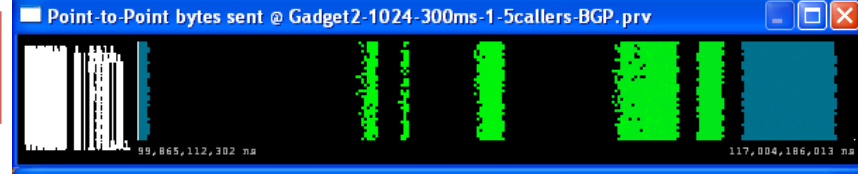
Collective bytes



p2p



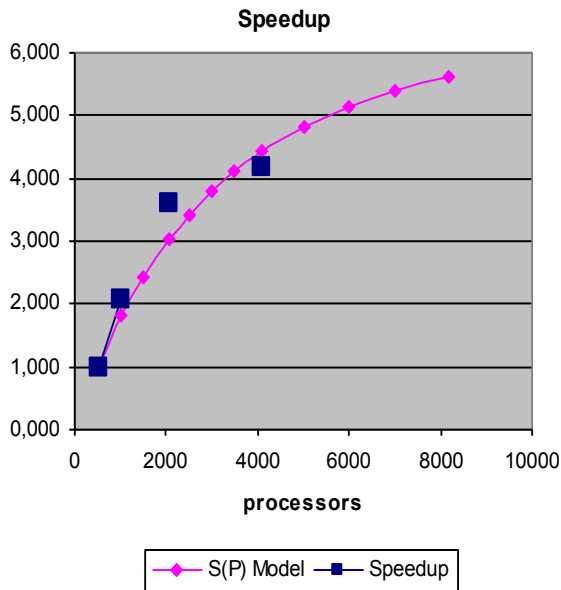
p2p bytes



p2p BW

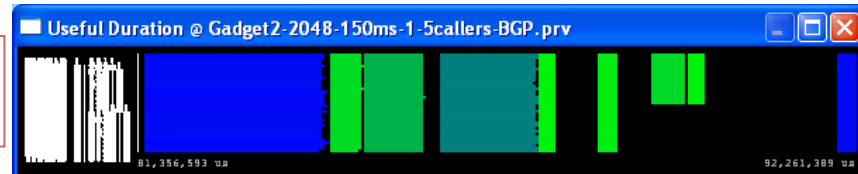


Software counters

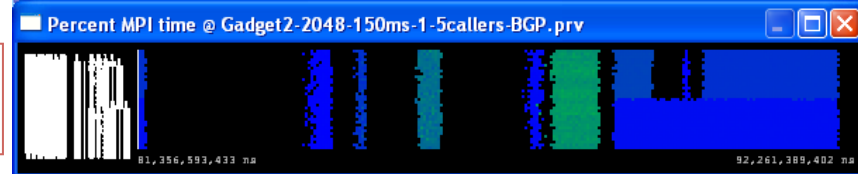


GADGET, PRACE Case A, 2048 procs

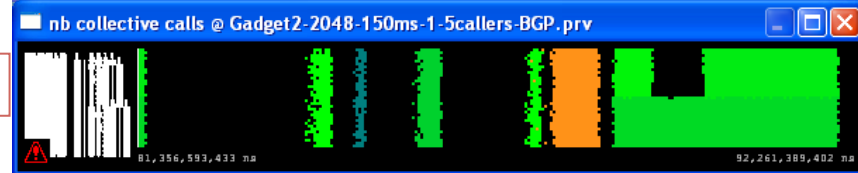
Useful duration



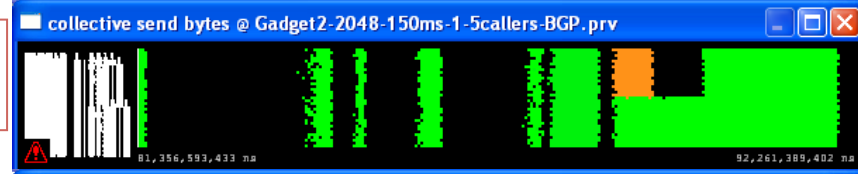
% MPI time



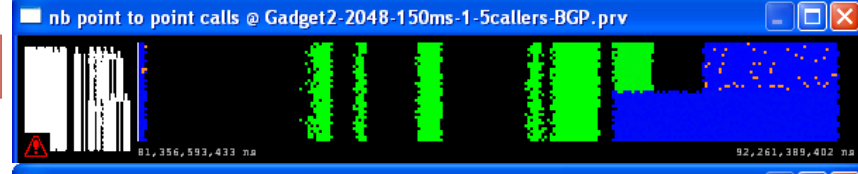
collectives



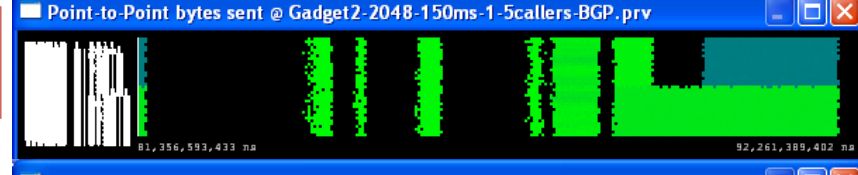
Collective bytes



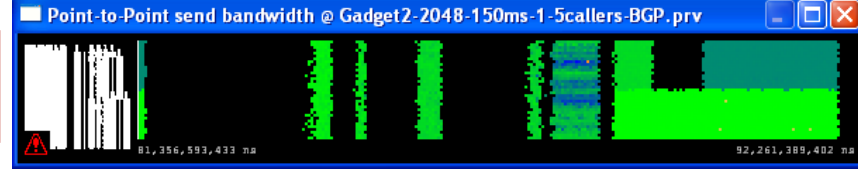
p2p



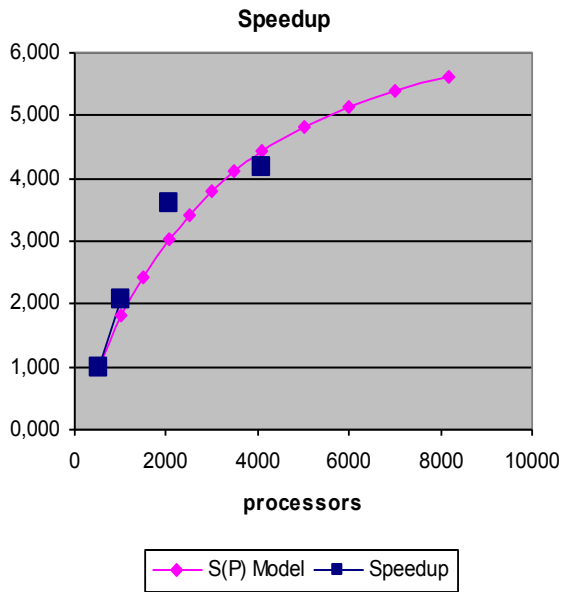
p2p bytes



p2p BW

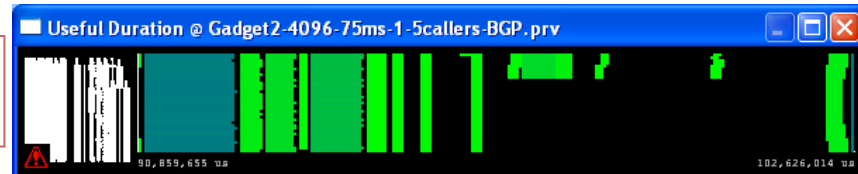


Software counters

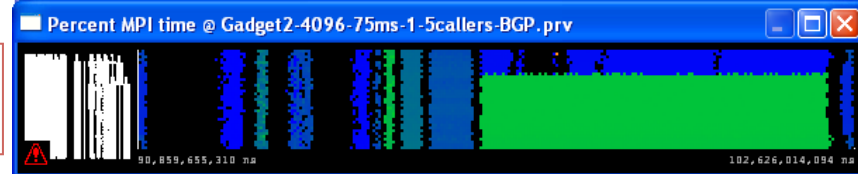


GADGET, PRACE Case A, 4096 procs

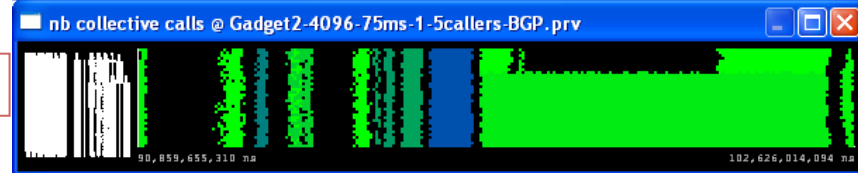
Useful duration



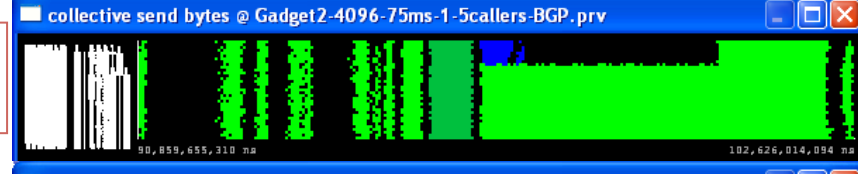
% MPI time



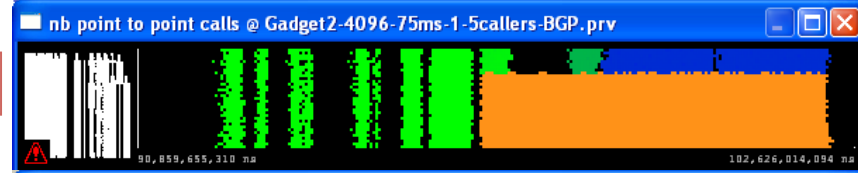
collectives



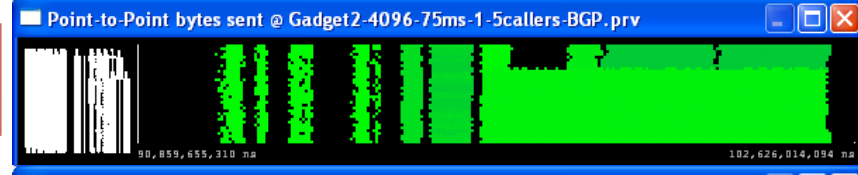
Collective bytes



p2p



p2p bytes



p2p BW

