

New techniques and integration efforts in the CEPBA-Tools environment

Jesus Labarta
Barcelona Supercomputing Center

Index



- General overview
- Recent developments
 - Time analysis
 - Modelling
 - Clustering
 - Sampling
- Tools integration
 - Peekperf
- Interconnection evaluation
- Hierarchical modelling / prediction



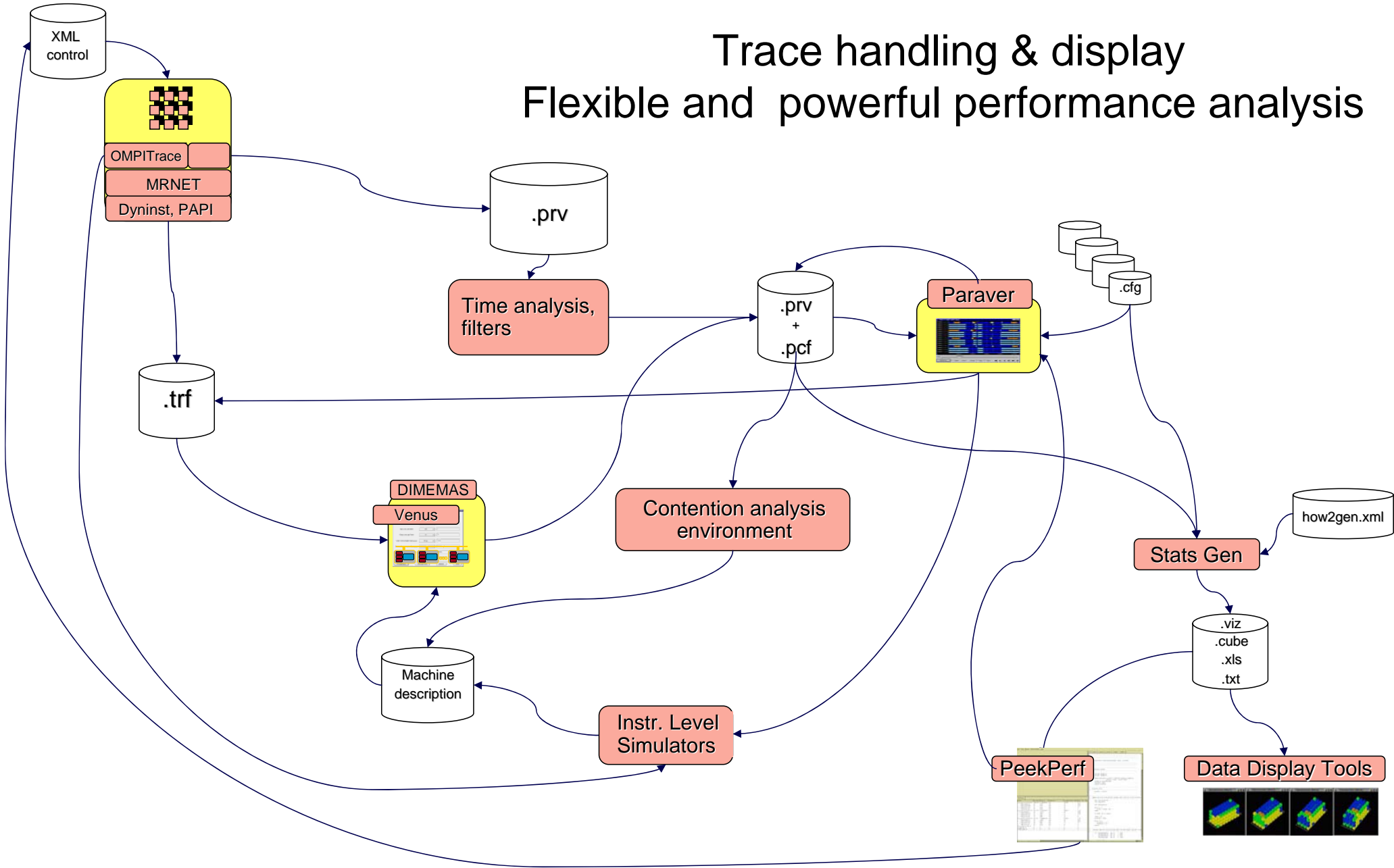


Paraver



Trace handling & display

Flexible and powerful performance analysis



Offer



- Paraver OPEN SOURCE by end of the year
- Dimemas. OPEN SOURCE by end of the year
- Instrumentation OPEN SOURCE by end of the year

- Structure analysis tools in development
 - Signal analysis
 - Clustering
 - Sampling + tracing

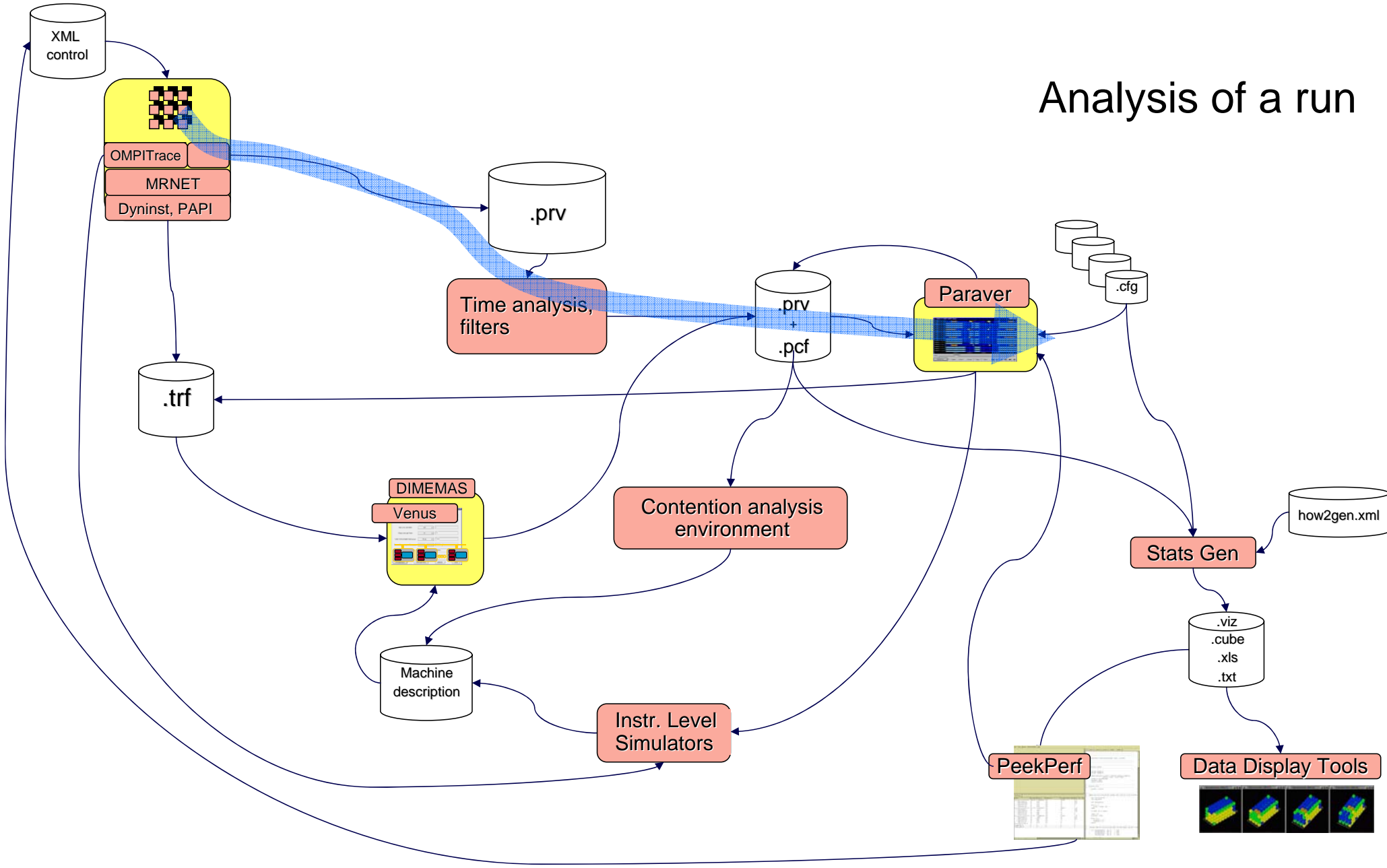
- OpenMP incl 3.0 tasks OPEN SOURCE
- StarSs: CellSs / SMPsSs /GPUSs OPEN SOURCE



CEPBA-Tools Environment

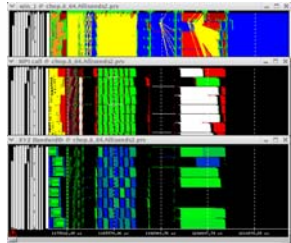


Analysis of a run

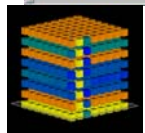


Parallel Program Instrumentation: Platforms

BGL



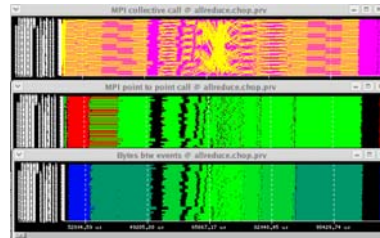
PAPI



Network counters

PERUSE

By Rainer Keller (HLRS) & (UTK)



MPICH collective internals

MPI Calls

MFLOPs

% Vector instr

Avg vector length

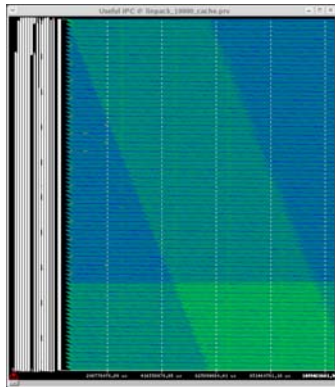
Bank conflicts/us

SX8

By Rainer Keller (HLRS)

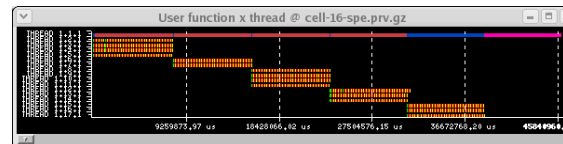


Marenostrum



MPI + OpenMP
Lib. Preload
PAPI

IPC



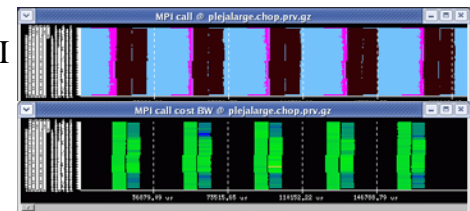
Cell BE

Power 5 AIX

MPI + OpenMP

DPCL

PMAPI



MPI calls

Bandwidth

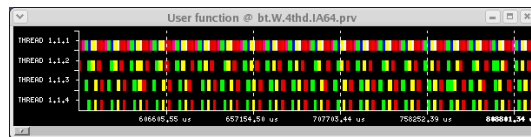
Altix

MPI + OpenMP

Dyninst

PAPI

(Supported by NASA AMES)



(O)MPI-Trace

- XML control specification

```
<trace enabled="yes" home="/gdfs/apps/CEPBATOOLS/64.hwc">  
<mpi enabled="yes">  
  <callers enabled="yes">1-3</callers>  
  <counters enabled="yes" />  
</mpi>
```

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

```
<user-functions enabled="yes">  
  <max-depth enabled="no" />  
  <counters enabled="yes" />  
</user-functions>
```

```
<counters enabled="yes">  
  <cpu enabled="yes" starting-set-distribution="1">  
    <set enabled="yes" domain="all" changeatglobalops="5">  
      PM_CYC,PM_DATA_FROM_MEM,PM_GCT_FULL_CYC,PM_INST_CMPL,PM_INST_DISP,PM_LD_MISS_L1,PM_LD_REF_L1,PM_ST_REF_L1  
    </set>  
    <set enabled="yes" domain="user" changeatglobalops="5">  
      PM_BRQ_FULL_CYC,PM_BR_MPRED_CR,PM_BR_MPRED_TA,PM_CYC,PM_GCT_FULL_CYC,PM_INST_CMPL,PM_INST_DISP,PM_LD_MISS_L1  
    </set>  
  </cpu>  
  <network enabled="yes" />  
  <resource-usage enabled="yes" />  
</counters>
```

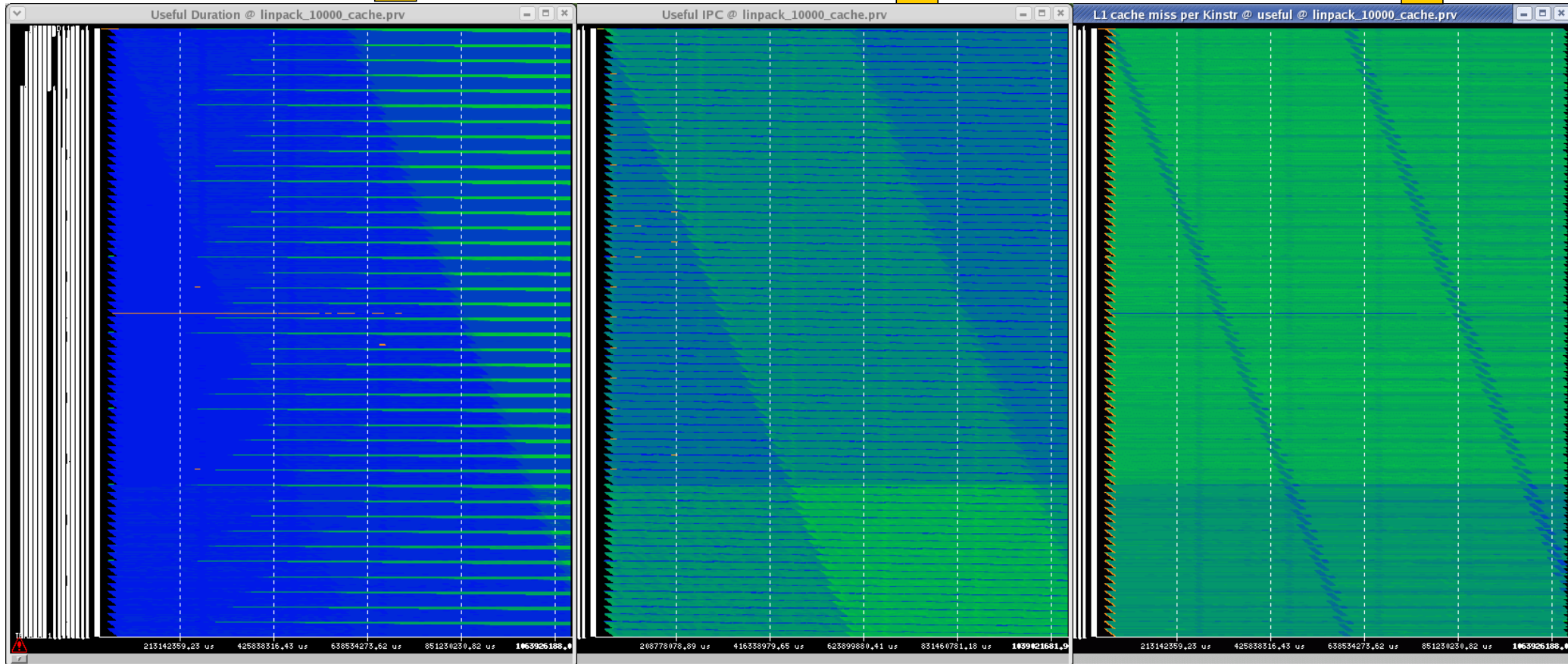
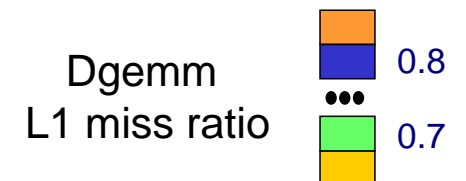
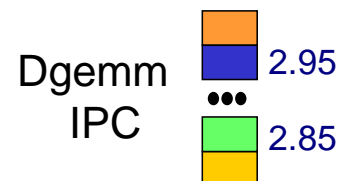
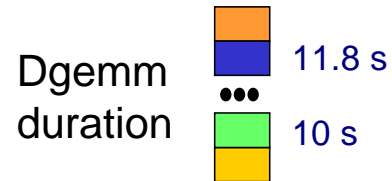
```
<bursts enabled="no">  
  <threshold enabled="yes">500u</threshold>  
  <counters enabled="yes" />  
  <mpi-statistics enabled="yes" />  
</bursts>
```



Scalability of Presentation: timelines



- Linpack @ Marenstrum: 10k cores x 1700 s



Scalability

- A dynamic range issue

- Real causes may be far away from observed effects
- Need to integrate measurement and modeling

MPI calls

MPI comms.

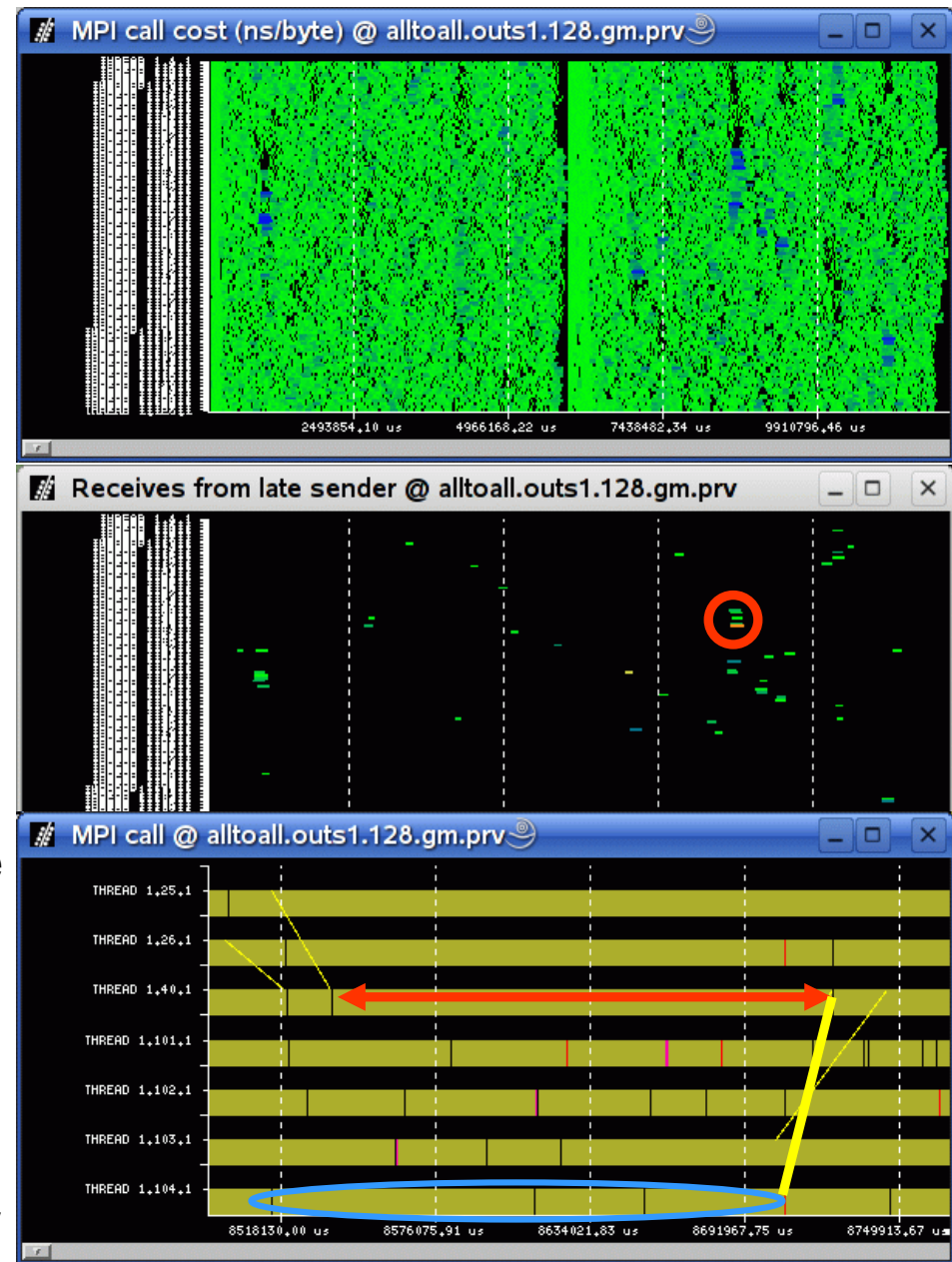
ns/byte

Early receivers

Severe
early receiver

Comms severe
early receiver

Detail severe
early receiver

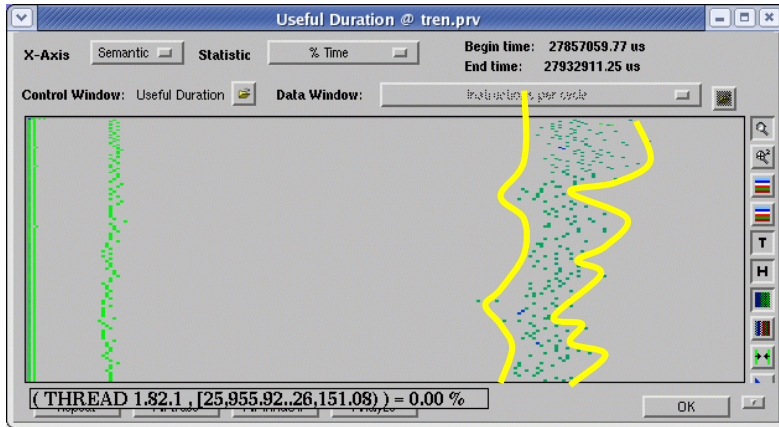


Scalability of Presentation: histograms

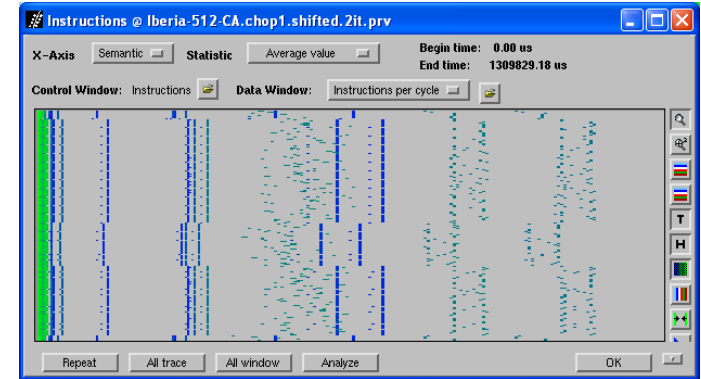


- There is more load imbalance than typically aware/accepted

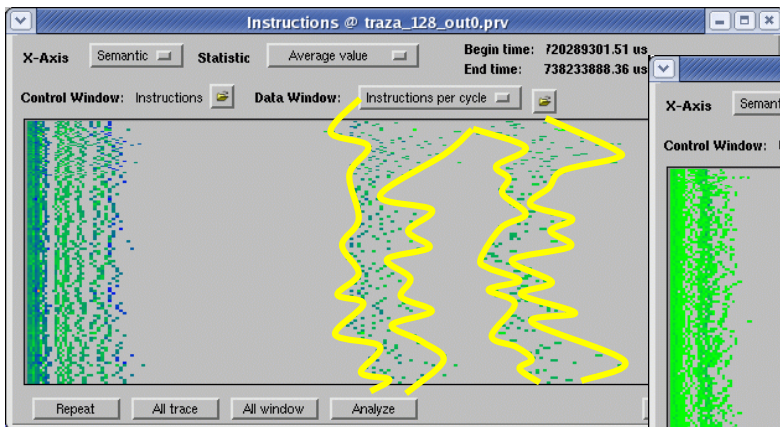
Duration Alya @ 128



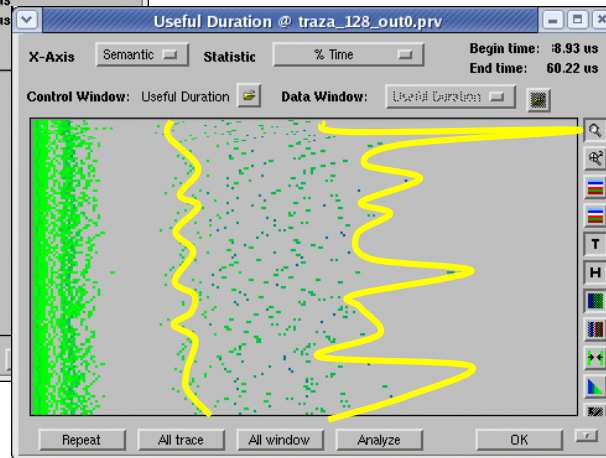
Instructions WRF @ 512 (Iberia12K)



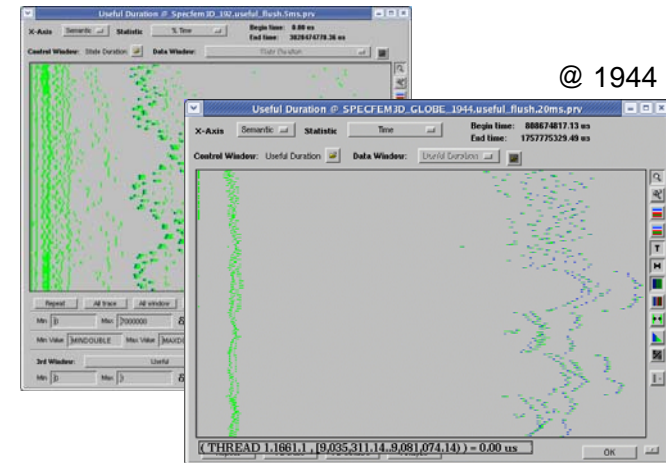
Instructions Airbus @ 128



Duration Airbus @ 128



Duration SPECfEM3D @ 192



@ 1944



Trace handling utilities



- Filter
 - Duration, Type of event, Size, Space.
- Cut
- Merge
- Shift
- Software counters
- Paramedir
- Trace format converters: OTF → Paraver, AIXtrace → Paraver, ...

- Combinations
 - Minimize backwards comms: optimization loop
 - Paramedir
 - Shift
 - Extract a communication phase sub-trace
 - Shift, cut, shift, cut





Signal processing



Signal processing and performance analysis

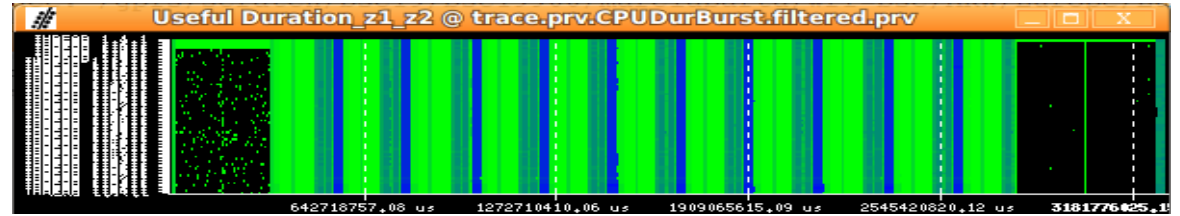
- Signals
 - Flushing processes, %preempted time, #msgs/BW (clogged system)
 - Sum burst duration, #processes in MPI, average IPC,...
- Mathematical morphology
- Spectral analysis
 - Autocorrelation
 - Wavelet transform
- Useful
 - Periodic structure: Reference focus for detailed analysis



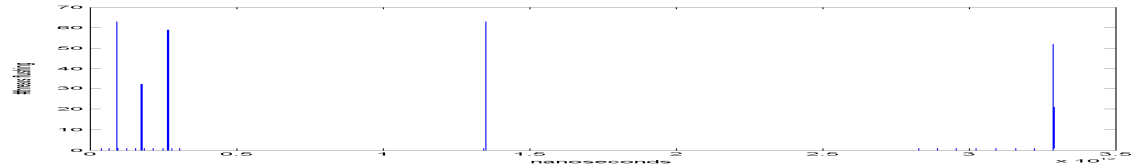
An example



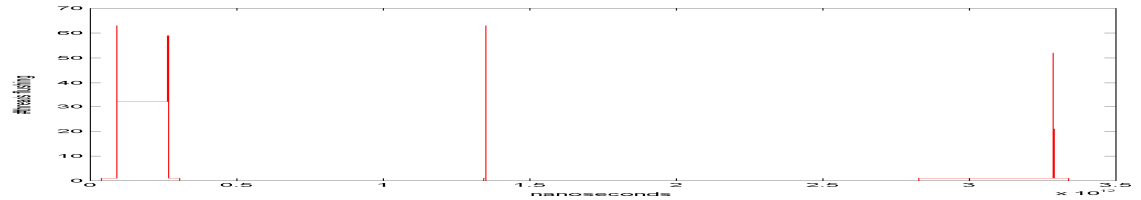
- CPMD



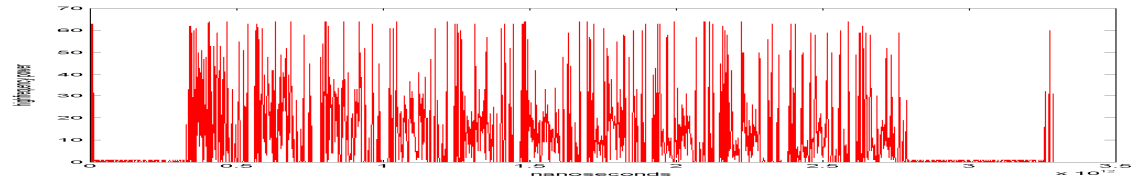
Flushing



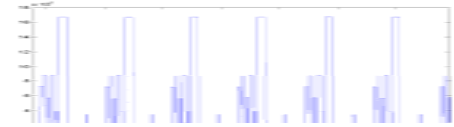
Flushing filtered



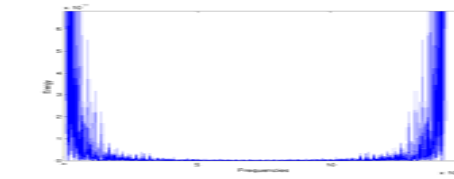
Wavelet
high frequency



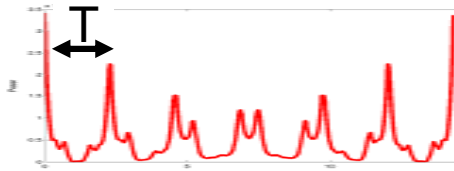
Σ Useful Duration



Spectral density



Autocorrelation

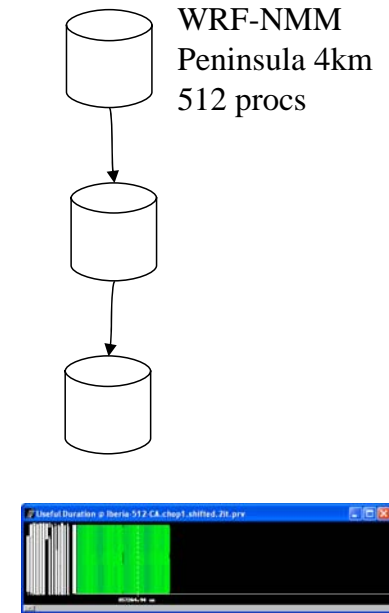
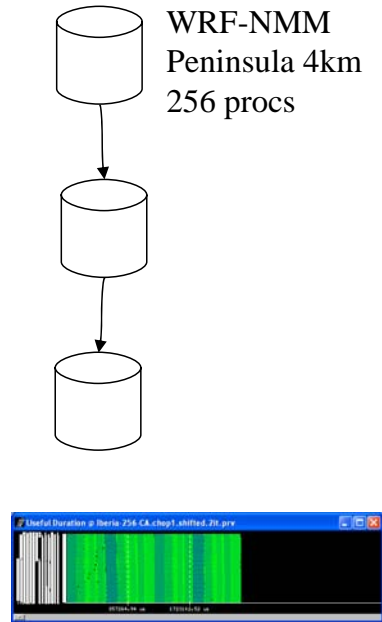
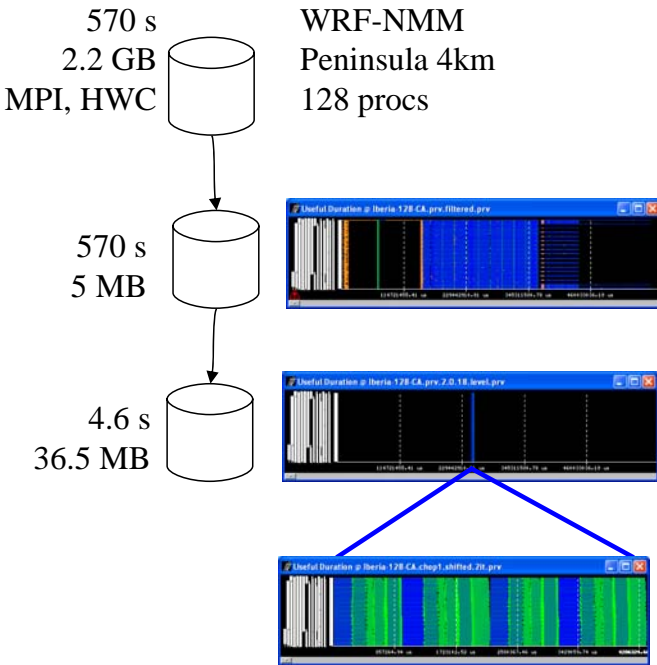




Models



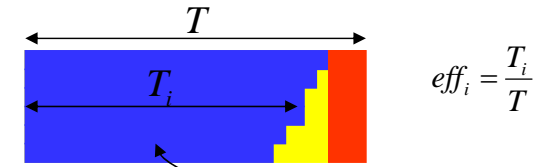
Methodology: automatic speedup analysis



Speedup model

$$Sup = \frac{P}{P_0} * \frac{LB}{LB_0} * \frac{CommEff}{CommEff_0} * \frac{IPC}{IPC_0} * \frac{\#instr_0}{\#instr}$$

Directly from real execution metrics



$$CommEff = \max(eff_i)$$

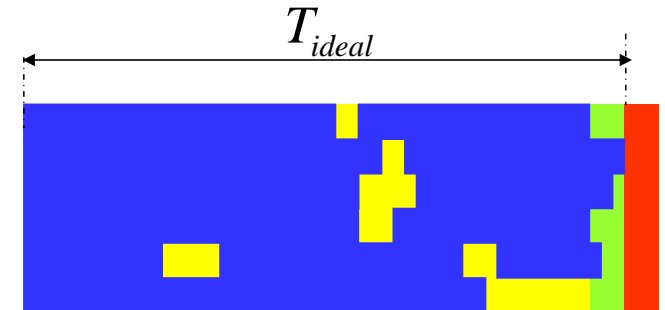
$$LB = \frac{\sum_{i=1}^P eff_i}{P * \max(eff_i)}$$

IPC
#instr

$$Sup = \frac{P}{P_0} * \frac{macroLB}{macroLB_0} * \frac{microLB}{microLB_0} * \frac{CommEff}{CommEff_0} * \frac{IPC}{IPC_0} * \frac{\#instr_0}{\#instr}$$

Requires Dimemas simulation

$$microLB = \frac{\max(T_i)}{T_{ideal}} \quad CommEff = \frac{T_{ideal}}{T}$$



Migrating/local load imbalance
Serialization



Speedup model



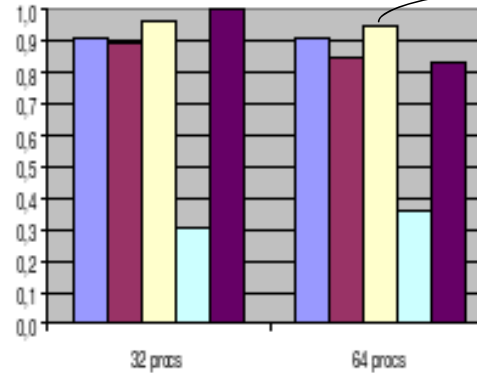
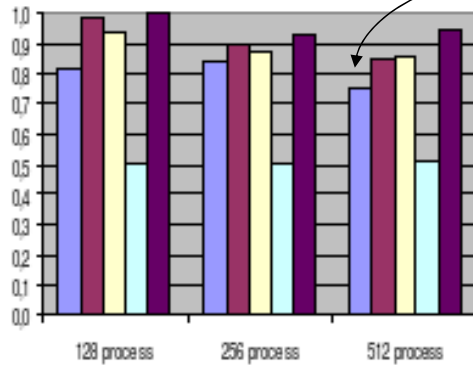
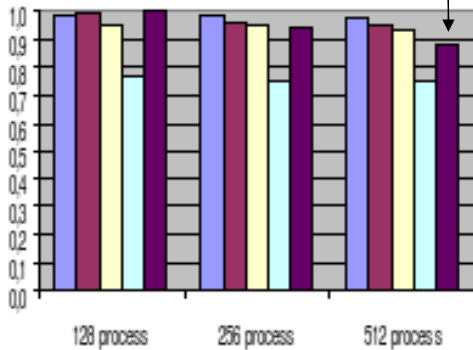
Replication of computation

Poor communication efficiency

WRF-NMM-Iberia

WRF-ARW-Iberia

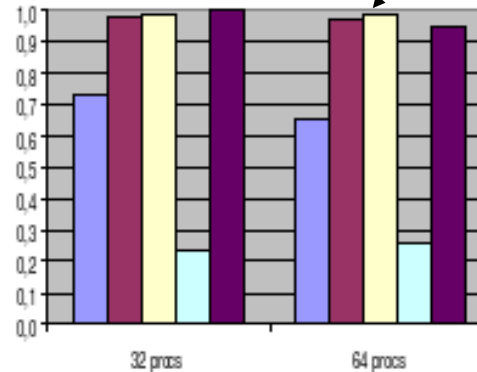
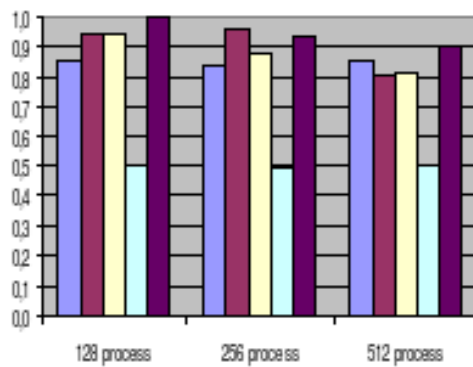
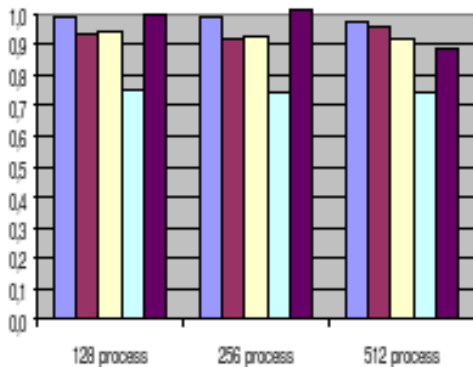
CPMD



WRF-NMM-Europe

WRF-ARW-Europe

CPMD-taskgroups



Improved macro and micro load balance

- Communication
- Micro Load Bal.
- Macro Load Bal
- IPC
- Computation



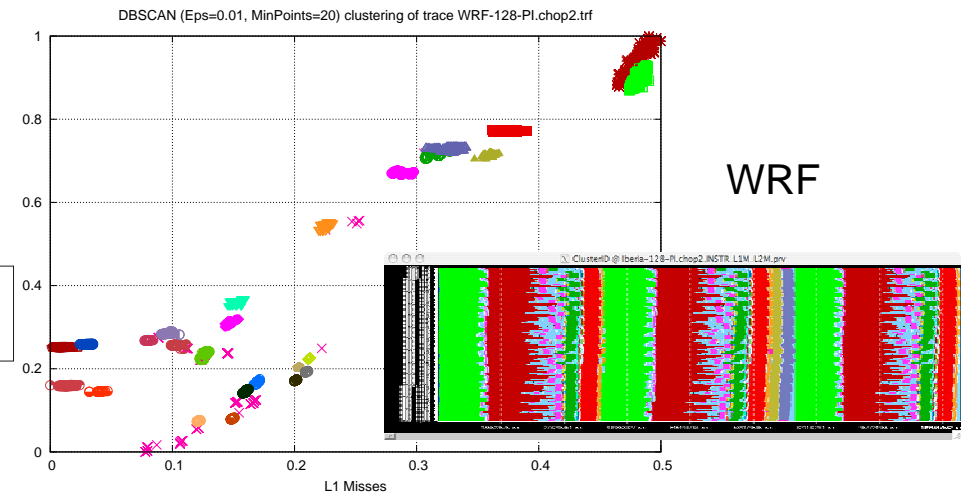
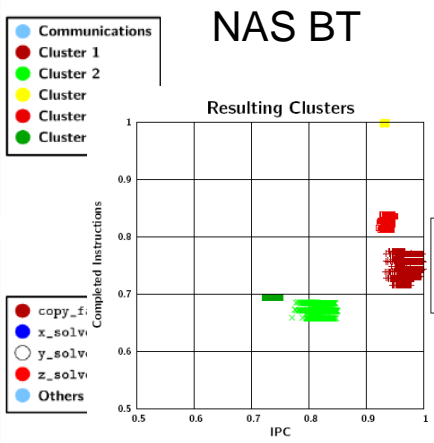
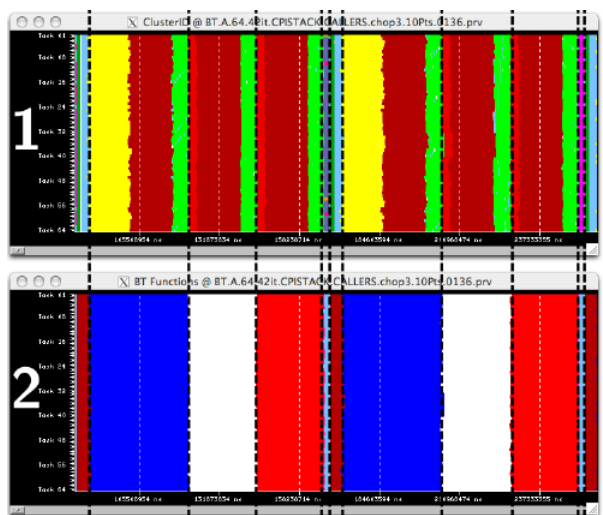
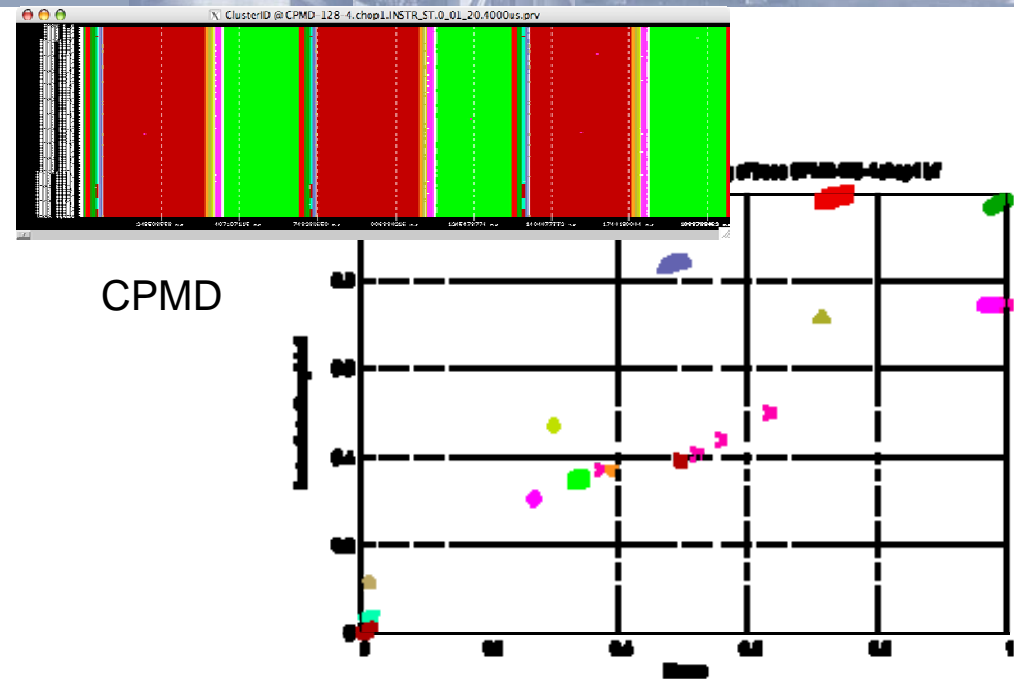


Clustering



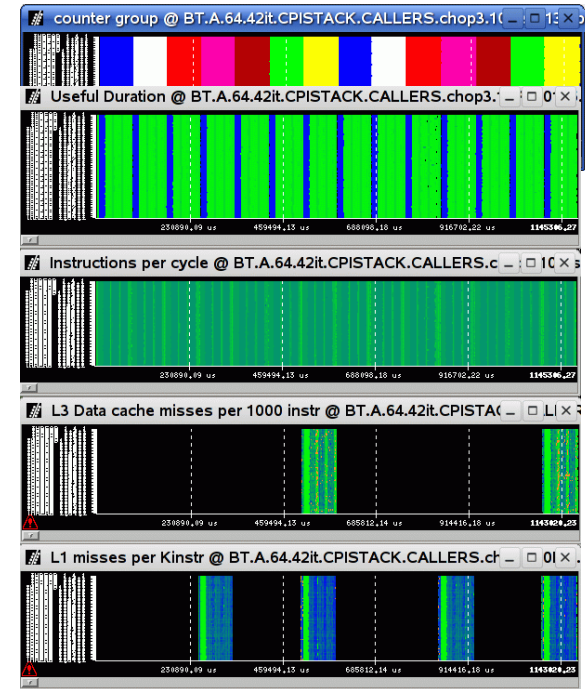
Clustering

- Useful for
 - Identifying and highlighting structure
 - Cluster information injected in trace
 - Phases within routines
 - Different routines may have similar behavior
 - Compact trace encoding
 - Input to time analysis

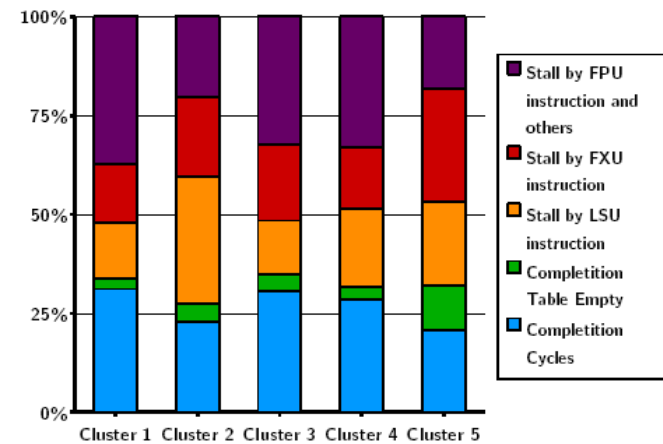


Clustering

- Useful for
 - Precise projection of hardware counters
 - Statistics
 - CPI stack model



CPI Stack Modelization



CLUSTER	1	2	3	4	5
%TIME	54.88	17.96	16.90	6.44	1.42
AVG. BURST DUR. (MS)	1.02	0.78	13.14	2.50	1.11
IPC	1.02	0.65	0.89	0.91	0.53
MIPS	2231.8	1423.3	1966.5	2001.8	1163.0
MFLOPS	339.2	46.3	191.6	269.2	23.6
L1M/KINSTR	0.92	1.53	1.19	1.17	2.88
L2M/KINSTR	0.06	1.26	0.06	0.35	0.21
MEM.BW (MB/s)	16.79	218.47	13.87	85.77	29.76





Sampling +



Sampling

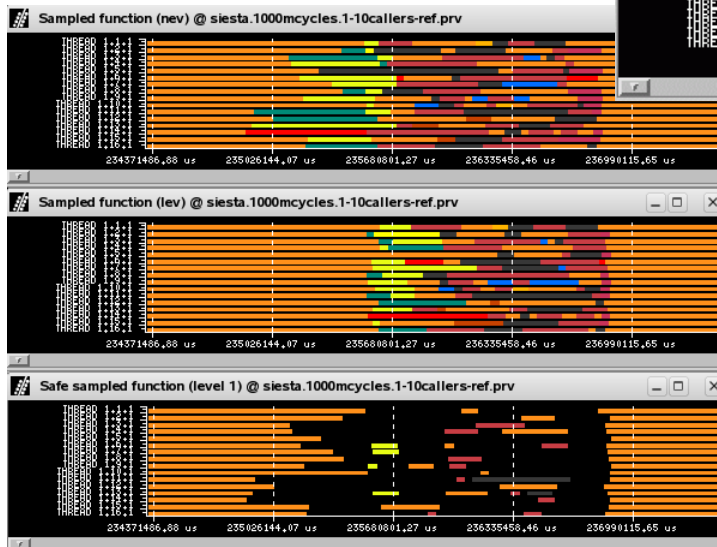
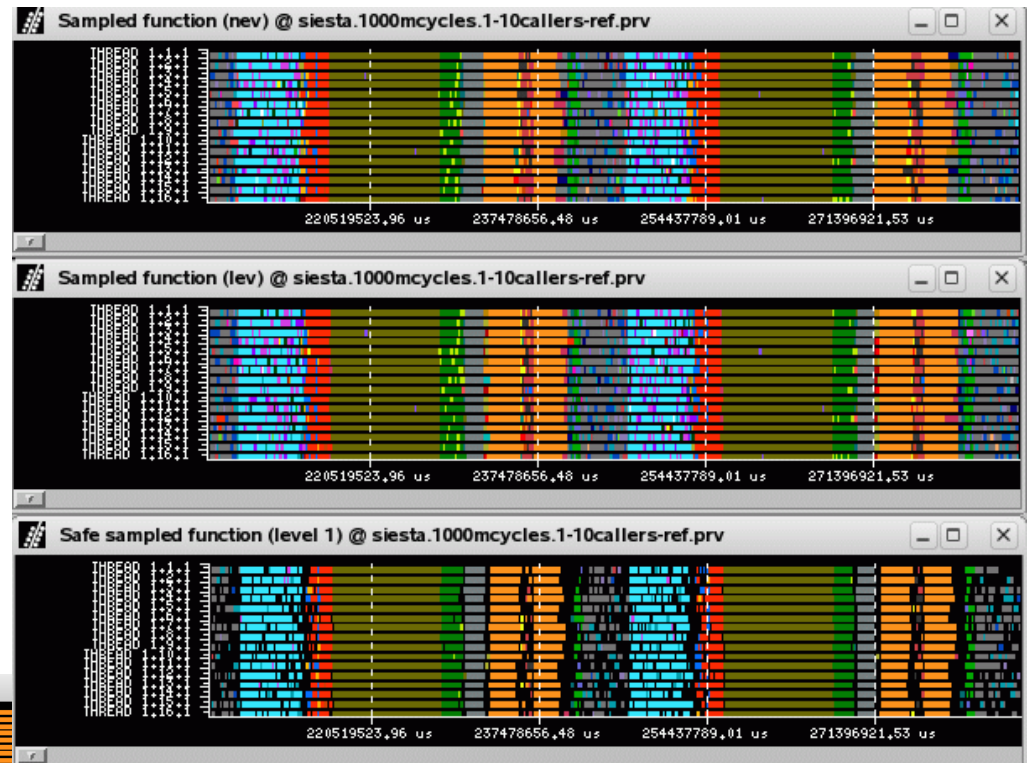


- Tesis:
 - Has much more potential than currently used
 - Mixed with: instrumentation, frequency analysis, clustering
- Example questions
 - Can we get Information of time distribution (i.e. same length of all instances)?
 - Can we get very fine grain information with little overhead?
- The issue: How to project information captured by sampling
- Periodic Sampling (cycles counter)
 - High frequency ($>$ Nyquist)
 - Low frequency ($<$ Nyquist)
- “non periodic” sampling
 - Correlated to specific architectural/application hardware counters



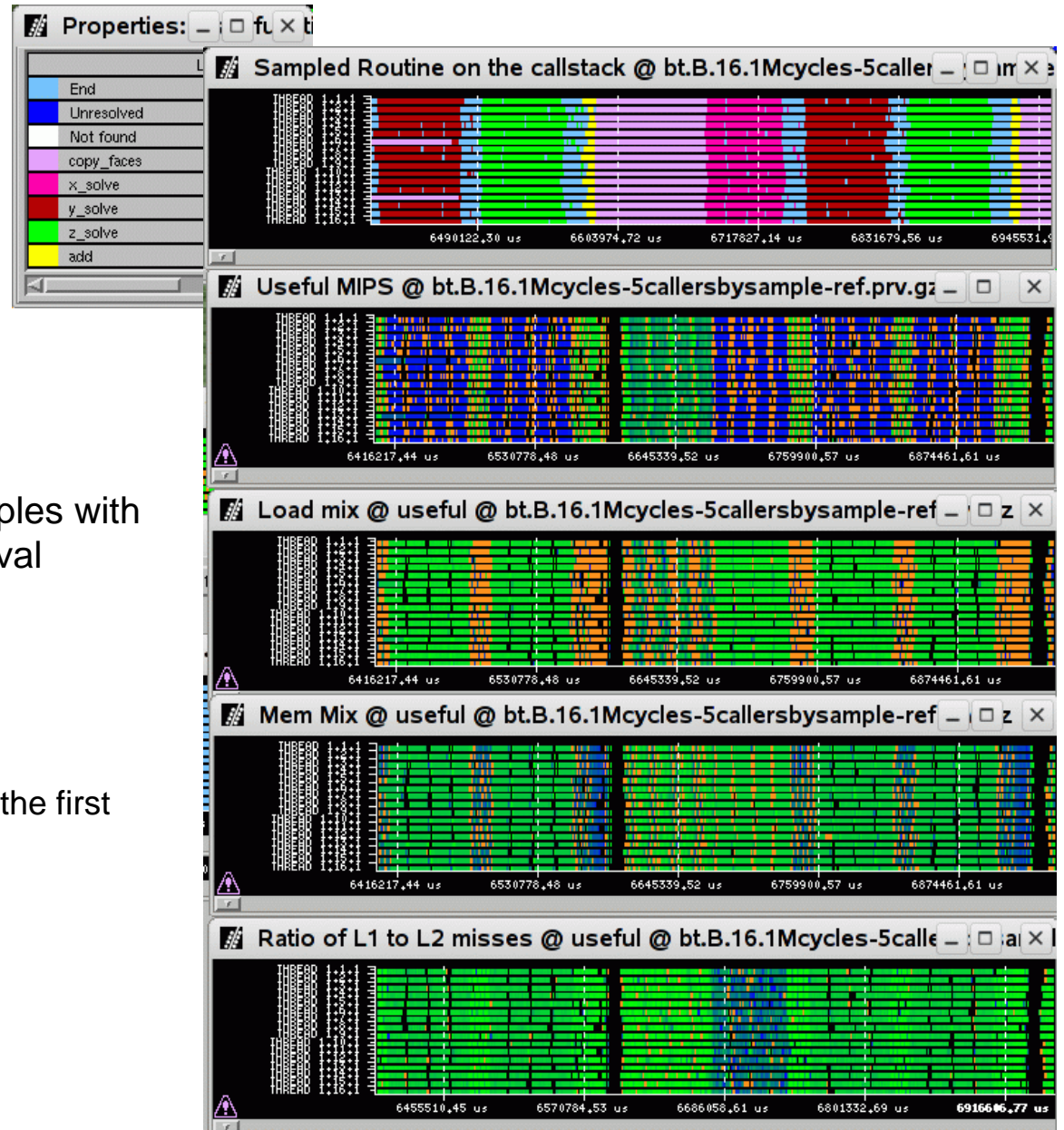
High frequency periodic sampling

- Sampling frequency $>$ Nyquist
- Function within interval?
 - Event at beginning
 - Event at end
 - “safe” when both are the same



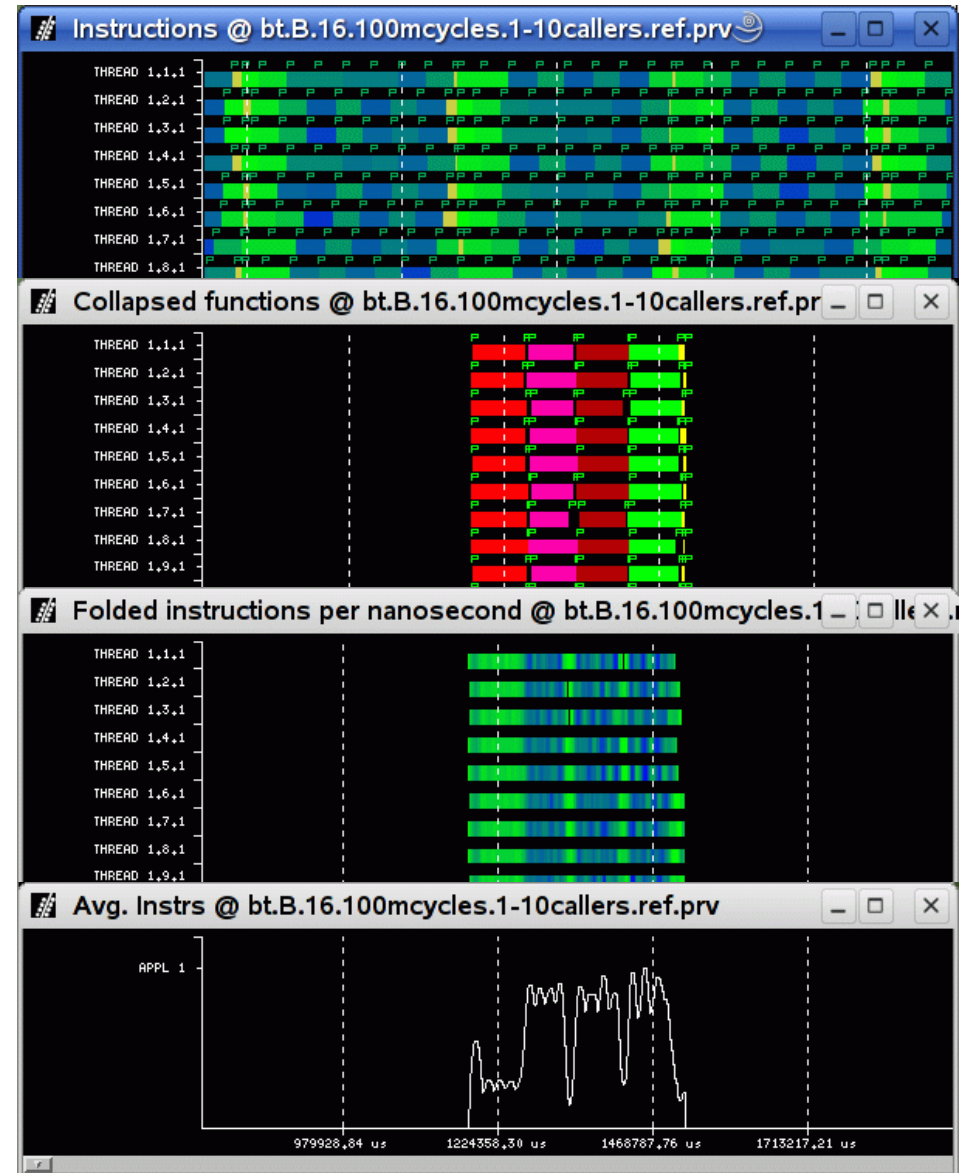
High frequency periodic sampling

- Sampling frequency > Nyquist
- Shows fine structure
- Identification of function span
 - Assumption: consecutive samples with same call stack \rightarrow whole interval assigned to function
 - All functions / specific subsets
 - Top of stack
 - Walk the stack searching for the first routine in the target set.



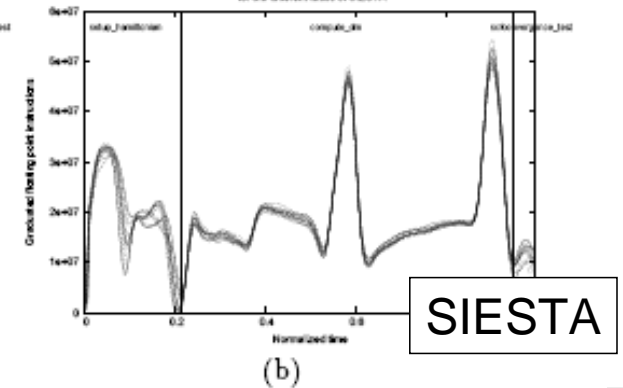
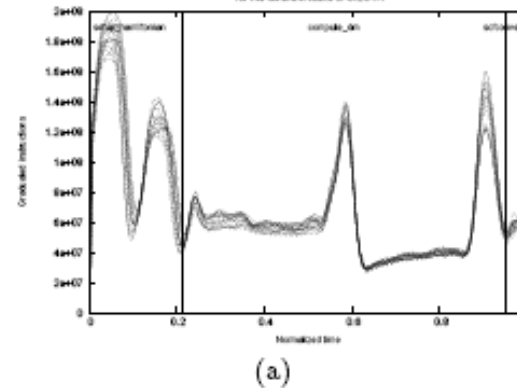
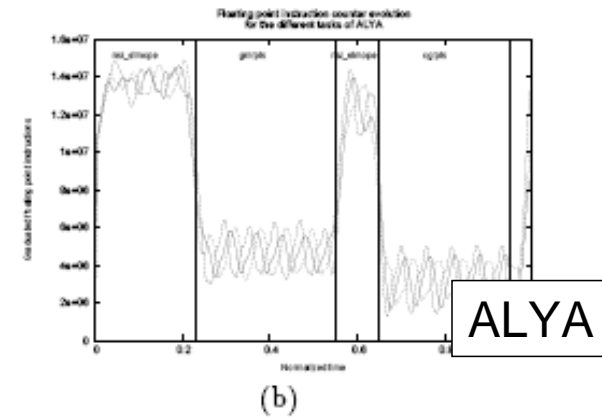
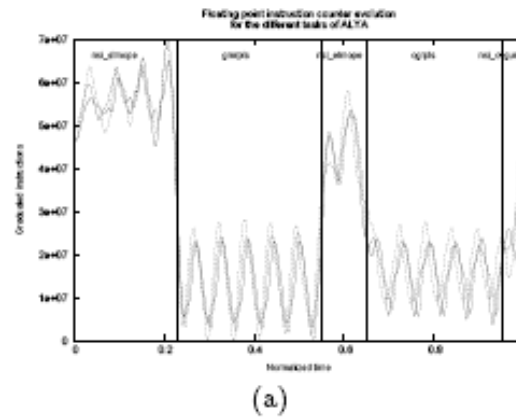
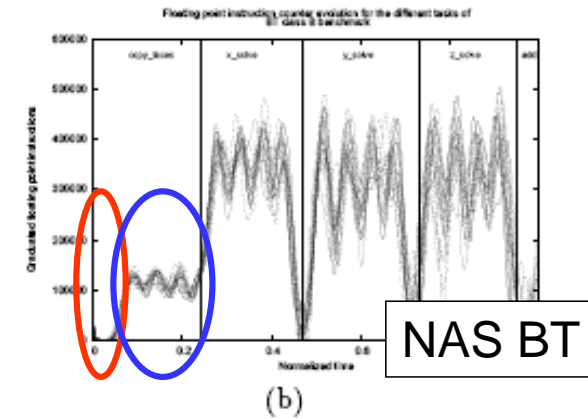
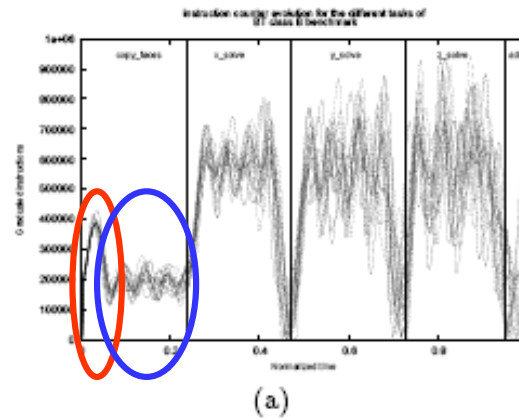
Low frequency periodic sampling

- How to increase precision?
- Folding based on known periodic structure of application (tagged iterations)
- Applies to stationary applications
- Result: trace for one iteration with synthetic paraver events
- Refer counts/timestamps to start of iteration
 - Call stack
 - Search for consecutive sequences of folded samples within same function and generate synthetic events
 - Hardware counters
 - Noise reduction
 - Fit folded samples
 - Sample fitting curve to generate synthetic events.



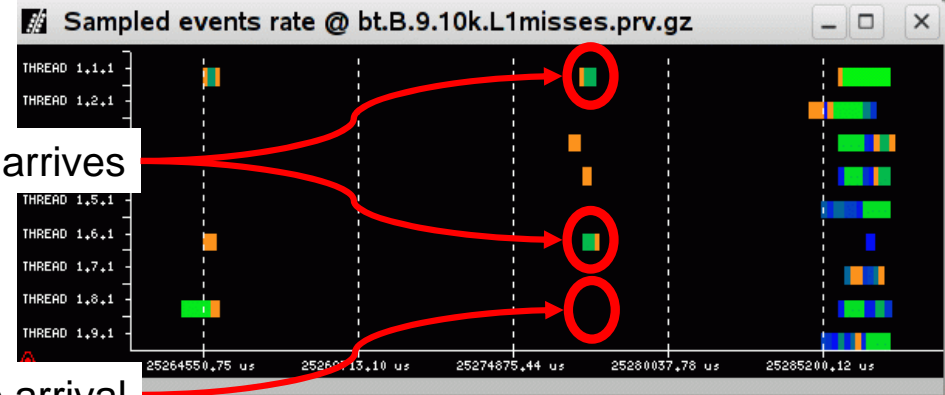
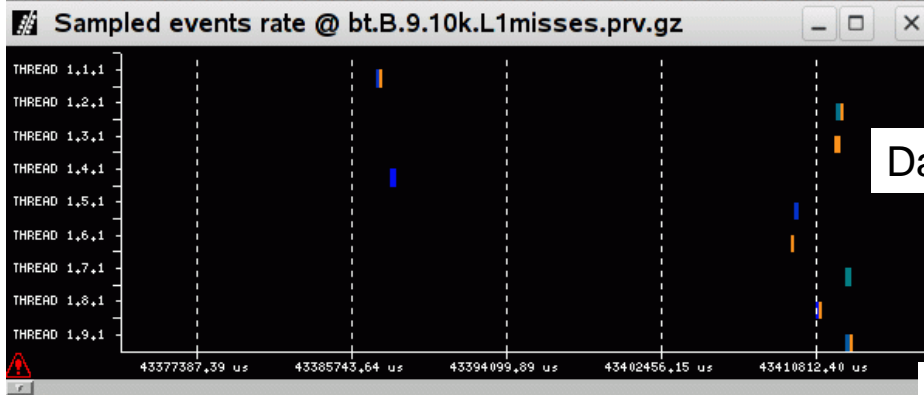
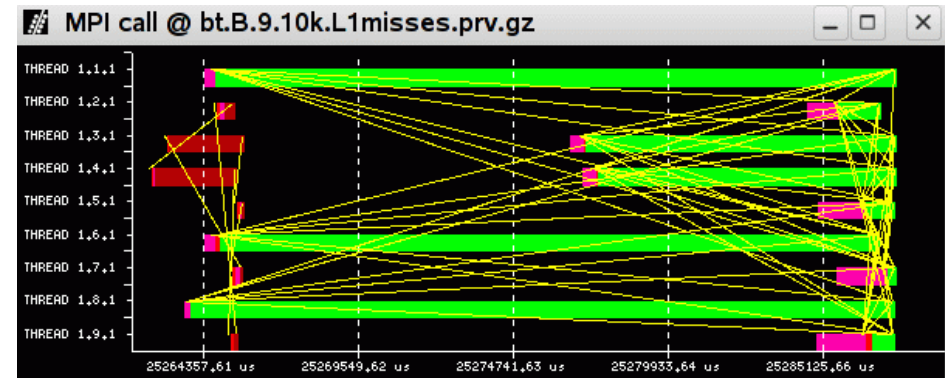
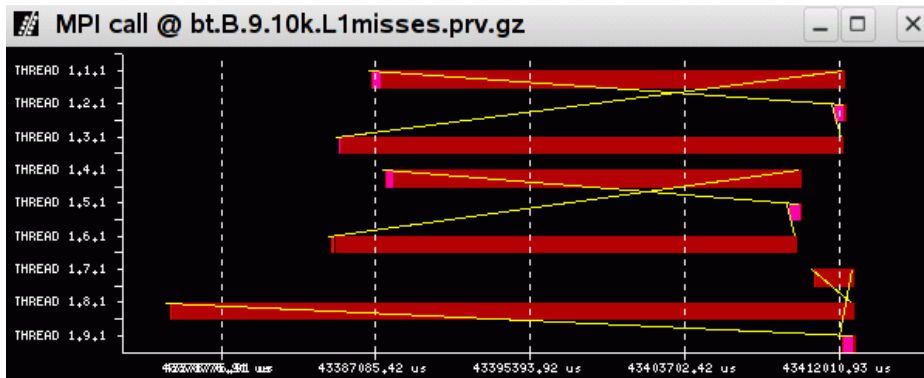


- Low pass filter ...
- but still possible to identify loops
- Scalable summarization ...
- but still high level of detail



Non periodic sampling

- Correlated with specific architectural/program feature
- Useful to identify
 - Internal behavior
 - Density of L1 misses within MPI in an SMP: when data actually arrives.



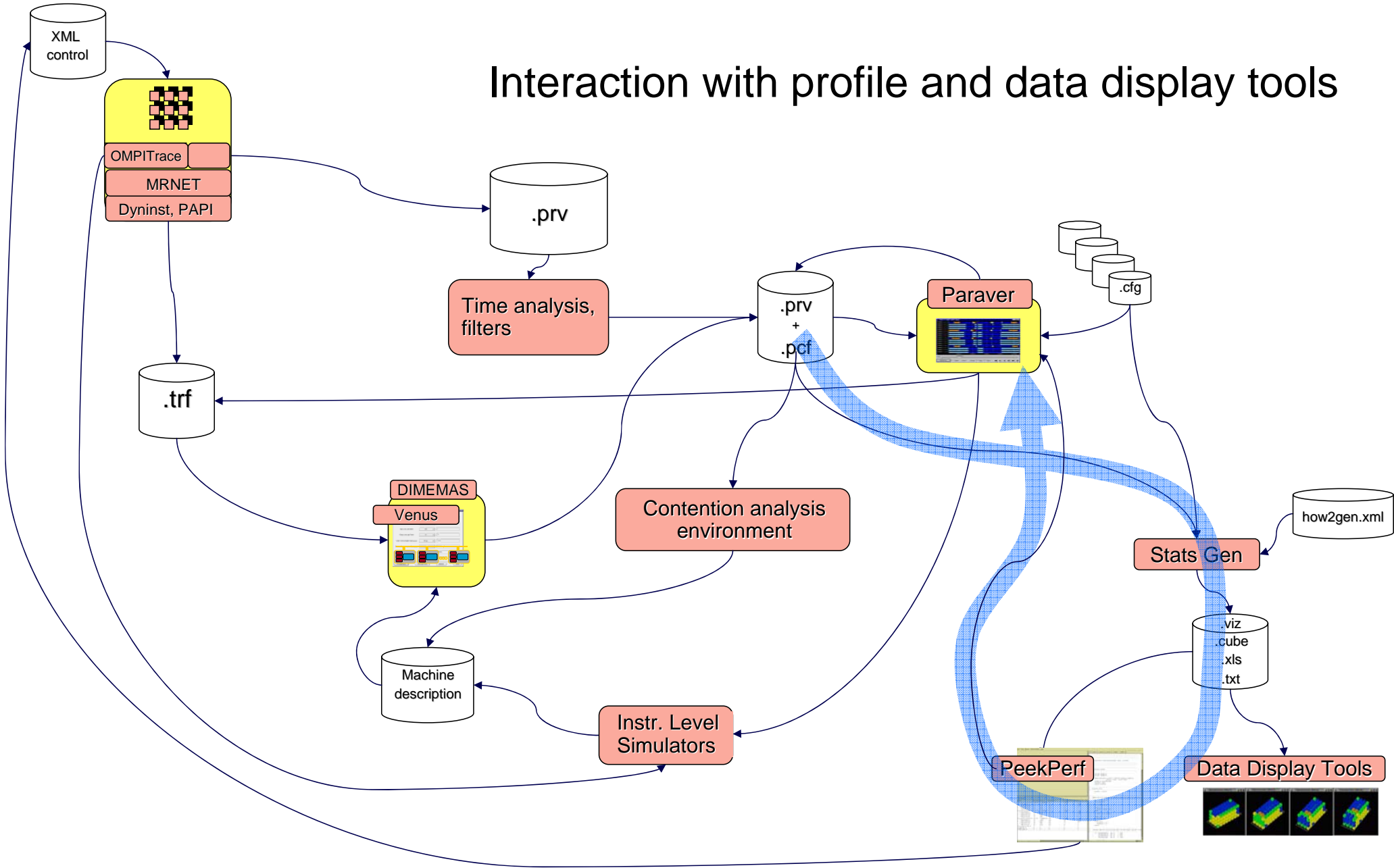


Profiles -- Timelines



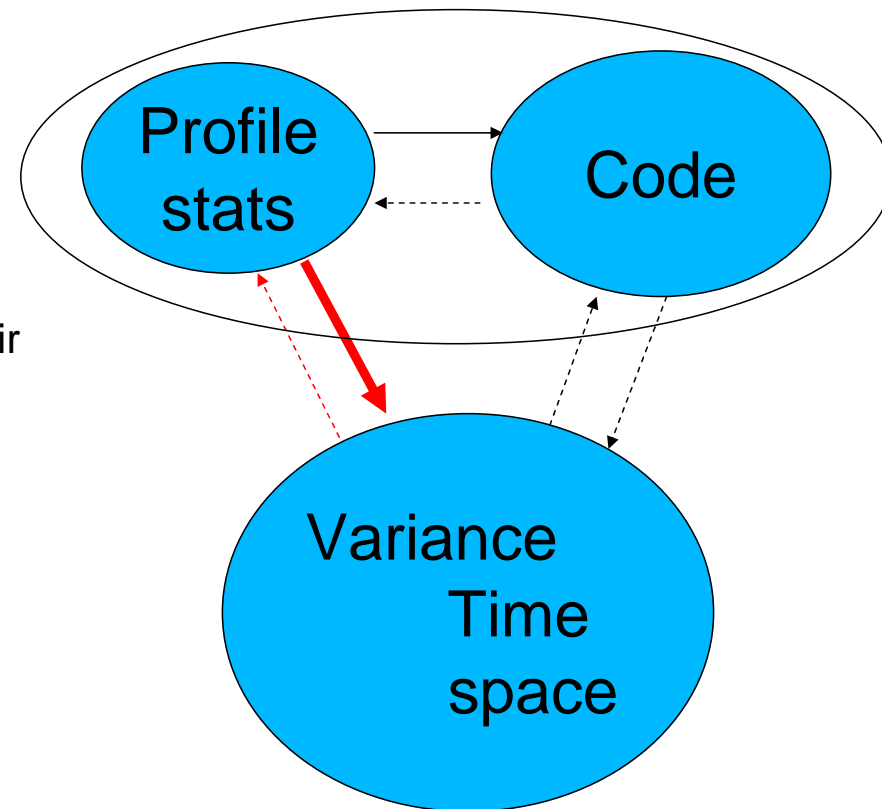
CEPBA-Tools Environment

Interaction with profile and data display tools



Integration with profile presentation tools

- Profile presentation tools
 - Reduction/aggregation of the performance dimensions
 - Time dimension disappeared / Space dimension sometimes
- Traces
 - “All” data is there
- Library
 - API to access statistics computed with Paramedir
- Translators
 - Configuration: Tool building environment
 - Prototypes
 - Peekperf prototype, Gprof like, ...
 - Starting
 - Paraprof, CUBE



Integration Paraver- Peekperf

- Statistics generator
 - Built a set of C++ libraries to access **paramedir** stats.

The screenshot displays the Paraver performance analysis tool interface. It consists of three main windows:

- Useful Duration @ NMM-Eur-12km-512-CA,shifted.chop1.1.prv**: A heatmap showing execution duration over time. The x-axis represents time in microseconds, with markers at 245098603.18 us, 245264209.77 us, 245429816.36 us, 245595422.95 us, and 245761029.54 us.
- Main Window - [DATA VISUALIZATION WINDOW]**: A table of performance metrics for various operations. The table has columns for Label, Time, # Bursts, IPC, #Instr, Cycles, L2 Data cache miss, Stores/ki, p2p. Bandwidth Mbytes/s, and p2p. Size (By).

Label	Time	# Bursts	IPC	#Instr	Cycles	L2 Data cache miss	Stores/ki	p2p. Bandwidth Mbytes/s	p2p. Size (By)
copy_faces!	1.53864e+07	202	0.44	1.77854e+10	1.00443e+11	8.80208e+07	318.47	1.24247e+07	9.9287e+07
error!	2908.93	1	0.59	3.91185e+07	1.75352e+08	142020	142.92	0	0
x_solve.f	4.23072e+07	201	0.43	3.5619e+10	1.76931e+11	9.2785e+07	351.77	1.35963e+06	5.25334e+07
x_receive_backsub_info(x_solve.f)	2442.72	201	0.23	5.96112	4.97582e+06	3659	212.69	0	0
x_receive_solve_info(x_solve.f)	3244.63	201	0.13	3.25657e+09	2.64391e+10	360771	215.65	0	0
x_send_backsub_info(x_solve.f)	50793.5	201	0.29	6.1894e+09	2.16842e+10	2.56367e+07	230.67	317142	8.75556e+06
x_send_solve_info(x_solve.f)	364870	201	0.22	1.67362e+10	7.59547e+10	3.51771e+07	211.16	1.35963e+06	5.25334e+07
x_solve(x_solve.f)	4.23072e+07	201	0.43	3.5619e+10	1.76931e+11	9.2785e+07	351.77	112983	5.25334e+07
Compute_x_solve	4.23072e+07	201	0.21	3.5619e+10	1.76931e+11	9.2785e+07	218.94	0	0
MPI_Wait_70	210146	201	0.28	1.32031e+08	4.71288e+08	85687	217.32	0	0
MPI_Wait_71	147992	201	0.06	1.91703e+07	3.45736e+08	12351	351.77	72813.4	5.25334e+07
MPI_Wait_98	17740.8	201	0.43	1.71969e+07	4.06413e+07	5784	245.4	0	0
MPI_Wait_99	15903.5	201	0.09	3.56576e+06	3.79221e+07	2757	334.17	112983	8.75556e+06
y_solve.f	4.6181e+07	201	0.36	2.83117e+10	1.68343e+11	1.15226e+08	351.98	1.1555e+06	5.25334e+07
z_solve.f	5.53727e+07	201	0.61	1.97999e+10	1.47329e+11	1.57546e+08	351.98	1.3265e+06	5.25334e+07
- Useful Instructions_3DZoom_ranges[58200000.00,88800000.00][0.80,1.00] ...**: A heatmap showing instruction counts over time, with the same x-axis markers as the Useful Duration window.
- SOURCE CODE WINDOW**: Displays the source code for the `x_solve.f` file. The code includes MPI-related operations and calls to `x_solve_cell`.


```

bt.f  x_solve.f  y_solve.f  z_solve.f
61      first = 0
62      call x_receive_solve_info(recv_id,c)
63      _____
64      c  overlap computations and communications
65      _____
66      call lhsx(c)
67      _____
68      c  wait for completion
69      _____
70      call mpi_wait(send_id,r_status,error)
71      call mpi_wait(recv_id,r_status,error)
72      _____
73      c  install C'(istart) and rhs'(istart) to be used in this cell
74      _____
75      call x_unpack_solve_info(c)
76      call x_solve_cell(first,last,c)
77      endif
78
            
```





Tracing???



Is it all about traces?



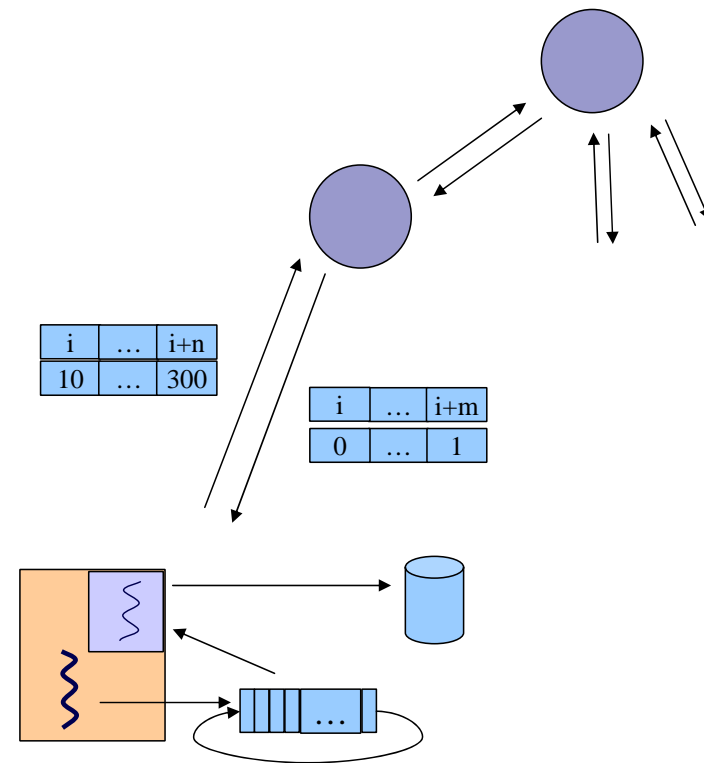
- Traces
 - See - colors
 - Scalable view of metrics and variance
 - Detailed analysis
 - Experience - Understanding - Insight
 - Ideas

- Then
 - Put the **intelligence** in online environments
 - Use scalable infrastructures



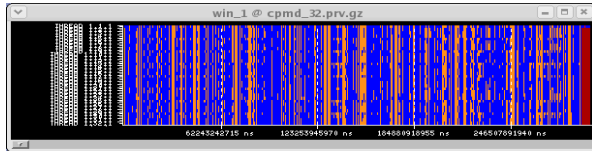
Distributed trace control

- MRNet based mechanism
 - Local instrumentation on a circular buffer
 - Periodic MRNet front-end initiation of collection process
 - Local algorithm
 - Reduction on tree
 - Selection at root propagated
 - Locally emit trace events

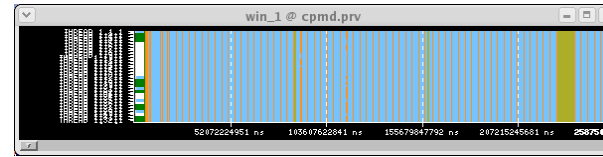


Online analysis

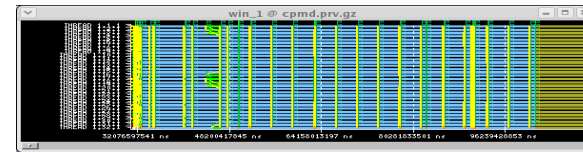
- Collective duration threshold



245MB, >15500 col

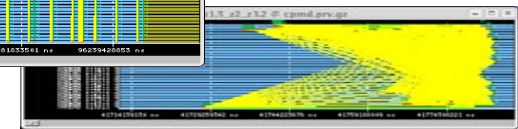


<1MB, <85 col

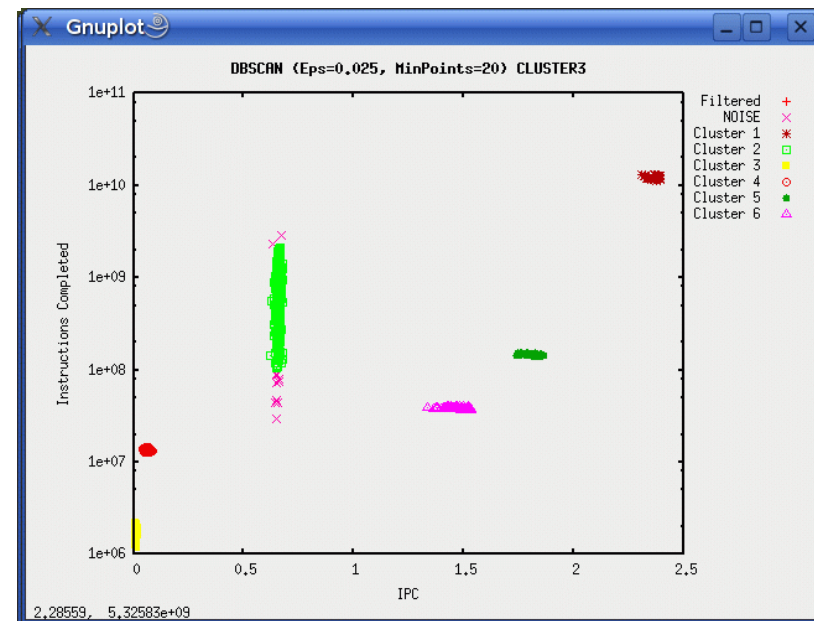


25MB, <85 col

Collective internals



- Periodic clustering snapshots
- Periodic frequency analysis
- Maximize information/data ratio
 - Direct report metrics,
 - Activate tracing (user/automatic)
 - Focused traces





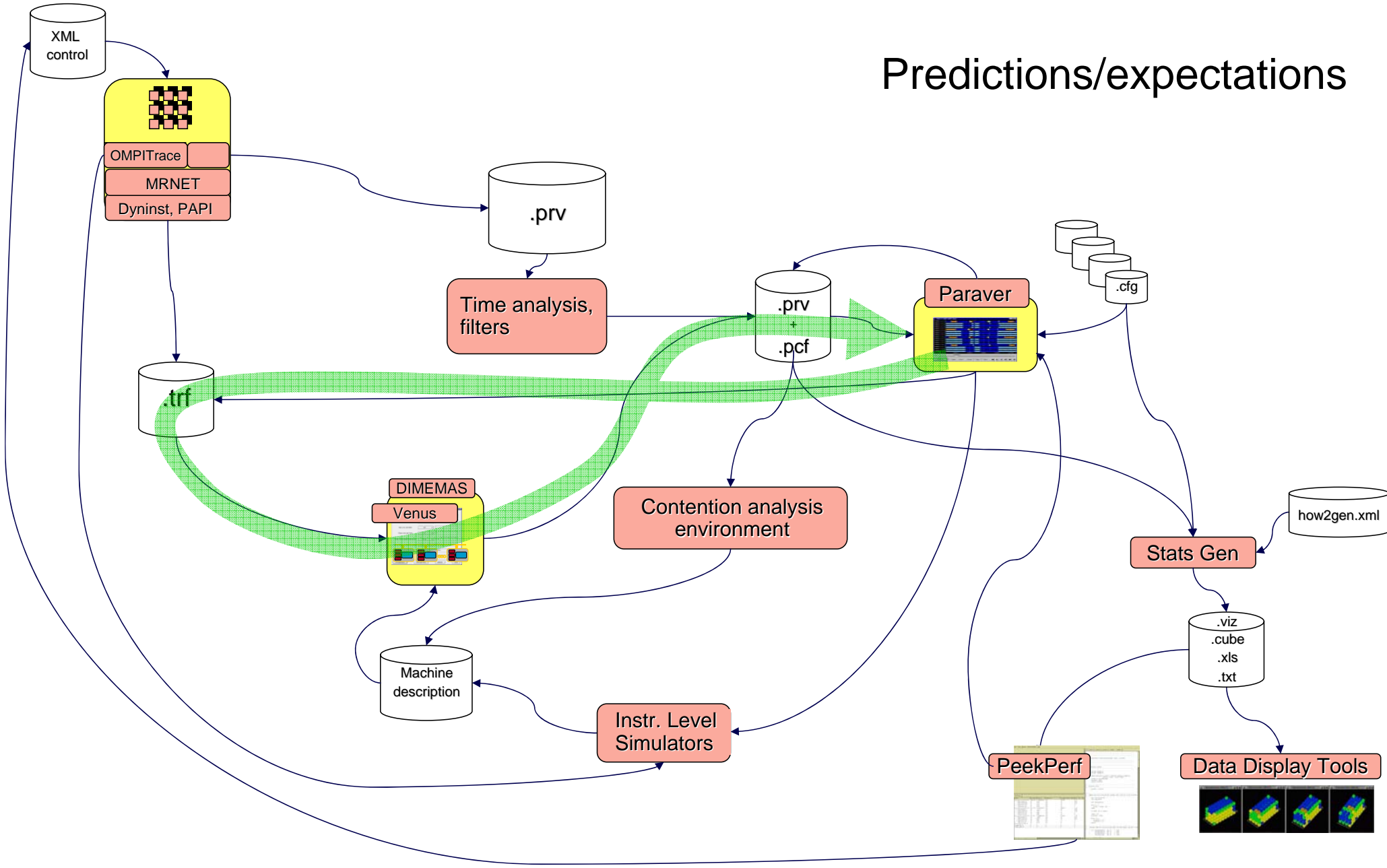
Dimemas



CEPBA-Tools Environment

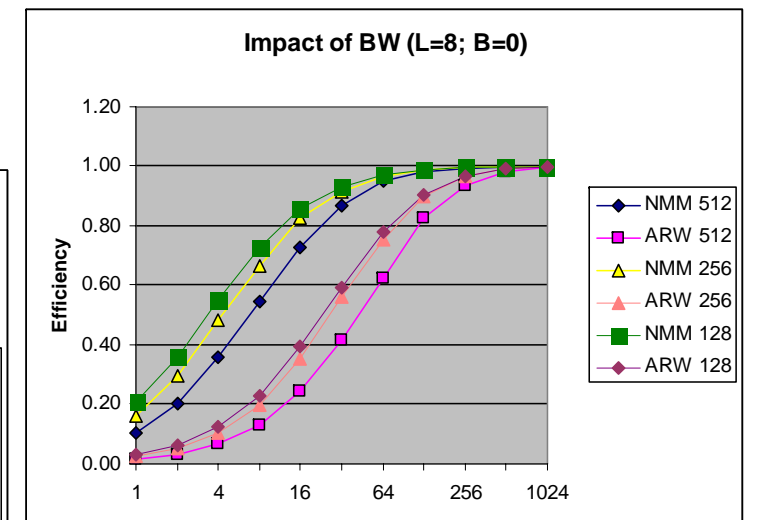
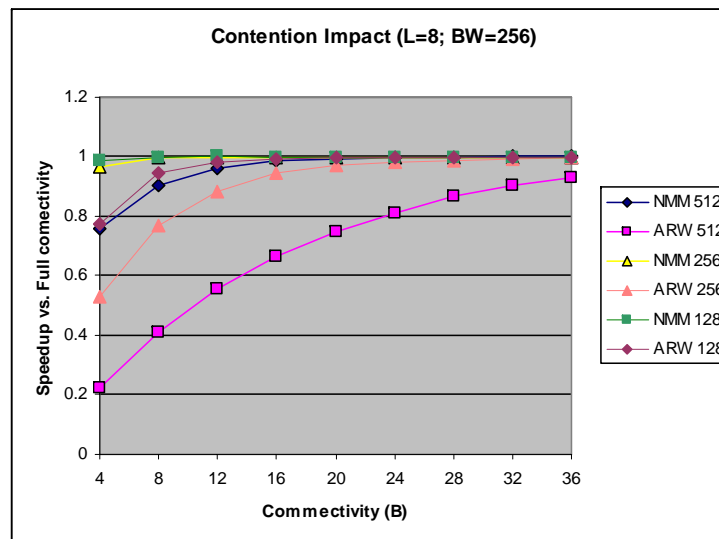
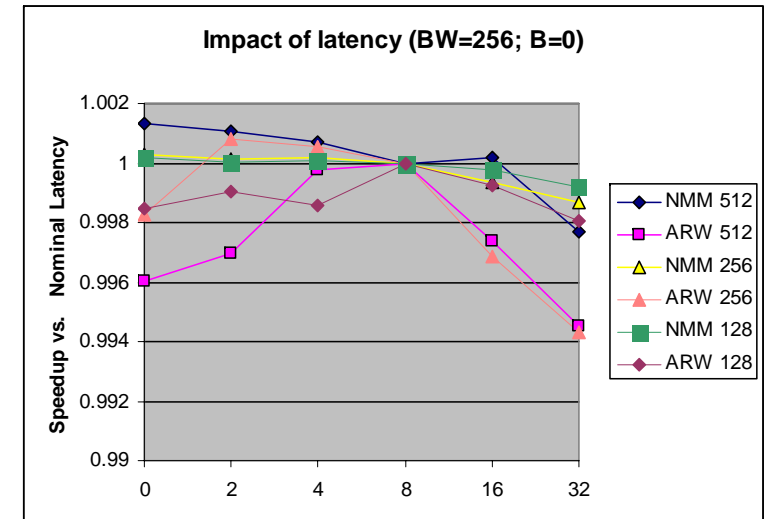


Predictions/expectations



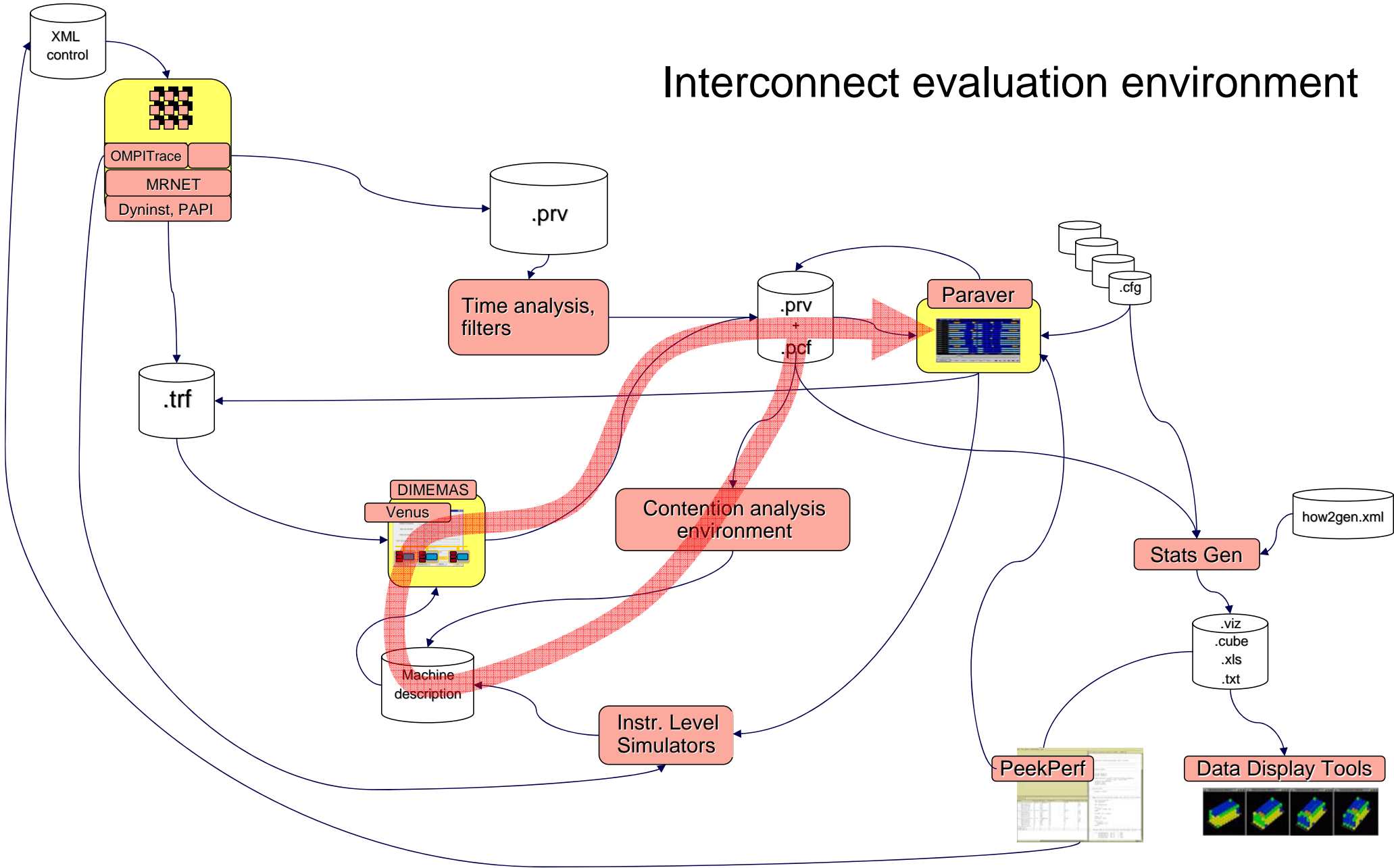
Network sensitivity

- Simulations with 4 processes per node
- NMM Iberia 4Km
 - Not sensitive to Latency
 - 512 sensitive to contention?
 - 256 MB/s OK
- ARW Iberia 4 Km
 - Not sensitive to Latency
 - sensitive to contention
 - Need 1GB/s



CEPBA-Tools Environment

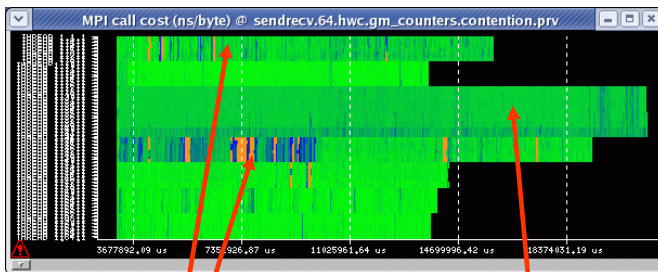
Interconnect evaluation environment



Contention impact

- on large systems?
- In multiuser environments

64 nodes, G=8, 4MB

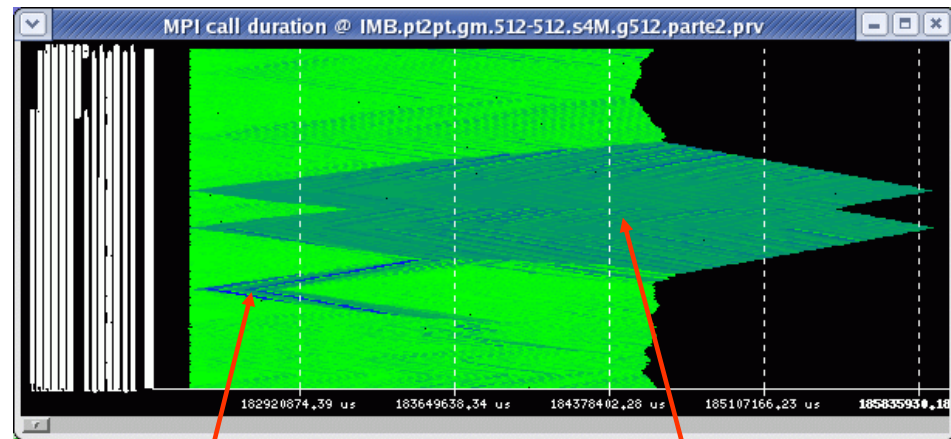
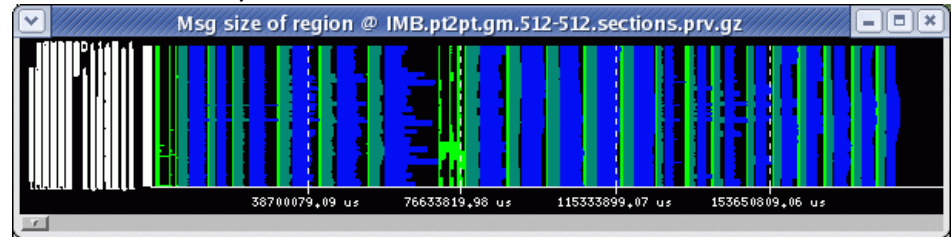


External contention

Internal contention

512 nodes, 4MB

Dependence on appl. phase (comm. Pattern)



Bubble propagation

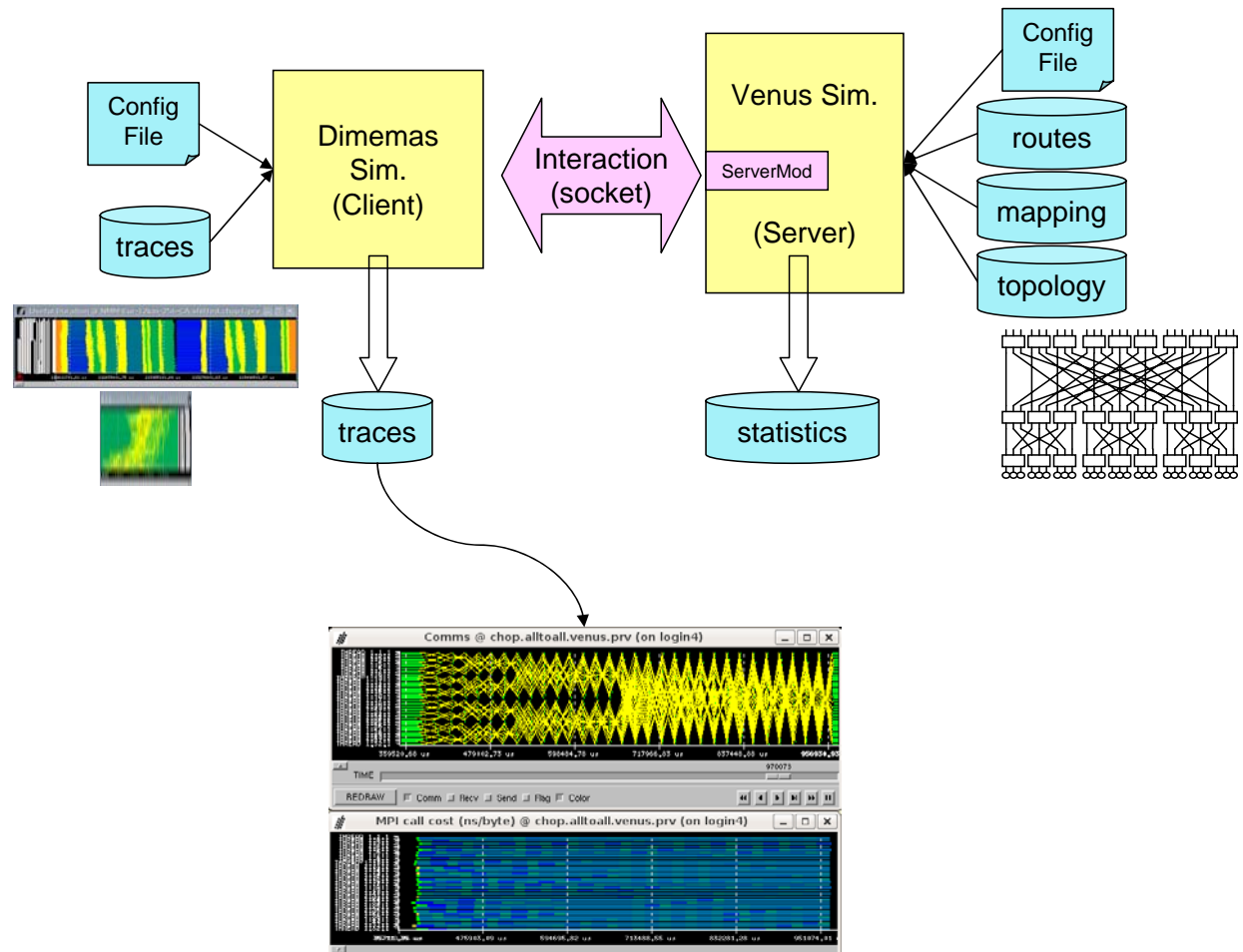
Propagation of internal contention

What is the benchmark measuring?
Appropriate number of iterations?



Interconnect simulation environment

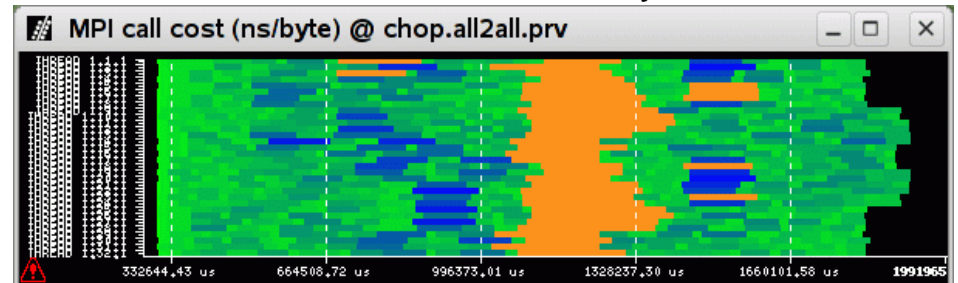
- Dimemas
 - Very fast
 - Sensitivity: identify coarse grain factors
 - identify relevant communication phases
- Venus
 - Very detailed network model
- Venus-Dimemas integration
 - Understand application usage of physical comm. resources
 - Communication subsystem design



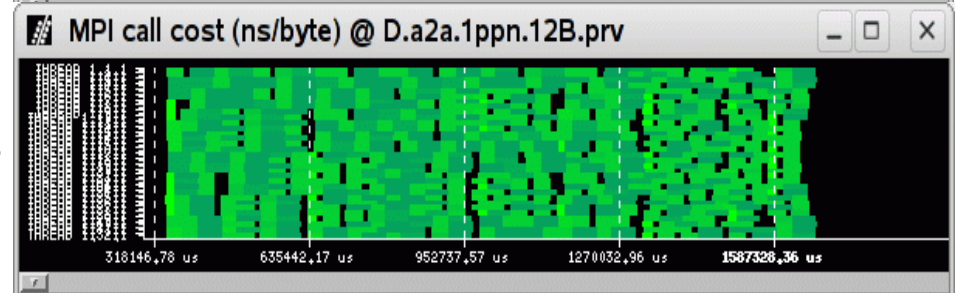
Contention

- P2p application trace chop
- @ collectives
- Evaluating
 - Routing
 - Oblivious
 - Pattern aware
 - Dynamic
 - Topology
 - Slimmed trees
 - Direct networks
 - Process mapping
 - Protocol
 - Eager limit
 - Packet segments

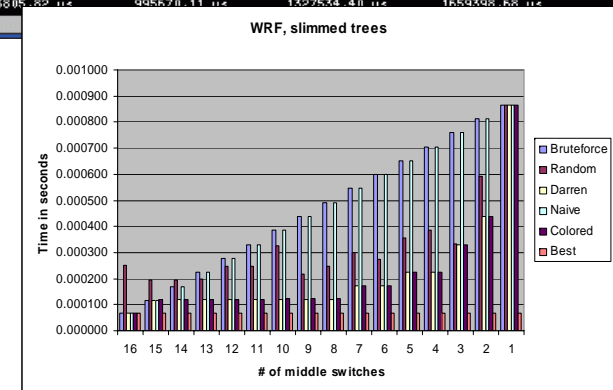
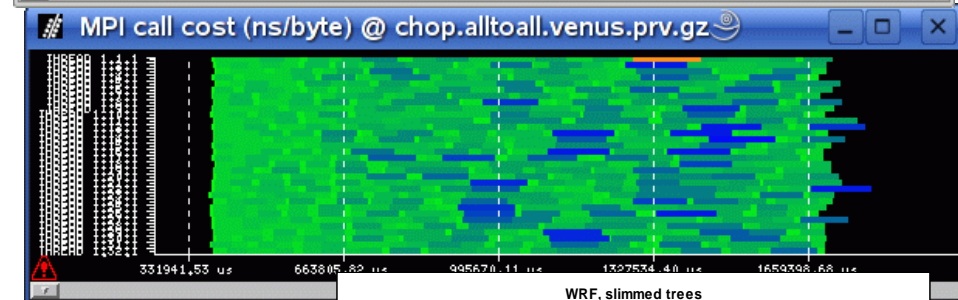
Actual run



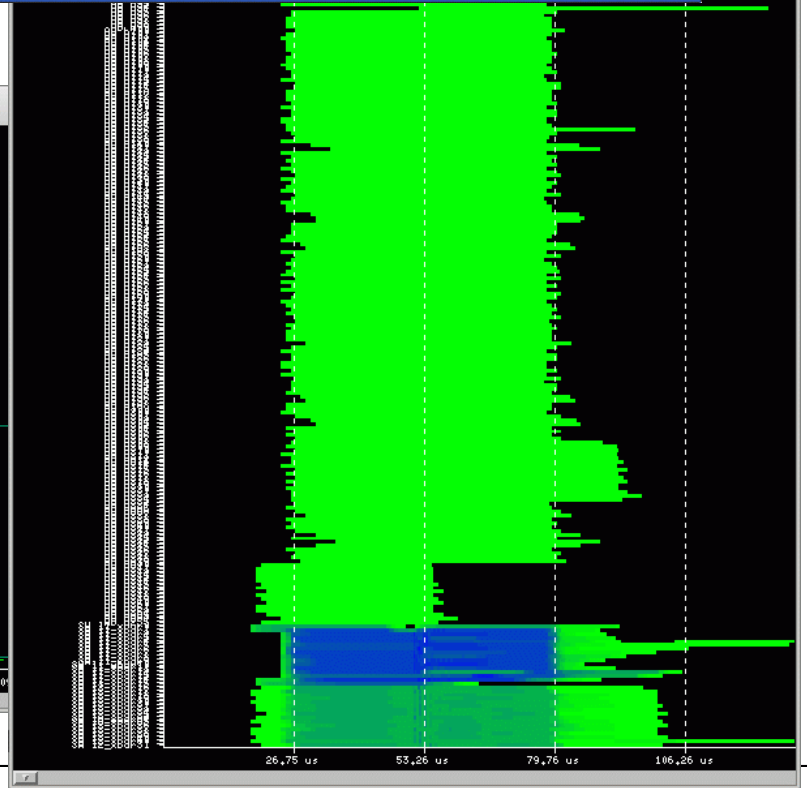
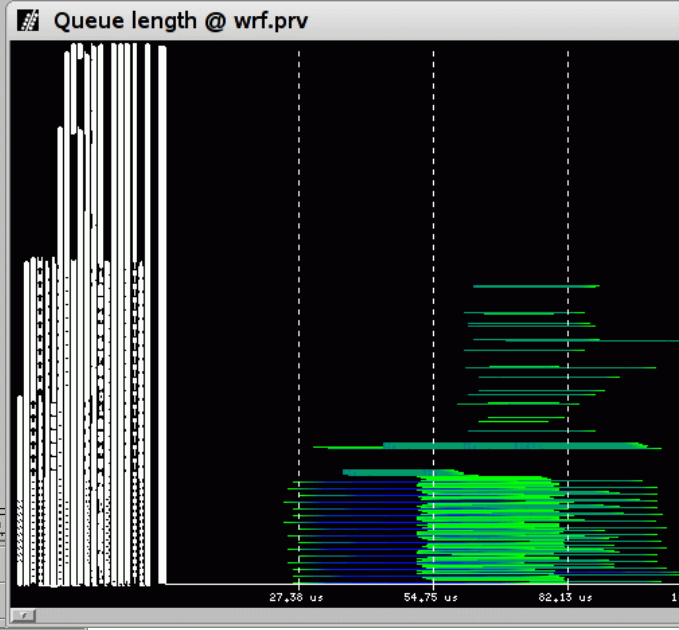
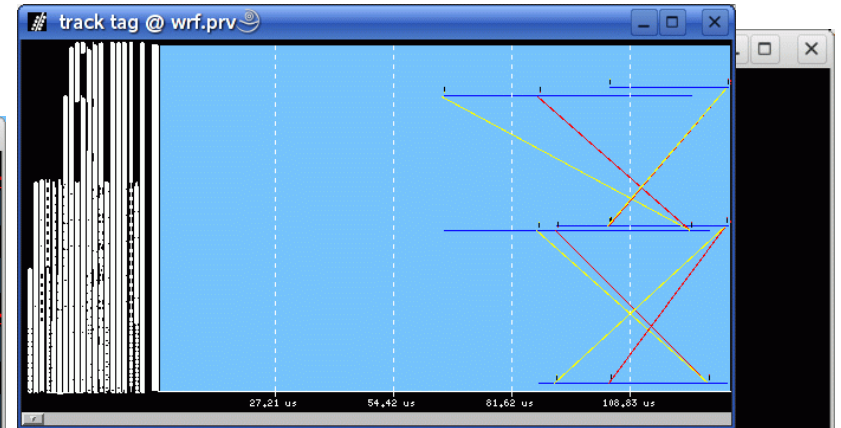
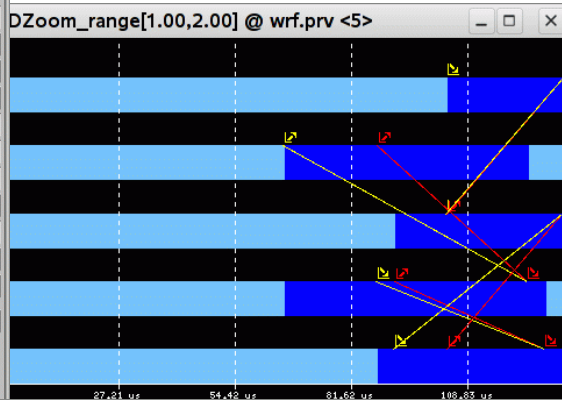
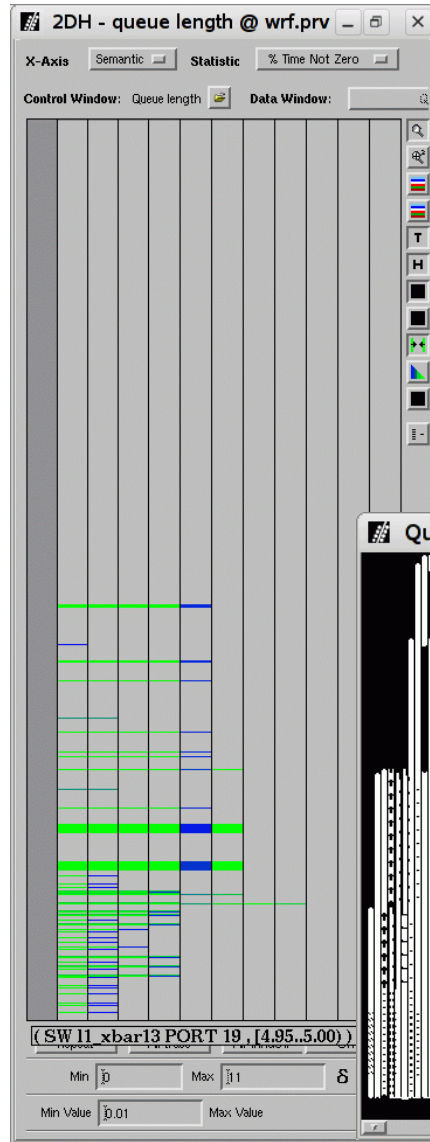
Dimemas



D&V



Internals network



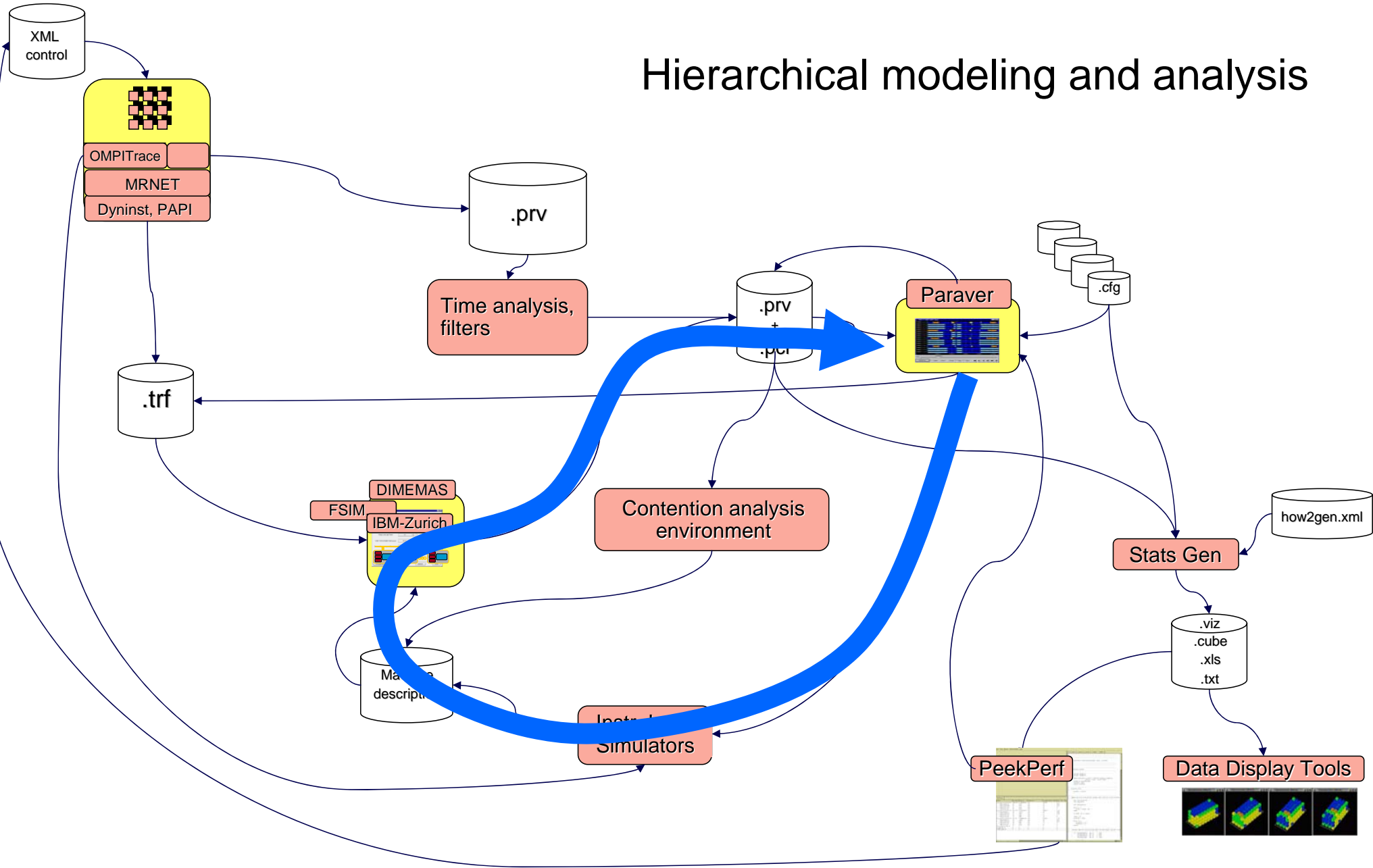


Hierarchical modeling

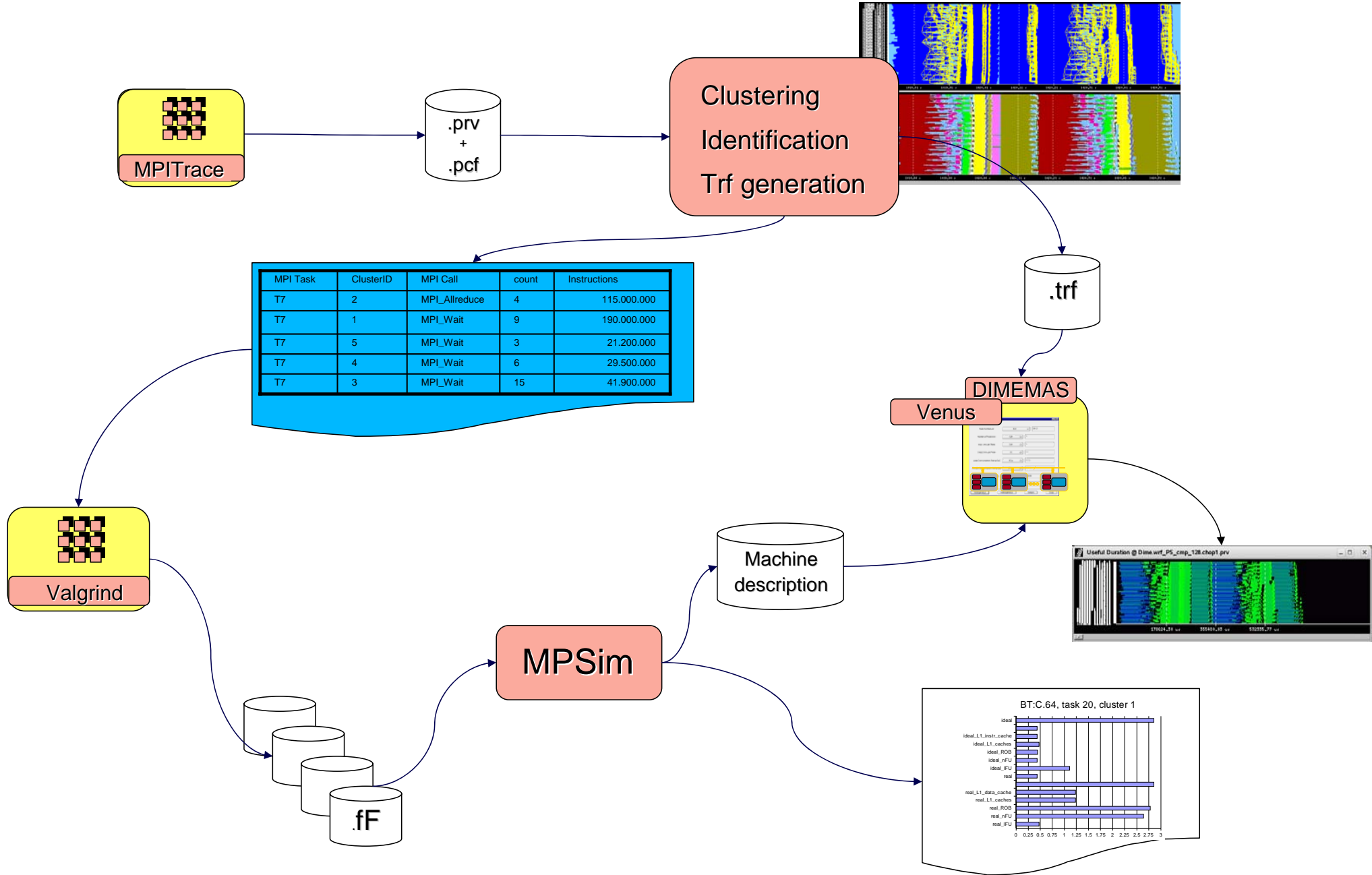


CEPBA-Tools Environment

Hierarchical modeling and analysis

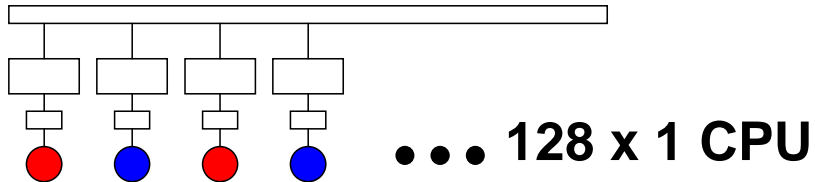


Hierarchical modeling and analysis

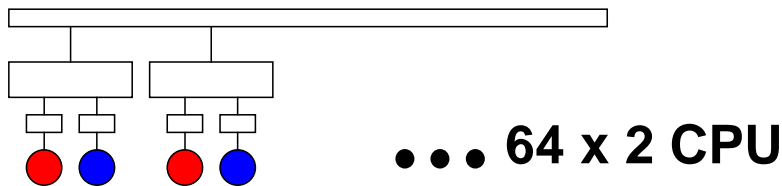


Turandot simulation results for WRF (table)

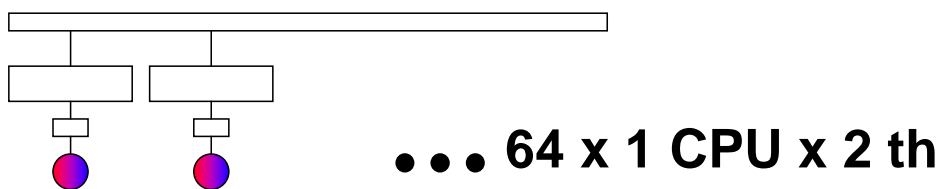
Single Core / Single Thread



Dual Core / Single Thread



Single Core / Dual Thread

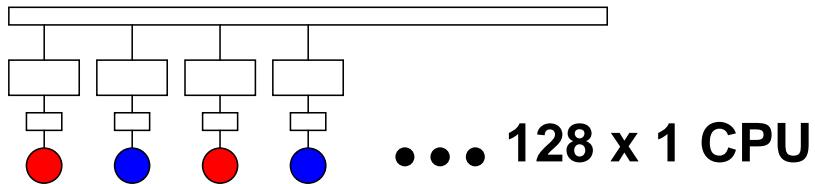


	IPC			
	Measured	Simulated P5 (TURANDOT)		
Cluster	MN	ST	CMP	SMT
1	0,54	0,75	0,75	0,56
2	0,49	0,72	0,72	0,52
3	0,62	0,83	0,72	0,43
4	0,72	0,94	0,83	0,50
5	0,79	1,07	0,84	0,52
	CPU ratio			
	Measured	Simulated P5 (TURANDOT)		
Cluster	MN	ST	CMP	SMT
1	1	1,38	1,38	1,04
2	1	1,47	1,47	1,05
3	1	1,34	1,16	0,70
4	1	1,30	1,15	0,69
5	1	1,36	1,06	0,66



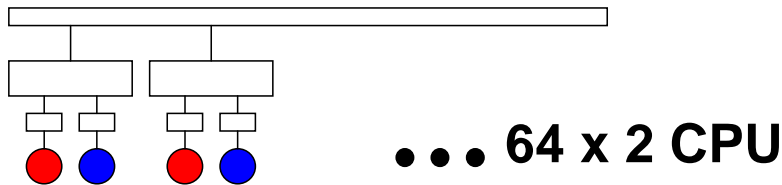
Dimemas simulation results for WRF

Baseline: ppc970 FX

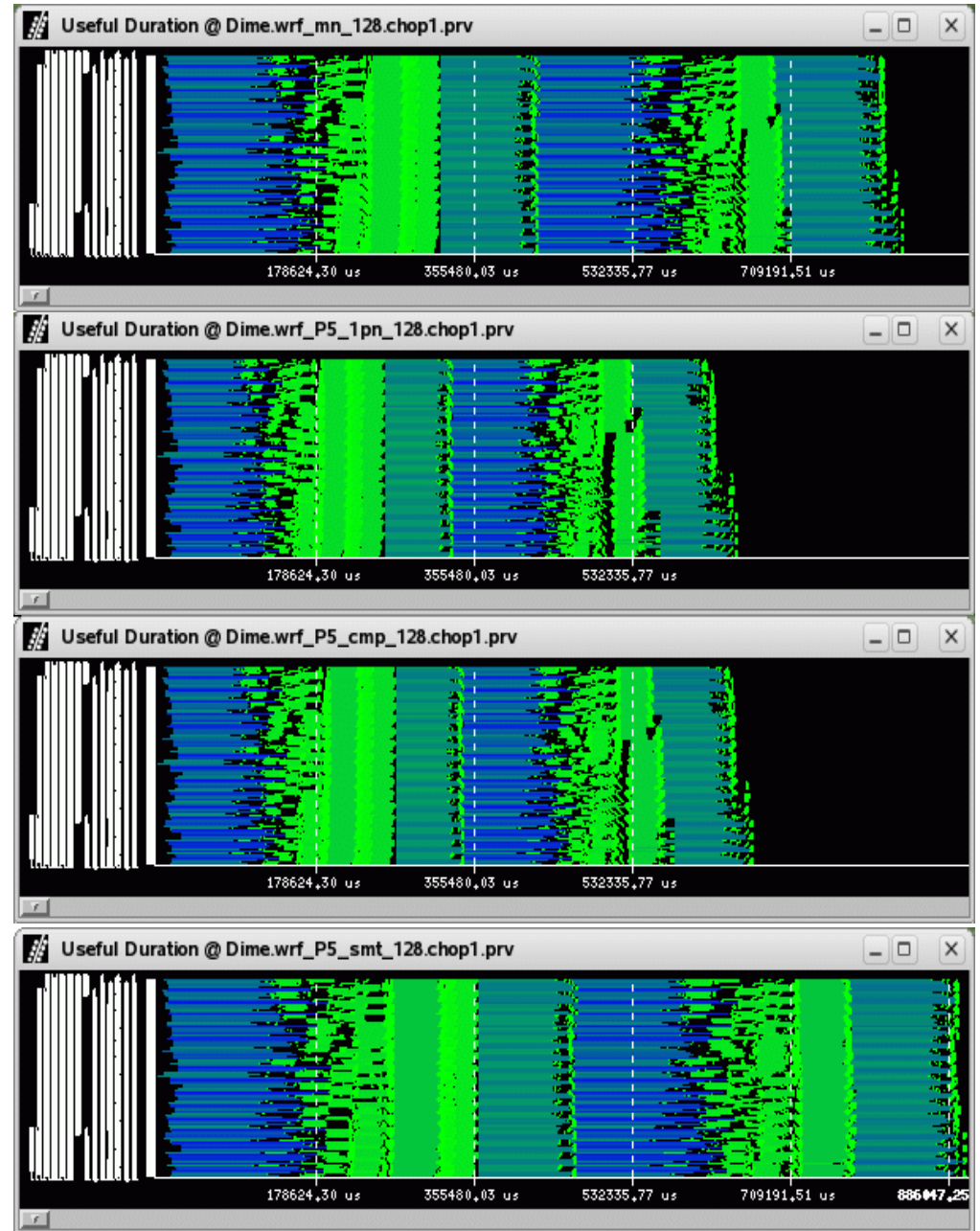
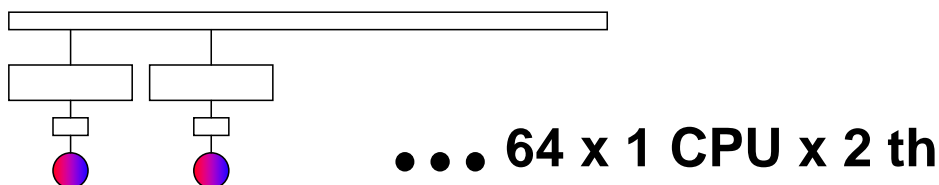


Replace 970 by Power5

Replace 970 by Power5 CMP



Replace 970 by Power5 SMT





Offers and interests



Offer



- Paraver OPEN SOURCE by end of the year
- Dimemas. OPEN SOURCE by end of the year
- Instrumentation. OPEN SOURCE by end of the year
- Structure analysis tools in development
 - Signal analysis
 - Clustering
 - Sampling + tracing
- OpenMP incl 3.0 tasks OPEN SOURCE
- CellSs /SMPSs OPEN SOURCE



Interest



- HWC
 - PAPI
 - CPIstack models from hardware counters → manufacturers
- Control infrastructure
 - MRNet, LaunchMON
- Instrumentation infrastructure
 - Dyninst, P^NMPI
- Profile display tools
 - File formats
 - Link to timeline mechanism
- Instrumentation control. Profiler based control
- Probe injection mechanisms
- Call stack API.
- Compiler information / model on application
- Simple code restructuring

