# Understanding applications with Paraver and Dimemas

tools@bsc.es

Feb 11th 2014

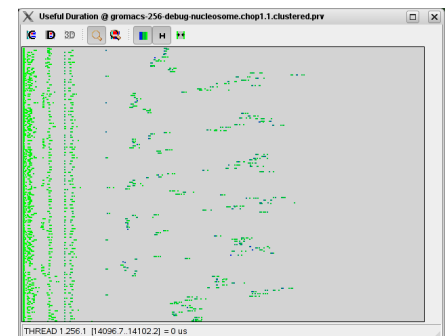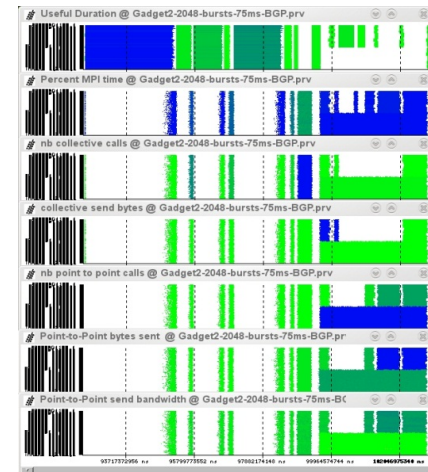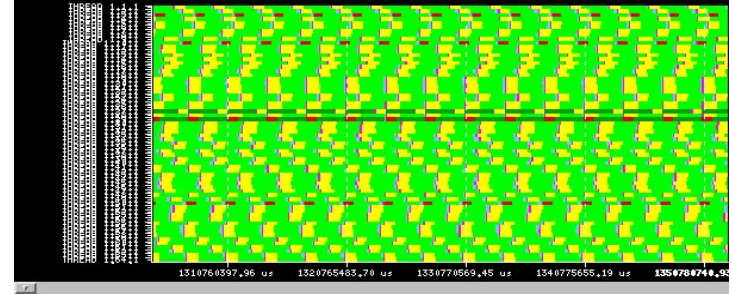- Tools presentation

- Demo: CG-POP analysis

- Hands-on

  - Core tools: Extrae, Paraver, Dimemas

  - Analytics modules: Clustering, Tracking, Folding

- 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
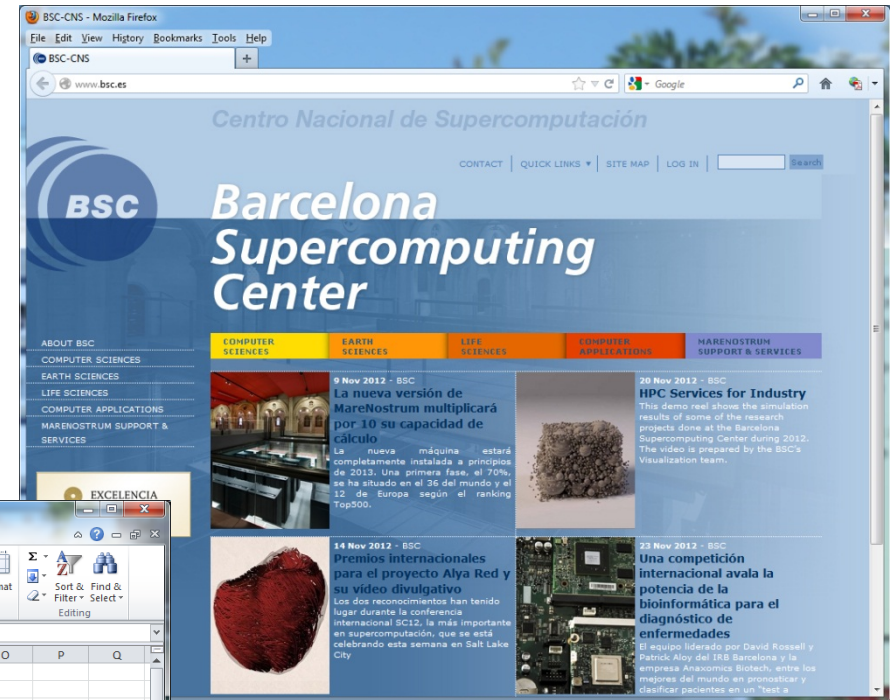    - Performance Analytics

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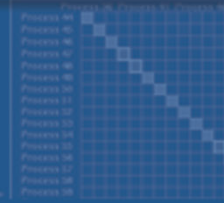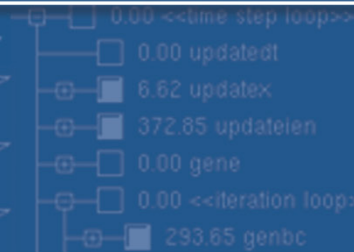  - Can combine into false colored but highly informative images

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

- **Behavioral structure** vs. syntactic structure
  - Algorithmic and performance
  - In space and time

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

- **Flexibility** to let analyst navigate the captured data and gain as much **insight** as possible from as **few application runs** as possible.

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**

# Paraver – Performance data browser



Raw data

Timelines

2/3D tables
(Statistics)

Trace visualization/analysis

+ trace manipulation

Goal = Flexibility

No semantics

Programmable

Comparative analyses

Multiple traces

Synchronize scales

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

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

- ## Types of functions
  - Categorical

  $$S_i \in [0, \mathrm{n}] \subset N, \qquad n <$$

    - State, user function, outlined routine

  - Logical

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

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

  - Numerical

  $$S_i \in R$$

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

- From timelines to tables

- By the way: six months later ….

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

- CESM: 16 processes, 2 simulated days

- Histogram useful computation duration shows high variability

- How is it distributed?

- Dynamic imbalance
  - In space and time
  - Day and night.
  - Season ? ☺

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



WRF-NMM
Peninsula 4km
128 procs

570 s
2.2 GB
MPI, HWC

570 s
5 MB

4.6 s
36.5 MB

See slides at end of presentation for details

# Performance Analytics

# Spectral analysis

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

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

# Scalability: online automatic interval selection



Back-end threads

T₀  T₁  Tₙ

Aggregate data

Broadcast results

MRNet Front-end

Clustering Analysis

Structure detection

Global periodicity @ ...RAE_Paraver_trace_2periods.prv (on login2)

0 ns    302,539,134,786 ns

Detailed trace for only small interval

*"Scalable tracing with dynamic levels of detail" ICPADS 2011*

# Clustering

# Using Clustering to identify structure



Automatic Detection of Parallel Applications Computation Phases. (IPDPS 2009)

SPECFEM3D

WRF 128 cores

GROMACS

Asynchronous SPMD

Balanced #instr
variability in IPC

SPMD
Repeated substructure

Coupled imbalance

MPMD structure

Different coupled
imbalance trends

duration vs. cluster

instr. vs. cluster

# Tracking

- OpenMX (strong scale from 64 to 512 tasks)



*"On the usefulness of object tracking techniques in performance analysis "  SC 2013*

# Folding

- ## Benefit from applications' repetitiveness

Iteration #1    Iteration #2    Iteration #3

Synthetic Iteration



Task 1 Thread 1 - copy_faces.26
Duration = 16.99 ms Counter = 69997.07 Kevents

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

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

- Performance of a sequential region  = 2000 MIPS


Is it good enough?


Is it easy to improve?



MIPS evolution for rrrmhd_mpi.16tasks.ppc.clustered.fused.extract.1.1.Cluster_1.0

# Instantaneous CPI stack



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

MRGENESIS

- Trivial fix.(loop interchange)
- Easy to locate?

- Next step?

- Availability of CPI stack models for production processors?
  - Provided by manufacturers?

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

# Correlating counters



CG-POP

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

CG-POP

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

# Dimemas

- ## Key factors influencing performance
  - Abstract architecture
  - Basic MPI protocols
  - No attempt to model details
- ## Objectives
  - Simple / general, Fast simulations
- ## Linear components
  - Point to point communication
  - Sequential processor performance (global CPU speed, per block/subroutine
- ## Non-linear components
  - Synchronization semantics (blocking receives, rendezvous)
  - Resources contention (CPU, links half/full duplex, busses)

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

- Paraver trace: what happens when
  - Actual wall clock time of events

- Dimemas trace: sequence of resource demands
  - Duration of computation bursts
  - Type of communication, partners and bytes

- Can be generated from Paraver trace
  - prv2dim input.prv output.dim

- Dimemas generates as output a Paraver file of the simulated run

- ## WRF, Iberia 4Km, 4 procs/node

  - ### No sensitive to latency

  - ### NMM

    - #### BW – 256MB/s
    - #### 512 – sensitive to contention

  - ### ARW

    - #### BW - 1GB/s
    - #### Sensitive to contention



Impact of latency (BW=256; B=0)



Impact of BW (L=8; B=0)



Contention Impact (L=8; BW=256)

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

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



GADGET @ Nehalem cluster 256 processes

Real run

Ideal network

Allgather + sendrecv

allreduce

alltoall

sendrec

waitall

Impact on practical machines?

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

GADGET



**Speedup**  —  64 procs

**Speedup**  —  128 procs

**Speedup**  —  256 procs

Bandwidth (MB/s), CPU ratio

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

GADGET, 128 procs

%elapsed time

Code region

Speedup

93.67%

Bandwdith (MB/s)

CPU ratio

Speedup

97.49%

Bandwdith (MB/s)

CPU ratio

Speedup

99.11%

Bandwdith (MB/s)

CPU ratio

# Methodology

# Help generate hypotheses

# Help validate hypotheses

**Qualitatively**

**Quantitatively**

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

- Serial efficiency – how far from peak performance?
  - IPC, correlate with other counters

- Scalability – code replication?
  - Total #instructions

- Behavioral structure? Variability?

Paraver Tutorial:
Introduction to Paraver and Dimemas methodology

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

$$C \quad o = \mu L m * T \quad B \quad m \, r$$



$T_{ideal}$

Migrating/local load imbalance
Serialization



**µLB**

$$\mu L\,B = \frac{m \, a \, x \, T()}{T_{ideal}} \qquad Transfer = \frac{T_{ideal}}{T}$$

$$\eta_{\parallel} = L \quad * \underset{t}{\mu}BL \quad * TB \quad r$$

CG-POP mpi2s1D - 180x120

Good scalability !!
Should we be happy?



speed up





$$\eta = \eta_{\parallel} \quad * \eta_{i \ n \ s \ t} \eta_{H \ P}$$

- ## 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

# Short Paraver Demo

# BSC tools hand-on

- Install the Paraver binaries in your laptop
  – Binaries for linux x86, x86-64, windows & mac (not tested on Mavericks)


- Configure Paraver package
  – Set-up the tutorials

- Copy ~nct00001/gpfs_projects/tuesday_material into your ${HOME}
  - `cp -r /gpfs/projects/nct00/nct00001/tuesday_material ~`

- Contents of tuesday_material
  - bin/
    - Some BSC tools scripts to make your life easy
  - slides/
    - All the slides wrt BSC tools
  - packages/
    - Paraver binaries
  - tutorials/
    - Paraver tutorials
  - documentation/
    - Documentation related to the BSC tools
  - jobscripts/
    - Modified jobscript of the course applications
  - extrae/ | dimemas/ | clustering/ | folding/ | tracking/
    - Files for the different tools used in the Hands-On session

Extrae

- Build the CGPOP application (already done)
  - Copy from `apps/CGPOP`
  - Issue `./build` in your copy


- Generate trace-file for the CGPOP application
  - Change dir to `~/apps/CGPOP/jobscript/marenostrum3`
  - Copy `tuesday_material/jobscripts/cgpop.extrae.lsf`
  - Edit `cgpop.extrae.lsf` and uncomment module load bsctools & choose appropriate TRACE
    - Optionally choose a tracefile name
  - Submit the job
    - `bsub < cgpop.extrae.lsf`

# Inside cgpop.extrae.lsf

```
#!/bin/bash

#BSUB -n 24
#BSUB -oo cgpop_%J.out
#BSUB -eo cgpop_%J.err
#BSUB -R"span[ptile=12]"
#BSUB -x # Exclusive use
#BSUB -J cgpop
#BSUB -W 00:10
#BSUB -U tools


# module load bsctools


# Override tracefile name using this environment variable
# export TRACE_NAME=cgpop.linux_icc.180x120.24tasks.prv


# Choose appropriate instrumentation type for your application
# export TRACE=~/tuesday_material/extrae/trace.mpi.c.sh # For C-based applications without sampling
(MPI)
# export TRACE=~/tuesday_material/extrae/trace.mpi.sampling.c.sh # For C-based applications with
sampling (MPI)
# export TRACE=~/tuesday_material/extrae/trace.mpi.f.sh # For Fortran-based applications without
sampling (MPI)
# export TRACE=~/tuesday_material/extrae/trace.mpi.sampling.f.sh # For Fortran-based applications with
sampling (MPI)


time mpirun -np 24 --npersocket 6 --bind-to-core ${TRACE} ../mpi2s1D/cgpop.linux_icc.180x120
```

- ## What is inside the shellscript files?

```
#!/bin/bash

# Workaround for MN3
export TMPDIR=$TMPDIR/extrae
mkdir -p ${TMPDIR}


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


$@
```

# Library alternatives

| Library | Serial | MPI | OpenMP | pthread |
|---|---|---|---|---|
| libseqtrace | ✔ | | | |
| libmpitrace[f][1] | | ✔ | | |
| libomptrace | | | ✔ | |
| libpttrace | | | | ✔ |
| libompitrace[f] [1] | | ✔ | ✔ | |
| libptmpitrace[f] [1] | | ✔ | | ✔ |

[1] for Fortran codes

# Extrae configuration (I)

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

      <set enabled="yes" domain="all" changeat-time="100000us">
PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_L1_DCM,PAPI_L2_DCM,PAPI_L3_TCM,PAPI_FP_INS,PAPI_BR_MSP
      </set>
      <set enabled="yes" domain="all" changeat-time="100000us">
PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_LD_INS,PAPI_SR_INS,RESOURCE_STALLS,PAPI_BR_UCN,PAPI_BR_CN,PAPI_VEC_SP
      </set>
      <set enabled="yes" domain="all" changeat-time="100000us">
PAPI_TOT_INS,PAPI_TOT_CYC,RESOURCE_STALLS:LB,RESOURCE_STALLS:RS,RESOURCE_STALLS:SB,RESOURCE_STALLS:ROB,PAPI_VEC_DP
      </set>                                                    Hardware counters
    </cpu>
    <network enabled="no" />
    <resource-usage enabled="no" />
    <memory-usage enabled="no" />
</counters>
```

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

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

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

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

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

```
<merge enabled="yes"
   synchronization="default"  tree-fan-out="16"  max-memory="512" joint-states="yes"  keep-mpits="yes"
   sort-addresses="yes"  overwrite="yes">  $TRACE_NAME$
</merge>
```
Trace generation

# Paraver

- Analyse the results with Paraver

- Follow tutorials

- Paraver navigation
    - Load configuration files
    - Generate new views
    - Contextual menus
        - Synchronize between windows
        - Zoom & Fit time-scale
        - Draw mode
    - Save configuration files

# Dimemas

- **Step 1**: CGPOP Paraver trace (chop) available at

  `cd ${HOME}/tuesday_material/dimemas/`


- **Step 2**: Define Dimemas configuration
  - We supply a basic Dimemas configuration (MN3.cfg)


- **Step 3**: Execute Dimemas

  **`${HOME}/tuesday_material/bin/dimemas-sim.sh <input_trace> <dimemas_cfg>`**


- **Step 4**: Analyse the results with Paraver

- Run the Dimemas GUI

  > `DimemasGUI`

- Load MN3.cfg file

- Tune target machine parameters
  - Processors 10x faster
  - Set network bandwidth to ¼ to the original one

```
[…]
#PARAVER_TRACE=${1}
PARAVER_TRACE=${1}.prv
DIMEMAS_TRACE=${1}.dim
DIMEMAS_CFG=${2}


[…]


# Translate from .prv to .dim
if [ ! -f ${DIMEMAS_TRACE} ]; then
  echo
  echo "Dimemas trace does not exist. Translating input Paraver trace"
  echo
  prv2dim ${PARAVER_TRACE} ${DIMEMAS_TRACE}
  echo
  echo "==================================================================="
fi


# Simulate parameter -S 32K fixed by default
echo
echo "Executing Dimemas"
echo
Dimemas -t --dim ${DIMEMAS_TRACE} -S 32K -p ${OUTPUT_PARAVER_TRACE} ${DIMEMAS_CFG}

echo "==================================================================="
```

**Paraver to Dimemas trace translation (if required)**

**Dimemas execution**

- A Paraver trace of the simulation
  - Output trace name: `<input_trace>.<cfg>.prv`

- Run Paraver and compare with the original trace

1. Follow guidelines for paraver basic navigation (Tut #1)
   Tracefile provided

2. Adapt scripts to instrument your application

3. Obtain a tracefile

4. Follow first steps of methodology guidelines (Tut #5)

   – Parallel efficiency

   – Distribution of computations

   – Instruction balance

   – …

5. Depending on the diagnosis use clustering, folding, dimemas…

# Performance Analytics

# Clustering

- **Step 1**: Previously generated trace SU3 trace available at:
  ```
  cd ${HOME}/tuesday_material/clustering/
  ```

- **Step 2**: Tune the configuration xml
  - We supply a configuration *xml* (cluster.xml)
    - Completed Instructions vs. IPC
    - DBSCAN parameter Eps = 0.01, MinPoints = 10
    - Adjusted filters for this trace

- **Step 3**: execute the cluster analysis:
  ```
  ${HOME}/tuesday_material/bin/clusterize.sh
  <trace_without_prv>
  ```

- To launch the clustering binary directly
  ```
  BurstClustering -a -d <xml_file> -i <in_trace> -o
  <out_trace>
  ```

# Clustering definition XML

```xml
<clustering_definition use_duration="no" apply_log="yes"
normalize_data="yes" duration_filter="1000" threshold_filter="0">

  <clustering_algorithm name="DBSCAN">
    <epsilon>.01</epsilon>
    <min_points>4</min_points>
  </clustering_algorithm>


  <clustering_parameters>

    <mixed_events apply_log="no" name="IPC" operation="/">
      <event_type_a>42000050</event_type_a>
      <event_type_b>42000059</event_type_b>
      <factor>1.0</factor>
    </mixed_events>

    <single_event apply_log="yes" name="PAPI_TOT_INS">
      <event_type>42000050</event_type>
<!--      <range_min>4e7</range_min> -->
      <factor>1.0</factor>
    </single_event>

  </clustering_parameters>
  […]
```

**Duration Filters and Normalizations**

**DBSCAN parameters**

**Clustering dimensions using Paraver events**

```
[…]

<extrapolation_parameters all_counters="yes"/>

<output_plots all_plots="no">

  <plot_definition raw_metrics="yes">
    <x_metric title="IPC">IPC</x_metric>
    <y_metric title="Instructions Completed">PAPI_TOT_INS</y_metric>
  </plot_definition>

</output_plots>

</clustering_definition>
```

**Events to characterize the bursts**

**Definition of output plots. Must be events previously defined**

- GNUplot scripts

  gnuplot **<input_trace>.clustered.[…].gnuplot**

- A Paraver trace with clusters information
  - Output trace name: **<input_trace>.clustered.prv**

- Clusters statistics (including the extrapolation)
  - Statistics file name:
    **<input_trace>.clustered.clusters_info.csv**
  - You can visualize it with an editor or import it to a spreadsheet

- What to do if…
  - … you find too many clusters?
    - Increase the value of Eps.

  - … there too many points (> 100K)
    - Increase the filters (duration / instructions)
    - You may need to cut the trace

# Tracking

- **Step 1:** Previously clustered traces (SU3) available at:

  ```
  cd ${HOME}/tuesday_material/tracking
  ```


- **Step 2**: Execute the tracking analysis:

  ```
  ${HOME}/tuesday_material/bin/track.sh <in_trace_1>
  <in_trace_2> … <in_trace_N>
  ```

- **Step 3**: Load the results with the visualizer

  `xtrack TRACKING.RESULTS.xtrack`

  - Other outputs
    - Scatter plots
      - gnuplot **<input_trace>.clustered.[…].gnuplot.scaled**
      - gnuplot **TRACKING.RESULTS.recolored.multiplot**

    - Paraver traces with clusters renamed
      - wxparaver **<input_trace>.clustered.tracked.prv**

# The 'xtrack' visualizer



Manipulate the graph area

Select / hide clusters

Correlation of all metrics for selected cluster across experiments

Select / hide metrics

Change axes

Switch between frame/trajectory view

Open Paraver trace

Run timed animation

Navigate to previous/next experiment

Average/dispersion for selected metric/cluster across experiments

- ## What to do if…

  - ### … there's too many objects?
    - Filter clusters that represent low percentage of time
      - Add argument "-m <time-percentage>" (i.e. –m 5)

  - ### … over aggregation?
    - Turn off the tracking heuristic based on callstack.
      - Remove argument: "-c <callstack-depth>"

  - ### … the application is not SPMD?
    - Turn off the tracking heuristics based on alignment.
      - Add argument: "-a 1"

# Folding

**Folding**

- **Step 1**: Previously clustered traces (SU3) with sampling available at:

  ```
  cd ${HOME}/tuesday_material/folding
  ```

- **Step 2**: Apply folding

  ```
  bsub < folding.lsf
  ```

- **Step 3**: Browse results

- The folding generates plots for the combinations of

  - Cluster
  - Performance counter
    - Plus all counter slopes combined
    - Plus MIPS and remaining counters in terms of ctr/instruction
    - *Architecture impact*
    - *Stall distribution*

- For instance

```
# gnuplot -persist *Cluster_1.*PAPI_TOT_INS*gnuplot
# gnuplot -persist *Cluster_1.*stalldistribution*gnuplot
# gnuplot -persist *Cluster_3.*stalldistribution*gnuplot
# gnuplot -persist *Cluster_3.*architectureimpact*gnuplot
```

- Execute in the directory results

```
# wxfolding-viewer *wxfolding
```

- In the GUI choose the cluster you want to analyze, and double click on the performance counter



Choose cluster to analyze

Choose metric or model

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

THANKS

# Detailed material

# Semantic Module

- 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

Trace

Filter

Semantic

Display

MPI call @ AMBER-64.Z.chop1.prv

4850856.78 us    4869052.03 us    4927247.29 us    4965442.54 us    5003637.79 us

Communications that pass through the filter

Events that pass through the filter

Show list of event types

- From Events to functions of time

    – Last event value $\quad S(i) = V(Ev_i)$

    – Next event value $\quad S(i) = V(Ev_{i+1})$

    – Average Next Event Value $\quad S(i) = \dfrac{V(Ev_{i+1})}{T(Ev_{i+1}) - T(Ev_i)}$

    – Interval btw. Events $\quad S(i) = T(Ev_{i+1}) - T(Ev_i)$

- From communication records to functions of time

    - Send Bytes

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

    - Send Bandwidth

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

    - Msgs in transit

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

    - Recv. Bandwidth

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

    - Rec. Negative Msgs

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

    - Comm. Partner

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

    - Bytes btw. Events

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

- S'(t) = f(S(t))                        S' = f ° 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

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

Loads

Stores

L2 Line Loads

x100

Mem Ops

L2 miss ratio

- Derived windows
  - Point wise operation
    - $S = \alpha * S^a \; <op> \; \beta * S^b$
    - $<op> : + , -, *, /, \ldots$



Interval between MPI events



In MPI call



MPI call duration

Ev: 2,1
Ev: 4,5

Ev: 2,0
Ev: 4,9

S: 3                    S: 5                    S: 3

0          1                          4                    5   Time

- ## Thread function: State as is

                          5
        3                              3

0          1                          4                    5   Time

- Useful for

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

Ev: 2,1          Ev: 2,0
Ev: 4,5          Ev: 4,9

S: 3                    S: 5                    S: 3

0          1                              4                    5   Time

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

                    1
  undefined                           0
  0          1                    4              5   Time

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

Ev: 2,1　　　　　　　　　　Ev: 2,0
Ev: 4,5　　　　　　　　　　Ev: 4,9

S: 3　　　　　　　　　S: 5　　　　　　　S: 3

0　　　　1　　　　　　　　　4　　　　5　Time

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

5　　　　9　　　　undefined

0　　　　1　　　　　　　　　4　　　　5　Time

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

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



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

# Semantic module: Examples

USR_FCT, idA USR_FCT, idB USR_FCT, 0 USR_FCT, idB USR_FCT, idC USR_FCT, 0 USR_FCT, 0 USR_FCT, 0

Timeline

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

idA | idB | idA | idB | idC | idB | idA

- Useful for
  - Routine

- # Process model
  - – Thread, task, application, workload

- # Resource model
  - – CPU, node, system

- Semantic value: S(t)

- $S = f_{comp2} \circ f_{comp1} \circ f_{Workload} \circ f_{Application} \circ f_{task} \circ S_{thread}$

- Semantic functions

  - $f_{comp2}$, $f_{comp1}$: sign, mod, div, in range, select range
  - $f_{Application}$, $f_{Workload}$ : add, average, max, select
  - $f_{task}$ : add, average, max, select
  - $S_{thread}$: in state, useful, given state,
  -         last event value,
  -         next event value,
  -         average next event value
  -         interval between events, …

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

- Semantic functions

  - $f_{System}$ : add, average, max, select
  - $f_{Node}$ : add, average, max, select
  - $f_{CPU}$ : active thread, select
  - $S_{thread}$: in state, useful, given state, next event value, thread_id

# Analysis Module

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



MPI call, user function,…

Thread

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

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

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

duration, instructions, BW, IPC, ...



Processors

THREAD 1.126.1 [4,09433e+08..4,11813e+08) = 0

Instructions

Thread

NULL entry

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

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

- Single flexible quantitative analysis mechanism
- Let

  For each window w

  - cw$_1$ and cw$_2$ two views we will call control views
  - dw a view we will call data window

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

- For each control window we define a set of bins

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

- And the discriminator functions

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

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

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

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

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

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

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

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

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

$$NumBurst(th, j, k) = i_{end} - i_{start} + 1$$

$$i_{start} = \min(i) \mid t_i > t_{start}, i_{end} = \max(i) \mid t_i < t_{end}$$

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

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

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

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

# Distributed Configurations

- 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



**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**

# Load configuration files



Select directory

Navigate through directory tree

List of directories and configuration files in current directory

**APPLIED TO THE CURRENT TRACEFILE**

**Shortcuts**:

Drag and move (D&M) – Zoom

Control D&M – Zoom XY

Shift D&M – Timing

# How to generate table and change statistic

VI-HPS

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



Paraver

File    Help

Window browser

C:\trazas\TEXT\LBC\trace_lbc_16x4_cores.prv

- Task
- 2DP - Tasks
- Useful MIPS
  - MIPS
  - Useful

Files & Window Properties

| Name | 2DP - Tasks |
| Begin time | 0 us |
| End time | 5.433.839,38 us |
| **Control** | |
| Window | Task |
| Minimum | 0 |
| Maximum | 200 |
| Delta | 1 |
| **Statistics** | |
| Type | Semantic |
| Statistic | % Time |
| Minimum Gradient | 0,107903 |
| Maximum Gradient | 100 |
| **Data** | |
| Window | Task |

2DP - Tasks @ trace_lbc_16...

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



**3D control window: determines planes**

**Actual Plane on display**

# Table information and control

# Table information and control

VI-HPS

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)

**Paraver**

File   Help

Window browser

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

Window properties

| | |
|---|---|
| Name | MPI call duration |
| Begin time | 15.500,13 us |
| End time | 25.867,11 us |
| **Control** | |
| Window | MPI call duration |
| Minimum | 25,719 |
| Maximum | 4024,328 |
| Delta | 19,993045 |
| **Statistics** | |
| Type | Semantic |
| Statistic | Time |
| Minimum Gradient | 25,849453 |
| Maximum Gradient | 2305,542 |
| **Data** | |
| Window | MPI call duration |
| **3D** | |
| 3rd Window | MPI call |
| Minimum | 6 |
| Maximum | 124 |
| Delta | 1 |
| Plane | MPI_Allreduce |

Selected plane

MPI call duration @ 480.chop1.prv
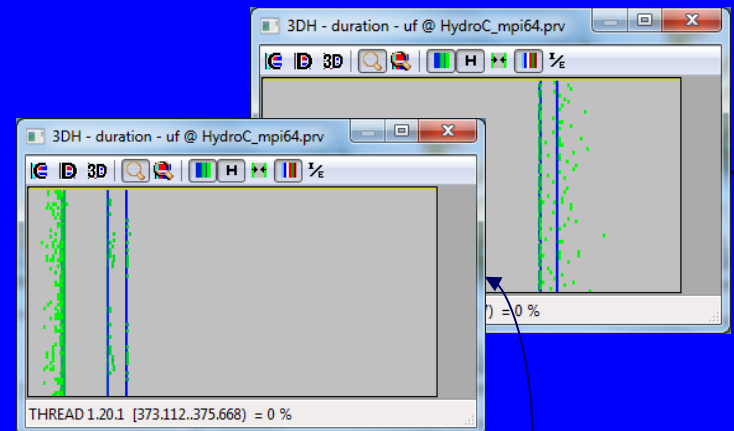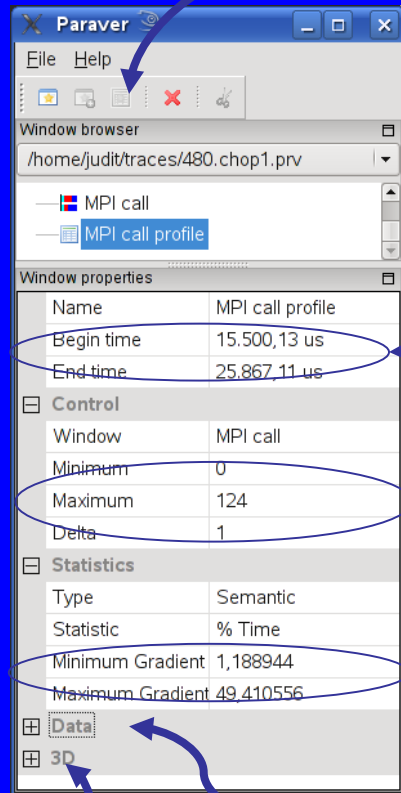
Right click

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

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...

Generate ASCII file with table data

Time
Objects
Size
Semantic Scale
Control scale
3D scale

**Paste Special**

Select properties to paste:

☑ Time
☑ Objects
☑ Size
☐ Semantic Scale

Cancelar    Aceptar

✓ Linear
Steps
Logarithmic
Exponential

Semantic
Objects
Both

**Shortcuts (zoom mode only)**:

Drag and move (D&M) – Zoom

Control D&M – Zoom XY

- 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**

- ## 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

- 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**



Cut & Filter

Trace
C:\trazas\Hydro\MPI+CUDA\cuHydroC.prv    [Browse]

☑ Load the resulting trace

[Load XML...]  [Save XML...]

Ordered execution
☐ 1.- Cutter
☑ 2.- Filter
☐ 3.- Software Counters

Cutter | **Filter** | Software Counters

Discard Records
☐ State  ☑ Event  ☑ Communication

States
☐ Idle
☑ Running
☐ Not created
☐ Waiting a message
☐ Bloking Send
☐ Thd. Synchr.
☐ Test/Probe

[Select all]
[Unselect all]

Min. burst time  3000

Events

[Add]
[Delete]
☐ Discard

Communications
Size  0                          MB

[OK]  [Cancel]

- Load a filtered trace and use the scissors tool

**Scissors tool**

**Browse to select file from which the cut will be obtained**

**Select cutter**

**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**

**Default setups**

# Extrae

**appl.job**

```bash
#!/bin/bash

export NP=8
export INPUT=$1

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

# Adapt job submission script

**appl.job**

```bash
#!/bin/bash

export NP=8
export INPUT=$1

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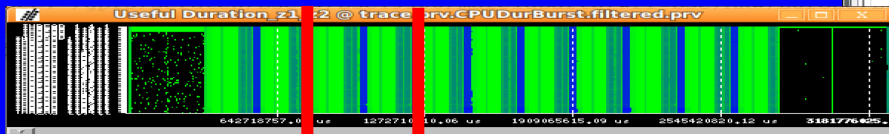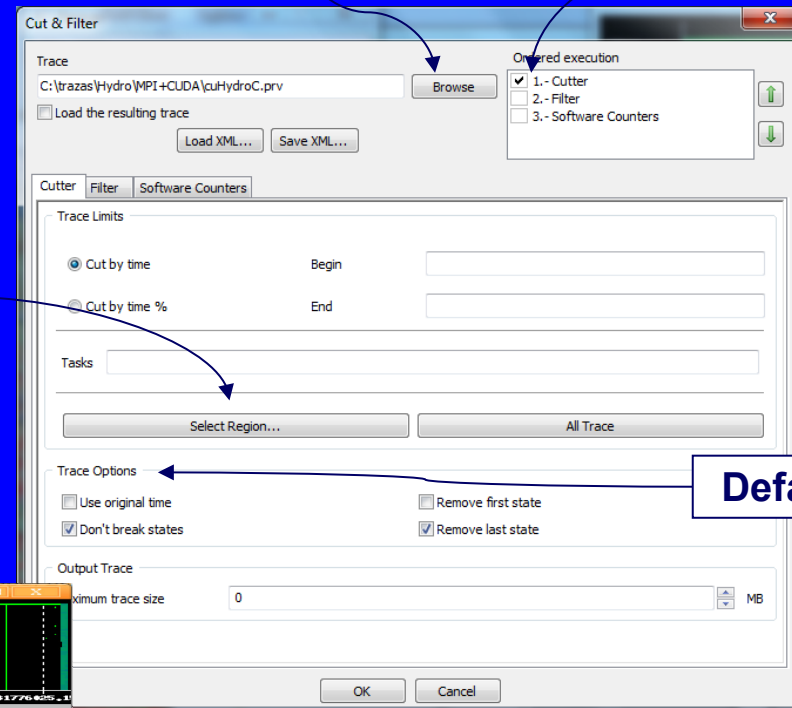
**trace.sh**

```bash
#!/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

$@
```

```
<?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>

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

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

...
```

**Activate MPI tracing and emit hardware counters at MPI calls**

**Do not activate OpenMP tracing**

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

Details in $EXTRAE_HOME/share/example/MPI/extrae_explained.xml

extrae.xml (cont)

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

…

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

**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,….)**

extrae.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)**

…

**extrae.xml   (cont)**

```
…

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



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

…
```

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

**Stop tracing after elapsed time …**

**… or when signal received**

extrae.xml   (cont)

…

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

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

**… plus summarized MPI events**

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

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

…

extrae.xml   (cont)

```
…

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


</trace>
```

**Merge individual traces into global application trace at end of run …**

**… into this trace name**

- Library depends on programming model

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

[1] for Fortran codes

# Using Dimemas

- Paraver → Dimemas trace Generation
  - prv2dim original.prv dimemas.dim

  - Default: duration of each computation region taken from .prv computation duration

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


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

Force synchronized start of all threads

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

# Dimemas GUI – Specify trace to simulate

# Dimemas GUI – Specify target machine

# Collective Communication Model

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

  - Specified in input file

# Scalability

- Linpack @ Marenostrum: 10k cores x 1700 s
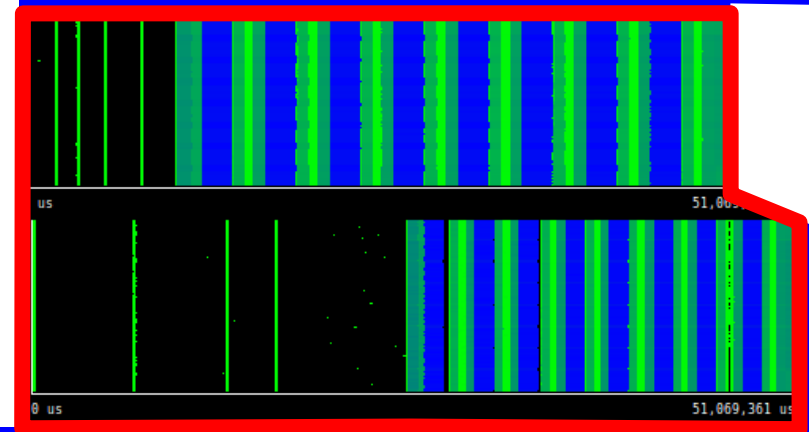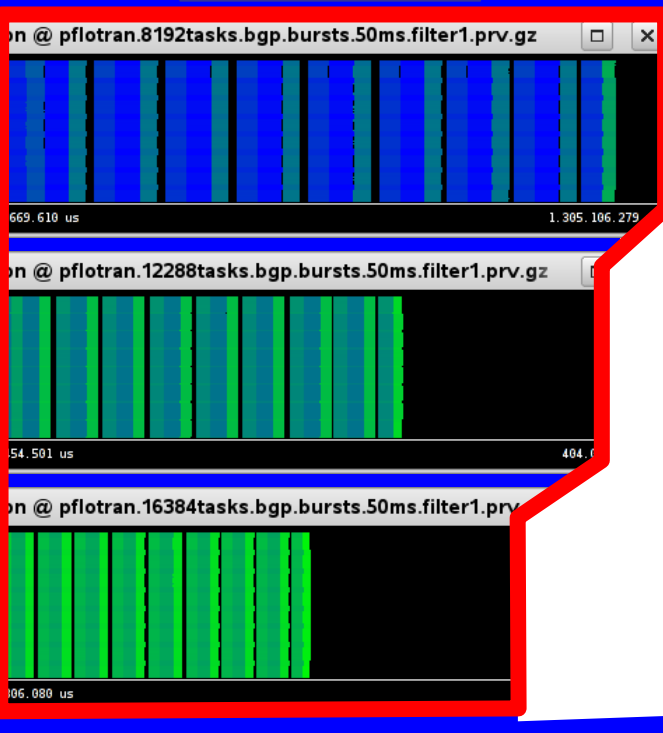
Jugene

Jaguar

~ 105 seconds

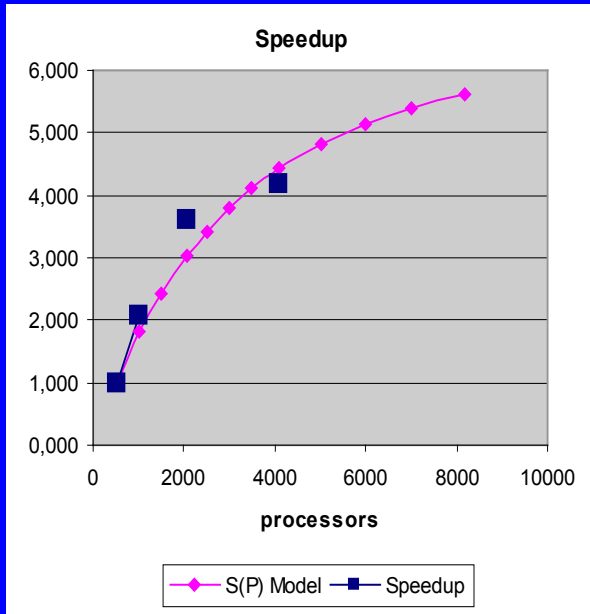~ 47 seconds

8K cores

12K cores

16K cores

Tran  Flow

Flow  Tran

PFLOTRAN

- ## 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

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

# Software counters



GADGET, PRACE Case A, 1024 procs

Useful duration — Useful Duration @ Gadget2-1024-300ms-1-5callers-BGP.prv

% MPI time — Percent MPI time @ Gadget2-1024-300ms-1-5callers-BGP.prv

# collectives — nb collective calls @ Gadget2-1024-300ms-1-5callers-BGP.prv

Collective bytes — collective send bytes @ Gadget2-1024-300ms-1-5callers-BGP.prv

# p2p — nb point to point calls @ Gadget2-1024-300ms-1-5callers-BGP.prv

p2p bytes — Point-to-Point bytes sent @ Gadget2-1024-300ms-1-5callers-BGP.prv

p2p BW — Point-to-Point send bandwidth @ Gadget2-1024-300ms-1-5callers-BGP.prv

# Software counters



GADGET, PRACE Case A, 2048 procs

# Software counters



GADGET, PRACE Case A, 4096 procs