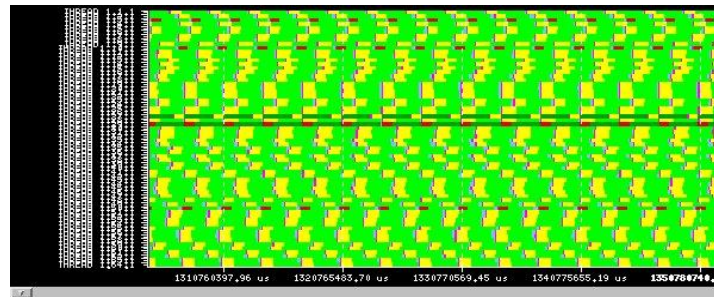# Humans are visual creatures

- Films or books?                                            PROCESS
    - Two hours vs. days (months)

- Memorizing a deck of playing cards                 STORE
    - Each card  translated to an image (person, action, location)

- Our brain loves pattern recognition                IDENTIFY
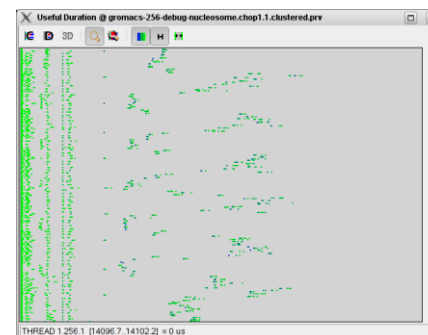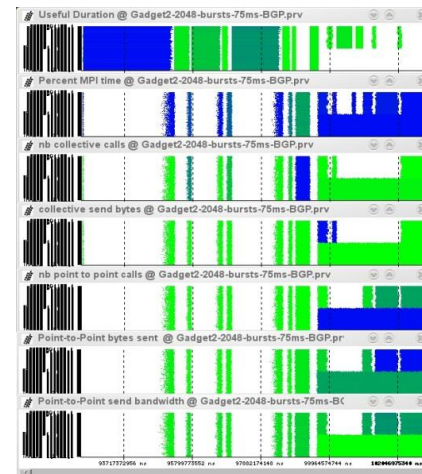    - What do you see on the pictures?

# Our Tools

- Since 1991

- Based on traces

- Open Source

- http://tools.bsc.es

- Core tools:

  - Paraver (paramedir) – offline trace analysis

  - Dimemas – message passing simulator

  - Extrae – instrumentation

- Focus

  - Detail, variability, flexibility

  - Behavioral structure vs. syntactic structure

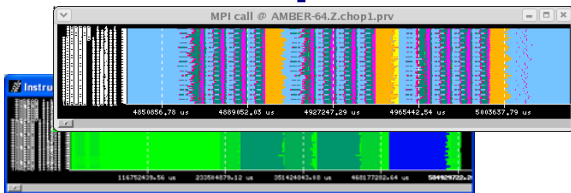  - Intelligence: Performance Analytics

# Paraver

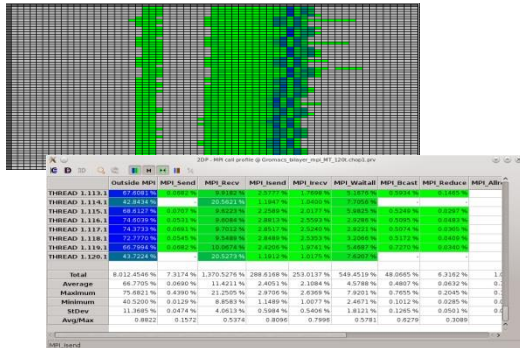# Paraver – Performance data browser

Raw data



**Timelines**

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

Trace visualization/analysis
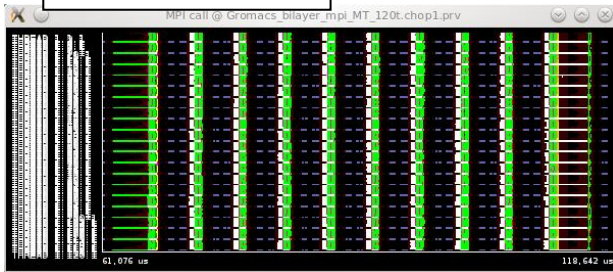
+ trace manipulation

Goal = Flexibility

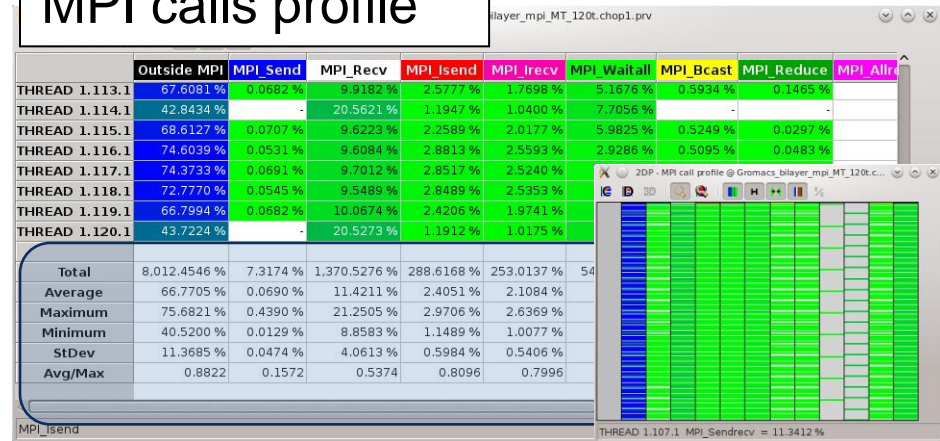No semantics

Programmable

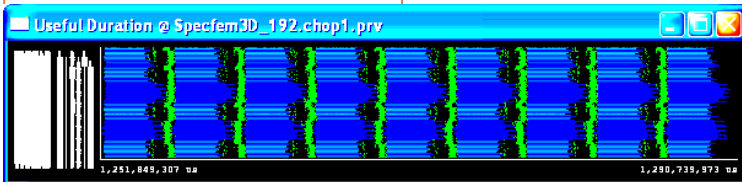Comparative analyses

Multiple traces

Synchronize scales

# From timelines to tables



MPI calls

MPI calls profile

Useful Duration

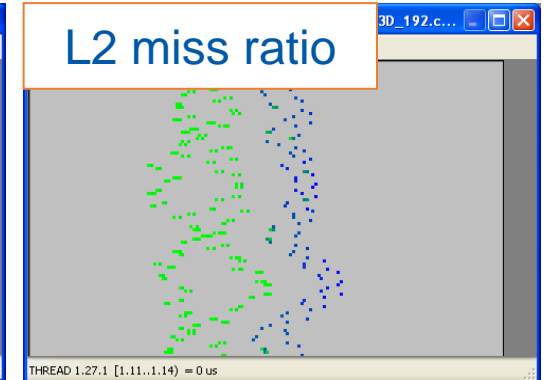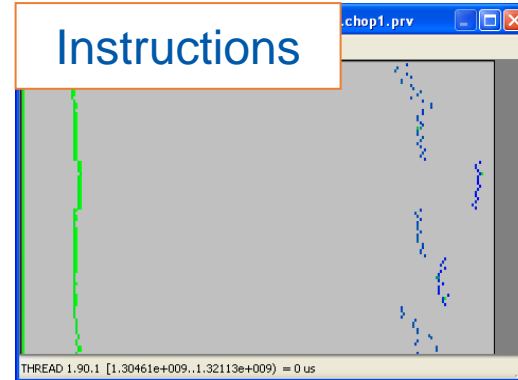Histogram Useful Duration

# Analyzing variability



Useful Duration

IPC

Instructions

L2 miss ratio

# Analyzing variability

- By the way: six months later ….

# From tables to timelines

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 ? ☺

# 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,...

570 s
2.2 GB
MPI, HWC

WRF-NMM
Peninsula 4km
128 procs

570 s
5 MB

4.6 s
36.5 MB

# Dimemas – Coarse grain, Trace driven simulation

- Simulation: Highly non linear model

    - MPI protocols, resource contention…

- Parametric sweeps
    - On abstract architectures
    - On application computational regions
- What if analysis
    - Ideal machine (instantaneous network)
    - Estimating impact of ports to MPI+OpenMP/CUDA/…
    - Should I use asynchronous communications?
    - Are all parts equally sensitive to network?
- MPI sanity check
    - Modeling nominal

- Paraver – Dimemas tandem
    - Analysis and prediction
    - What-if from selected time window





**Detailed feedback on simulation (trace)**

# Network sensitivity

- MPIRE 32 tasks, no network contention



L = 5µs – BW = 1 GB/s



L = 1000µs – BW = 1 GB/s



L = 5µs – BW = 100MB/s

**All windows same scale**

# Network sensitivity

- WRF, Iberia 4Km, 4 procs/node
  - Not 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)

# Would I will benefit from asynchronous communications?

SPECFEM3D



Courtesy Dimitri Komatitsch

Real

Ideal

Prediction MN

Prediction 100MB/s

Prediction 10MB/s

Prediction 5MB/s

Prediction 1MB/s

# Ideal machine

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

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



GADGET @ Nehalem cluster 256 processes

Allgather + sendrecv

allreduce

alltoall

sendrec

waitall

Real run

Ideal network

Impact on practical machines?

# Impact of architectural parameters

- **Ideal speeding up ALL** the computation bursts by the CPUratio factor

    - The more processes the less speedup (higher impact of bandwidth limitations) !!



GADGET

# Hybrid parallelization

- Hybrid/accelerator parallelization

  - Speed-up SELECTED regions by the CPUratio factor



Profile

% Computation Time

128 procs.

Code regions

Speedup

93.67%

Speedup

97.49%

Speedup

99.11%

Bandwidth (MB/s)

CPU ratio

Bandwidth (MB/s)

CPU ratio

Bandwidth (MB/s)

CPU ratio

(Previous slide: speedups **up to 100x**)

# Efficiency Model

# Parallel efficiency model



Computation   Communication

MPI_Send

MPI_Recv

Do not blame MPI

LB          Comm

- Parallel efficiency = LB eff * Comm eff



| | Outside MPI | MPI_Recv | MPI_Isend | MPI_Irecv |
|---|---|---|---|---|
| THREAD 1.130.1 | 87,95 % | 9,51 % | 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 | 67,62 % | 54,97 % | |
| Minimum | 30,67 % | 0,00 % | 0,00 % | |
| StDev | 15,27 % | 6,06 % | 21,40 % | 0,00 % |
| Avg/Max | 0,7 | 0,03 | 0,19 | 0,81 |

2DP - MPI call profile @ trace_24h_atmos_symbols.cho...

η

CommEff

LB

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Parallel efficiency refinement: LB * µLB * Tr



- Serializations / dependences (µLB)
- Dimemas ideal network → Transfer (efficiency) = 1

# Why scaling?

$$\eta_{\parallel} = LB * Ser * Trf$$

CG-POP mpi2s1D - 180x120

Good scalability !!
Should we be happy?



speed up



Parallel eff
LB
uLB
transfer



Efficiency
Parallel eff
instr. eff
IPC eff

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Why efficient?

Parallel efficiency =93.28
Communication = 93.84

Parallel efficiency = 77.93
Communication = 79.79

Parallel efficiency = 28.84
Communication eff = 30.42

# Analytics

# Using Clustering to identify structure



Completed Instructions

IPC

# What should I improve?

What if ….

PEPC

… we increase the IPC of Cluster1?

Aggregative Cluster Refinement

13% gain

… we balance Clusters 1 & 2?

19% gain

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Tracking scability through clustering

- OpenMX (strong scale from 64 to 512 tasks)

# Folding

- Instantaneous metrics with minimum overhead
  - Combine instrumentation and sampling
    - Instrumentation delimits regions (routines, loops, …)
    - Sampling exposes progression within a region
  - Captures performance counters and call-stack references

Initialization     Iteration #1     Iteration #2     Iteration #3     Finalization

Synth Iteration

# "Blind" optimization

- From folded samples of a few levels to timeline structure of "relevant" routines

Recommendation without access to source code

# CG-POP multicore MN3 study

- Unbalanced MPI application
- Same code
- Different duration
- Different performance



Instruction mix model for the unbalanced CGPOP on different cores of the same hexacore chip

ClusterID @ cgpop.linux_icc.180x120.chop2.clustered.prv

THREAD 1.13.1
THREAD 1.14.1
THREAD 1.15.1
THREAD 1.16.1
THREAD 1.17.1
THREAD 1.18.1

9.903.628.171 ns

# Methodology

# Performance analysis tools objective

## Help generate hypotheses

## Help validate hypotheses

Qualitatively

Quantitatively

Barcelona
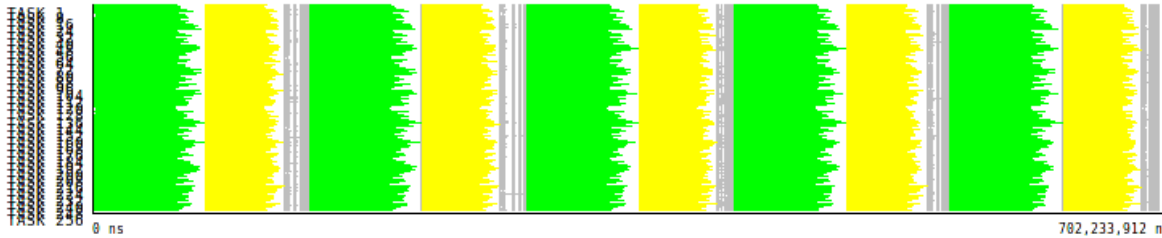Supercomputing
Center
Centro Nacional de Supercomputación

# 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, correlate with other counters

- Scalability – code replication?
  - Total #instructions

- Behavioral structure? Variability?

Paraver Tutorial:

Introduction to Paraver and Dimemas methodology

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

# BSC Tools web site

- tools.bsc.es
  - 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

# Demo

# Same code, different behaviour

| Code | Parallel efficiency | Communication eff. | Load Balance eff. |
|---|---|---|---|
| lulesh@mn3 | 90.55 | **99.22** | 91.26 |
| lulesh@leftraru | **69.15** | 99.12 | **69.76** |
| lulesh@uv2 (mpt) | 70.55 | 96.56 | 73.06 |
| lulesh@uv2 (impi) | 85.65 | 95.09 | 90.07 |
| lulesh@mt | 83.68 | 95.48 | 87.64 |
| lulesh@cori | 90.92 | 98.59 | 92.20 |
| lulesh@thunderX | 73.96 | 97.56 | 75.81 |
| lulesh@jetson | 75.48 | **88.84** | 84.06 |
| lulesh@claix | 77.28 | 92.33 | 83.70 |
| lulesh@jureca | 88.20 | 98.45 | 89.57 |
| lulesh@mn4 | 86.59 | 98.77 | 87.67 |
| lulesh@inti | 88.16 | 98.65 | 89.36 |
| lulesh@archer | 88.01 | 97.95 | 89.86 |
| lulesh@romeo | 89.56 | 99.01 | 90.45 |
| lulesh@mn4 | **91.02** | 98.38 | **92.52** |
| lulesh@ stampede2 (skl) | 85.76 | 97.63 | 87.84 |
| lulesh@ stampede2 (knl) | 89.21 | 98.42 | 90.64 |
| lulesh@isambard | 90.32 | 97.16 | 92.96 |

Warning::: Higher parallel efficiency does not mean faster!