



Barcelona Supercomputing Center Centro Nacional de Supercomputación

BSC Tools

Jesús Labarta, Judit Gimenez BSC

Moscow, Nov 27th 2012

Index

- (Performance tools
- (Extrae
- (Paraver
 - Philosophy
 - Timelines and tables
 - Finding needles in haystacks
 - Scalability
 - Scaling model
 - More examples
 - Analysis recommendations

(Dimemas

- Concept and uses
- Using Dimemas
- Scaling model

(Performance analytics

- Spectral analysis
- Clustering
- Folding
- (Tareador
- (Additional slides
 - Paraver internals



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









BSC – tools framework



Performance analysis tools objective

Help generate hypotheses

Help validate hypotheses

Qualitatively

Quantitatively







Barcelona Supercomputing Center Centro Nacional de Supercomputación



Jesús Labarta, Judit Gimenez BSC

Moscow, Nov 27th 2012

Extrae

(Parallel programming model runtime

- MPI, OpenMP, pthreads, OmpSs, CUDA, MIC...
- (Counters
 - CPU counters
 - Using PAPI and PMAPI interfaces
 - Network counters
 - OS counters
- (Link to source code
 - Callstack at MPI
 - OpenMP outlined routines and their containers
 - User functions selected
- (Periodic samples

(User events



How does Extrae intercepts your app?

(LD_PRELOAD

- Specific libraries for each combination of runtimes

- MPI
- OpenMP
- OpenMP+MPI
- ...

(Dynamic instrumentation (not available in Graphit)

- Based on DynInst (developed by U.Wisconsin/U.Maryland)
 - Instrumentation in memory
 - Binary rewriting

(Other possibilities

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



How to use Extrae?

(Build normal production binary of your app.

(Adapt job submission script

(If special features required select/adapt .xml configuration file

- Bunch of examples distributed in the package
 - Look at \$EXTRAE_HOME/share/example
- (Run it and get the trace!!



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



Trace control .xml





Trace control .xml (cont)





Trace control .xml (cont) mpitrace.xml (cont) ... <trace-control enabled="yes"> <file enabled="no" frequency="5m">/qpfs/scratch/bsc41/bsc41273/control</file> <global-ops enabled="no"></global-ops> <remote-control enabled="no"> <signal enabled="no" which="USR1"/> </remote-control> **External activation of tracing** </trace-control> (creation of file will start tracing) Stop tracing after elapsed time ... <others enabled="no"> <minimum-time enabled="no">10M</minimum-time> <terminate-on-signal enabled="no">USR2</terminate-on-signal> </others> ... or when signal received ... Barcelona Supercomputing Center ntro Nacional de Supercomputación

Trace control .xml (cont)





Trace control .xml (cont)





LD_PRELOAD library selection

(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





Barcelona Supercomputing Center Centro Nacional de Supercomputación

Paraver

Jesús Labarta, Judit Gimenez BSC

Moscow, Nov 27th 2012

"That what is simple is rarely understood"

my iPads Shangai cookies



Barcelona Supercomputing Center Centro Nacional de Supercomputación

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







Instruments

- (Lots of analysis
- (To obtain sufficient information/insight
 - Avoid flying blind
 - Identification of productive next steps









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





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)
- Need a language to describe how to compute such functions of time.
 - 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



Timelines

(Each window displays one view

- Piecewise constant function of time

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

- (Types of functions
 - Categorical
 - State, user function, outlined routine
 - Logical
 - In specific user function, In MPI call, In long MPI call

- Numerical

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



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

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

 $S_i \in R$

Timelines

(Representation

- Function of time

L2 miss ratio @ bt.4.mpi.prv	JX
THREAD 1.1.1	
	_
THREAD 1.3.1	-1
	1
16291582.07 us 16448220.79 us 16604859.51 us 16761498.23 us 16918136.95 us	

- Colour encoding



- Not null gradient
 - Black for zero value
 - Light green \rightarrow Dark blue



(Non linear rendering to address scalability



Basic functions of time

- (The filter module presents a subset of the trace to the semantic module. Each thread h 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



Basic functions of time

Semantic module Semantic module (From communication records to functions of time (From Events to functions of time - Last event value $S(i) = V(Ev_i)$ Send Bytes $s(t) - \sum Sz(C_j), j \mid (T_z(C_j) < t) \land (T_z(Cj) > t) \land (Dir(Cj) - send)$ - Next event value $S(t) = V(Ev_{co})$ $s(t) = \sum_{j} \frac{Sz(C_j)}{T_z(C_j) - T_z(C_j)}, j \mid (T_z(C_j) < t) \land (T_z(Cj) > t) \land (Dir(Cj) - send)$ Send Bandwidth - Average Next Event Value $S(i) = \frac{V(E_{Y_{i,i}})}{T(E_{Y_{i,i}}) - T(E_{Y_{i,i}})}$ Msgs in transit $s(t) = \sum sign(j), j \mid (T_s(C_j) < t) \land (T_s(Cj) > t) \land (Dir(Cj) \longrightarrow send)$ - Interval btw. Events $S(t) = T(Ev_{ab}) - T(Ev_{bb})$ Recv. Bandwidth $s(t) = \sum_{i} \frac{Sz(C_i)}{T_x(C_i) - T_x(C_i)}, j | (T_x(C_j) < t) \land (T_x(C_j) > t) \land (Dir(C_j) - recv)$ - Rec. Negative Msgs $s(t) = \sum sign(j), j \mid (T_s(C_j) < t) \land (T_s(Cj) > t) \land (Dir(Cj) = -recv)$ Comm. Partner $s(t) = Partner(C_j), j \mid (T_s(C_j) < t) \land (T_s(C_j) > t)$ Bytes btw. Events $S(t) = \sum Sz(C_j), j \mid T_z(C_j) \in \left[T(Ev_i), T(Ev_{i-1})\right] \lor T_z(C_j) \in \left[T(Ev_i), T(Ev_{i-2})\right)$ OSC Receiver (050 Animitan Animitan Conter



Tables: Profiles, histograms, correlations

(Huge number of statistics computed from timelines



Tables: Profiles, histograms, correlations

(By the way: six months later





How to read profiles

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

🗙 MPI profile @	🥑 lberia-128-0	A.chop1.1it.	shifted.prv			
i c id 3d 🔍	🔍 🛛 🔳 н	• • • • •				
	End	MPI_lsend	MPI_lrecv	MPI_Wait	MPI_Allreduce	MPI_Com
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 %	
1						•



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**



3D

Barcelona

Center

(An additional control dimension

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


Tables

Tables

2D analysis module

- (Single flexible quantitative analysis mechanism
- (Let
 - cw1 and cw2 two views we will call control views
 - dw a view we will call data window
- (For each control window we define a set of bins $bin_{i}^{cw} = range_{i}^{cw}, range_{i+1}^{cw}$ $range_{i+1}^{c_{w}} = range_{i}^{c_{w}} + delta^{c_{w}}$
- (And the discriminator functions
 - $\delta_{t}^{\sigma_{w}}(t) = ((S^{\sigma_{w}}(t) \in bin_{t}^{\sigma_{w}})?1:0)$ $\delta_{i,z}(t) = \delta_i^{coe}(t) * \delta_z^{coe}(t)$
- Identify regions with cw's within the (j,k) bin

For each window w

- (The 3D analysis module computes a cube (or plane in the case of 2D) of statistics $M(thread, j, k) = statistic(S_{4k}^{dw}(t) * \delta_{4k+1}(t))$
- Where the statistic can represent the average value, the number of intervals,....



See slides at end of presentation for details



Environmentations Approximations Contains Contains

Paraver: finding needles in haystacks



Barcelona Supercomputing Center Centro Nacional de Supercomputación

Comparative analyses

- (Possible to load several traces
- (Copy and paste
 - right click menu
 - From one window to another: time, duration, size, objects displayed,...
 - Time between windows and tables: analysis will be computed only for the selected time interval



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?





Scalability



Barcelona Supercomputing Center Centro Nacional de Supercomputación

Scalability of Presentation



Scalability of analysis



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



Barcelona Supercomputing Center Centro Nacional de Supercomputaciór J. Labarta, et al.: "Scalability of tracing and visualization tools", PARCO 2005

Software counters







Software counters



Barcelona Supercomputing

Centro Nacional de Supercomputación

Center

BSC



Useful Duration @ Gadget2-2048-150ms-1-5callers-BGP.prv

Software counters



GADGET, PRACE Case A, 4096 procs





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







Analyze coarse grain structure

- (Filtered trace is a Paraver trace
- (Can be analyzed with standard cfgs as long as the information they require is still in the trace
 - A typical view that shows a lot of the structure of a trace is useful_duration.cfg
 - Repetitive structure is often apparent
 - Perturbations can also be typically identifed
 - A clean/representative interval can be identified

📕 Useful Du	Useful Duration_z1_z2 @ trace.prv.CPUDurBurst.filtered.prv				
	642718757.08 us	1272710410.06 us	1909065615.09 us	2545420820.12 us	3181776425.1





Scaling model



Barcelona Supercomputing Center Centro Nacional de Supercomputación

Presenting application performance

- (Factors modeling parallel efficiency
 - Load balance (LB)
 - Micro load balance (µLB) or serialization
 - Transfer

$$= LB * \mu LB * Transfer$$

CommEff

(Factors describing serial behavior

- Computational complexity: #instr
- Performance: IPC

$$T_{Comp} = \frac{\#instr}{IPC}$$

(Scaling model

ro Nacional de Supercomputació

Barcelona Supercomputing Center

$$Sup = \frac{P}{P_0} * \frac{\eta}{\eta_0} * \frac{IPC}{IPC_0} * \frac{\#instr_0}{\#instr}$$

η



M. Casas et al, "Automatic analysis of speedup of MPI applications". ICS 2008.

Scaling model

 $\eta = LB * CommEff$

Directly from real execution metrics



55



Parallel Performance Models

$$\eta = LB * \mu LB * Transfer$$







(Old \rightarrow new version: -43% instr, -52% time



Examples



Barcelona Supercomputing Center Centro Nacional de Supercomputación





Barcelona Supercomputing Center Centro Nacional de Supercomputación

Comparison between versions

(**1** 24 processes, 180x120

х MPI call @ mpi2s1D.180x120.clustered.prv 1D 38.531.820 us 38.884.807 u MPI call @ mpi2s1D_overlap.180x120.clustered.prv THREAD 1.1.1 1D overlap 31.553.754 us 31,906,741 x MPI call @ mpi2s2D.180x120.clustered.prv 2D 32.485.935 us 32.838.922 us Barcelona Supercomputing BSC Center Centro Nacional de Supercomputación

Version	η	Time LB	Comm
1D	0.76	0.77	0.99
1D-ovlp	0.75	0.76	0.99
2D	0.93	0.98	0.95

Version	IPC	Comp LB
1D	0.32	0.75
1D-ovlp	0.34	0.75
2D	0.17	0.99

- Load Imbalance
 - Significant improvement in 2D ...
- Poor IPC
 - …compensating Load balance improvements in 2D ☺

CGPOP: MPI connectivity

- (Who sends to whom
- (Leverage full capabilities of Paraver table mechanism
 - Control window: To whom I am sending







Short diagnosis: FrontFlowRed – 256 vs 512 tasks

- Unbalance and MPI time increases
- Sequence of MPI_Allreduce calls
- Poor IPC, high cache misses



Modifications in the source code

- Improved load balance weight nodes based on their connectivity, balance total weight
- Reduction of all to all communications half of MPI_Allreduce calls were eliminated
- . Reduce cache misses reordering node number

Half an hour meeting for the analysis, few days of work to modify code \rightarrow 25% of improvement



256 tasks – Load balance



256 tasks – MPI time





	End	MPI_lsend	MPI_irecv	MPI_Waitall	MPI_Allreduce
Total	7.617,81 %	2.432,05 %	1.772,21 %	2.245,19 %	11.532,74 %
Average	29,76 %	9,50 %	6,92 %	8,77 %	45,05 %
Maximum	51,85 %	18,11 %	13,00 %	28,59 %	67,69 %
Minimum	17,59 %	1,65 %	1,37 %	1,46 %	14,88 %
StDev	8,04 %	3,33 %	2,39 %	6,05 %	11,03 %
Avg/Max	0,57	0,52	0,53	0,31	0,67

	End	MPI_Isend	MPI_irecv	MPI_Waitall	MPI_Allreduce
Total	9.089,05 %	3.100,13 %	2.107,96 %	1.511,02 %	9.791,83 %
Average	35,50 %	12,11 %	8,23 %	5,90 %	38,25 %
Maximum	51,05 %	20,77 %	14,59 %	21,31 %	64,34 %
Minimum	22,99 %	3,47 %	3,21 %	1,72 %	10,90 %
StDev	5,92 %	3,62 %	2,61 %	3,89 %	10,13 %
Avg/Max	0,70	0,58	0,56	0,28	0,59



256 tasks – Memory access



Analysis recommendations



Barcelona Supercomputing Center Centro Nacional de Supercomputación

Navigation

(Both for timelines and tables

(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.

(Third: attitude

(Top down: Do not get obsessed about details!!!! Look at the forest, not the trees !!!



Analysis

(Loop:

- Question/hypothesis
- What metric I need to quantify/validate/reject ?
- What view (timeline or table) would provide the information to answer the question?
- Do I have a configuration file to show such view?
 - If not, how do I generate it?

(Considerations:

- May be useful to identify a proper region where to apply the analysis
 - Avoid initialization, perturbed regions, focus on a couple of iterations,....
- The more iterations you can make out of a single trace the more understanding you will get !!!



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



Analysis of MPI programs

- (Typical initial questions:
 - What is the global efficiency, Load balance, and communication efficiency?
 - If efficiency is good, still need to look at sequential performance.
 - Is IPC good? Everywhere?
 - If bad, is if due to cache misses? Other?
 - If imbalanced
 - is it due to computation or IPC imbalance?
 - If communication efficiency low
 - are message sizes large?
 - Too many messages?
 - Is the effective bandwidth achieved reasonable? For all transfers?
 - Is there serialization of the computations?


Paraver Summary



Barcelona Supercomputing Center Centro Nacional de Supercomputación

Summary

(Performance tools are needed more and more !!!!!

- To tune our applications, to design our system software.
- To understand what really happens, how our systems really behave,...
- (Paraver, ...
 - Extremely flexible an powerful.
 - "professional" equipment
 - Attitude: need curiosity to understand, confidence on measurement capability of the tool.
 - A Personal Challenge: squeeze information from the data

I have seen things you people wouldn't believe...

Roy Batty – Blade Runner

Seeing is believing ... measuring is better



Barcelona Supercomputing Center Centro Nacional de Supercomputaciór

Free adaptation of a Spanish saying

Tools web site

- II <u>www.bsc.es/paraver</u>
 - downloads
 - Sources / Binaries
 - documentation
 - Training guides
 - Tutorial slides

(Assignment ©

- get from Graphit: /export/hopsa/BSCtools/wxparaver_pkgs
 - Paraver (binaries) for your windows/linux/MAC laptop
- get tutorials from Graphit: /export/hopsa/BSCtools/tutorials
- Start wxparaver
- Help \rightarrow tutorials and follow instructions
- Follow tutorials 1 (Introduction to Paraver and MPI) and 2 (Tutorial on HydroC)







Barcelona Supercomputing Center Centro Nacional de Supercomputación

Dimemas

Jesús Labarta, Judit Gimenez BSC

Moscow, Nov 27th 2012

BSC – tools framework



Dimemas: Coarse grain, Trace driven simulation

(Simulation: Highly non linear model

- Linear components
 - Point to point communication
 - Sequential processor performance
 - Global CPU speed
 - Per block/subroutine
- Non linear components
 - Synchronization semantics
 - Blocking receives
 - Rendezvous
 - Resource contention
 - CPU
 - Communication subsystem
 - » links (half/full duplex), busses





P2P communication model



Dimemas

(A model of a hypothetical machine

- Reference
- Identification of importance of different factors

(Use examples:

- Are all parts of an app. equally sensitive to network?
- Ideal machine
- Architecture parametric sweeps
- Estimating impact of ports to MPI+OpenMP/CUDA/...
- Endpoint contention
- Core architecture vs. system architecture trade offs



Are all parts of an app. equally sensitive to network?

MPIRE 32 tasks, no network contention



Ideal machine

- (The impossible machine: $BW = \infty$, L = 0
- (Actually describes/characterizes Intrinsic application behavior
 - Load balance problems?
 - Dependence problems?



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) !!



The potential of hybrid/accelerator parallelization





Intrinsic application behavior

- (End point contention
 - Simulation with Dimemas
 - Very low BW
 - 1 output link, ∞ input links
- (Recommendation:
 - Important to schedule communications.

Everybody sending by destination rank order Endpoint contention at low ranked processes





Multiscale Simulation



Multiscale simulation: L2 size vs network bw

- (Left: clusters IPC with different cache sizes
 - 64KB 512MB
- (Right: execution time with different network bw and cache size
 - 125Mb/s 500Mb/s

((NAS BT

- Can compensate cache reduction with more network bw
- ((VAC, WRF
 - Dominated by comp.





Using Dimemas



Barcelona Supercomputing Center Centro Nacional de Supercomputación



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

- Internal collective operations _ 🗆 ×			
Machine number <<< 1 >>>			
COLLECTIVE OP.	FAN OUT		
Name	Model Size	Model	I
MPI_Barrier	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Bcast	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Gather	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Gatherv	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Scatter	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Scatterv	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Allgather	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Allgatherv	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Alltoall	LIN 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Alltoallv	LIN 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Reduce	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Allreduce	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Reduce_Scatter	LOG 🔻 MAX	▼ 0 ▼ MAX	-
MPI_Scan	LOG 🔻 MAX	▼ 0 ▼ MAX	-
Apply to all:	Select 🔻 Select	▼ Select ▼ Sele	ct 🔻
Save Do all the same Close			



Scaling model



Barcelona Supercomputing Center Centro Nacional de Supercomputación

Presenting application performance

- (Factors modeling parallel efficiency
 - Load balance (LB)
 - Micro load balance (μLB) or serialization
 - Transfer

$$\eta = LB * \mu LB * Transfer$$

CommEff

(Factors describing serial behavior

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

$$T_{Comp} = \frac{\#instr}{IPC}$$

(Scaling model

ro Nacional de Supercomputació

$$Sup = \frac{P}{P_0} * \frac{\eta}{\eta_0} * \frac{IPC}{IPC_0} * \frac{\#instr_0}{\#instr}$$



enter

M. Casas et al, "Automatic analysis of speedup of MPI applications". ICS 2008.

Scaling model

 $\eta = LB * CommEff$

Directly from real execution metrics



96









Barcelona Supercomputing Center Centro Nacional de Supercomputación

Performance Analytics

Jesús Labarta, Judit Gimenez BSC

Moscow, Nov 27th 2012

BSC – tools framework



Spectral analysis



Barcelona Supercomputing Center Centro Nacional de Supercomputación

(... in Paraver are functions of time

(Natural target for signal processing techniques

(Relevant functions of time at global application level

- # processes in MPI, outside MPI, ...
- Sum (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

(Useful

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





Signal processing applied to performance analysis

(Hierarchical structure identification





Scalability: online automatic interval selection



Supercomputing

Centro Nacional de Supercomputación

Center

BSC



Clustering



Barcelona Supercomputing Center Centro Nacional de Supercomputación

Clustering: analysis of performance @ serial computation bursts

(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





Clustering: analysis of performance @ serial computation bursts

(**Good** clustering algorithms ?

- Data not necessarily Gaussian \rightarrow K-means type of algorithms not very good
- Density based algorithms "better" —
 - DBSCAN

5e+008

4.5e+008

4e+008

3.5e+008

3e+008

2.5e+008

2e+008

1.5e+008

1e+008

Completed Instructions



Sup Cei BSC

J. Gonzalez et al, "Automatic Detection of Parallel Applications Computation" Phases. (IPDPS 2009)

Performance @ serial computation bursts



DBSCAN (Eps=0.015, MinPoints=10)
Automatic clustering quality assessment

- (Leverage Multiple Sequence alignment tools from Life Sciences
- Process == Sequence of clusters \leftrightarrow sequence of amino acids == DNA
- CLUSTAL W, T-Coffee, Kalign2
- Cluster Sequence Score (0..1)
- (Per cluster / Global
 - Weighted average





Quality driven clustering algorithms



Hierarchical clustering guided by SPMDness

(Automatic improvement of clustering "quality"





Using clustering

(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



. . .

Using clustering: Identifying code regions and variability



instr. vs. cluster















Using clusters: to understand apps behavior



Using clusters: HWC projection

(Complete hardware counter characterization

Precise Statistics (>> multiplexing)

PAPI_TOT_INS PAPI L2 DCM PAPI FP INS

Formula

Neural Net

+ Cach

111

- CPI stack model



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



CPI Stack Modelization





J. Gonzalez et all, "Performance Data Extrapolation in Parallel Codes", ICPADS 2010

Centro Nacional de Supercomputación

Folding



Folding: Instrumentation + sampling

- (Extremely detailed time evolution of hardware counts, rates and callstack
- (Minimal overhead
- (Based on
 - Instrumentation events (iteration, MPI, ...) and periodic samples.
 - Application structure: manual iteration instrumentation, routines, clusters
- (Folding
 - Post processing to project all samples into one instance



Sampling and Folding

(Instructions

(L1 data cache misses





Sampling and Folding





Folding \rightarrow profiles of rates and ratios (Call-site sampling information is folded Correlation between hwc and call-sites - GVIM add-on to show performance within source code Timeless but useful to point performance issues Edita Eines Sintaxi Buffers Finestra Ajuda Folded source code line loop over all cells owned by this node THREAD 1.1.1 21 do c = 1, ncells 23 C---24 c compute the reciprocal of density, and the kinetic energy 25 c and the speed of sound. rmΠ do k = -1, cell_size(3,c) 479.707.396 ns 495.315.648 ns do j = -1, cell_size(2,c) cell size(1.) (k,c) = rho invTHREAD 1.1.1 34 u(2,i,j,k,c)*u(2,i,j,k,c) + u(3,1,j,k,c)*u(3,1,j,k,c) + 38 u(4, i, j, k, c)*u(4, i, j, k, c)) * rho_inv 495.315.648 ns 479.707.396 ns 40 enddo 41 enddo **Folded instructions** 42 enddo 42,1 4% H. Servat et al. "Unveiling Internal Evolution of Parallel Application Computation Phases" ICPP 2011 Barcelona Supercomputing Center

Centro Nacional de Supercomputación

Examples



Instantaneus hardware counter rates & CPI stack models

(Normalized rates

 $M(c,t) = \frac{Counter _rate(c,t) / max(Counter _rate(c,t))}{MIPS(t)}$

(Different compilers and architectures \rightarrow slightly different characteristics



Instantaneus hardware counter rates & CPI stack models



- ((180 x 120
- (24 processes
- ((IBM PPC-970
- (Three main program regions







CGPOP

1D vs. 2D





CGPOP

Sequential code optimization



Clustering + Folding + CPI Stack

- (Unmodified production binary
 - Instrumentation of MPI calls plus periodic samples.
 Rotating hardware counters.
 - Clustering detects structure to enable folding
 - Folding of multiple hardware counters to compute CPI stack model

(Instantaneous MIPS and CPI Stack model

- Clean abstract identification of performance problem
- Good identification of phases

╘ ଅ ▦ ଵ ଵ ଵ ၖ ?







Instantaneous CPI stack



Instantaneous CPI stack



Tareador



Predicting performance of task based models

(Performance prediction

- Predict MPI/StarSs multithreaded from a pure MPI run
- Leveraging other tools in environment



Predicting performance

- (Potential concurrency between task (for a suggested taskification)...
- (... as number of cores increases.





Predicting performance









Conclusion



Conclusion

(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



Conclusion

- BSC tools: Extreme flexibility and analysis power
 - Maximize insight obtained form single experiment
 - Learning curve
 - "Don't ask whether something can be done, ask how can it be done"
- Huge:
 - Needs for such detailed and precise analysis
 - Systems are complex, variability is everywhere, ...
 - Insight, not speculation to fly in the midst and "maximize" productivity
 - Potential of data analysis techniques applied to performance data

(Use them, stay tuned:

www.bsc.es/paraver





Barcelona Supercomputing Center Centro Nacional de Supercomputación

THANKS

Detailed material



Semantic Module



Basic functions of time

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



Filter module







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



Semantic module

- (From communication records to functions of time
 - Send Bytes
 - Send Bandwidth
 - Msgs in transit

Recv. Bandwidth

Rec. Negative Msgs

$$s(t) = \sum_{j} \frac{Sz(C_{j})}{T_{R}(C_{j}) - T_{S}(C_{j})}, j \mid (T_{S}(C_{j}) < t) \land (T_{R}(C_{j}) > t) \land (Dir(C_{j}) = send)$$

 $s(t) = \sum_{j} Sz(C_j), j \mid (T_s(C_j) < t) \land (T_R(C_j) > t) \land (Dir(C_j) = send)$

$$s(t) = \sum_{j} sign(j), j \mid (T_{s}(C_{j}) < t) \land (T_{R}(C_{j}) > t) \land (Dir(C_{j}) = send)$$

$$s(t) = \sum_{j} \frac{Sz(C_{j})}{T_{R}(C_{j}) - T_{S}(C_{j})}, j \mid (T_{S}(C_{j}) < t) \land (T_{R}(C_{j}) > t) \land (Dir(C_{j}) = recv$$

$$s(t) = \sum_{j} sign(j), j \mid (T_R(C_j) < t) \land (T_S(C_j) > t) \land (Dir(C_j) = recv)$$

 $s(t) = Partner(C_j), j \mid (T_s(C_j) < t) \land (T_R(C_j) > t)$

– Comm. Partner

- Bytes btw. Events

$$S(i) = \sum_{i} Sz(C_{j}), j | T_{S}(C_{j}) \in [T(Ev_{i}), T(Ev_{i+1})] \lor T_{R}(C_{j}) \in [T(Ev_{i}), T(Ev_{i+1})]$$


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 ° ls equal S'(t) = sign (S(t) = a ? S(t) : 0)
 - Delta $S'(t) = S_{i+1} S_i$
 - Stacked value



Semantic module



Semantic module

(Derived windows

- Point wise operation

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





(Thread function: State as is



Useful for

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







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





• Hwc events (TLB, L1 misses,...) within interval





• Useful for

• Hwc events (TLB, L1 misses,...) per time unit within interval





Semantic module





Barcelona Supercomputing Center Centro Nacional de Supercomputación

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

Process view



IIIS mantic mod

- 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, ...





IIIS mantic mod

• 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



Barcelona Supercomputing Center Centro Nacional de Supercomputación

Tables

- Single flexible quantitative analysis mechanism
- Let
 - cw₁ and cw₂ two views we will call control views
 - dw a view we will call data window
- (For each control window we define a set of bins

$$bin_{j}^{cw} = \left[range_{j}^{cw}, range_{j+1}^{cw}\right]$$
 $range_{j+1}^{cw} = range_{j}^{cw} + delta^{cw}$

And the discriminator functions

$$\delta_{j}^{cw}(t) = ((S^{cw}(t) \in bin_{j}^{cw})?1:0)$$

$$\delta_{j,k}(t) = \delta_{j}^{cw_{1}}(t) * \delta_{k}^{cw_{2}}(t)$$

For each window w

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

$$range_{j+1}^{cw} = range_{j}^{cw} + delta^{cw}$$

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



entro Nacional de Supercomputación

How to ...



Barcelona Supercomputing Center Centro Nacional de Supercomputación

Main Paraver window



Load configuration files



Navigation





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



Table information and control



Table information and control

