

BSC Tools Hands-On

Lau Mercadal, Germán Llort (tools@bsc.es) Barcelona Supercomputing Center





Getting a trace with Extrae



Extrae Features

Platforms

- Intel, Cray, BlueGene, Intel MIC, ARM, Android, Fujitsu Sparc, ...
- Parallel programming model
 - MPI, OpenMP, pthreads, OmpSs, CUDA, OpenCL, Java, Python, ...
- Performance counters
 - Using PAPI interface
- Link to source code
 - Callstack at MPI routines
 - OpenMP outlined routines
 - Selected user functions
- Periodic samples
- User events (Extrae API)

No need to recompile or relink!

Extrae Overheads

	Average values	Jureca Cluster
Event	150 – 200ns	123ns
Event + PAPI	750 – 1000ns	920ns
Event + callstack (1 level)	1µs	1.439µs
Event + callstack (6 levels)	2µs	3.447µs

How does Extrae work?

- Symbol substitution through LD_PRELOAD
 - Specific libraries for each combination of runtimes
 - MPI, OpenMP, MPI+OpenMP, ...

Dynamic instrumentation

- Based on DynInst (developed by U.Wisconsin/U.Maryland)
 - Instrumentation in memory
 - Binary rewriting
- Static link (i.e., PMPI, Extrae API)



Using Extrae in 3 steps

- **1. Adapt** your job submission script
 - Append Extrae loader script

2. Configure what to trace

- Modify extrae.xml
- **3. Run** it!
- For further reference check the Extrae User Guide
 - <u>https://tools.bsc.es</u> -> Documentation -> Tools manuals
 - Also distributed with Extrae in \$EXTRAE_HOME/share/doc

Using Extrae in JURECA

```
laptop> ssh <USER>@jureca.fz-juelich.de
jrlXX> cp -r tools-material $HOME
jrlXX> ls $HOME/tools-material
    apps/
    clustering/
    extrae/
    slides/
    traces/
```







Choose depending on application type

Library	Serial	MPI	OpenMP	pthread	CUDA
libseqtrace	\checkmark				
libmpitrace[f] ¹		\checkmark			
libomptrace			\checkmark		
libpttrace				\checkmark	
libcudatrace					\checkmark
libompitrace[f] ¹		\checkmark	\checkmark		
libptmpitrace[f] ¹		\checkmark		\checkmark	
libcudampitrace[f] 1		\checkmark			\checkmark

¹ Append "f" suffix for Fortran codes

Step 3: Run it!

Submit your job & check status

jrlXX> cd \$HOME/tools-material/extrae

jrlXX> sbatch job.slurm

jrlXX> squeue -u \$USER

- Once finished the trace will be in the same folder
 - lulesh2.0_27p.{pcf,prv,row}
- Any issue?
 - Traces already generated in \$HOME/tools-material/traces

Step 2: Configure what to trace



Step 2: Configure what to trace

jrlXX> vi \$HOME/tools-material/extrae/extrae.xml

```
<counters enabled="yes">
        <cpu enabled="yes" starting-set-distribution="1">
            <set enabled="yes" domain="all" changeat-time="0">
                 PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_L1_DCM, PAPI_L2_DCM
            </set>
        </cpu>
        <network enabled="no" />
        <resource-usage enabled="no" />
        <memory-usage enabled="no" />
        </counters>
```

Select which HW counters are measured

(How's the machine doing?)

Step 2: Configure what to trace





Installing Paraver & First analysis steps



Install Paraver in your laptop



Install Paraver tutorials



Uncompress, rename & move

Paraver

laptop> tar xf wxparaver-4.8.1-Linux_x86_64.tar.bz2

laptop> mv wxparaver-4.8.1-Linux_x86_64 paraver

Tutorials

laptop> tar xf paraver-tutorials-20150526.tar.bz2

laptop> mv paraver-tutorials-20150526 paraver/tutorials

Tell Paraver where to find the tutorials

Start Paraver

laptop> \$HOME/paraver/bin/wxparaver &

• Check the tutorials are available



View full pa	th in trace selector able trace size (MB) 500	÷	
Default directo	ries		
Traces	/home/emercada	Browse	
CFGs	/home/emercada/soft/wxparaver/4.8.1/cfgs	Browse	Browse to
Filters XML	/home/emercada/soft/wxparaver/4.8.1/share/filters-config	Browse	
Tutorials root	/home/emercada/soft/wxparaver/4.8.1/tutorials	Browse	
Tmp dir	/home/emercada	Browse	
Behaviour —			

Check that everything works

• Check the tutorials are available





Also available in JURECA

jrlXX> module load Paraver/4.8.1

First steps of the analysis

Copy the trace to your laptop

laptop> scp <USER>@jureca:\$HOME/tools-material/extrae/lulesh2.0_27p.{prv,pcf,row} ./



Measure the parallel efficiency

• Click on mpi_stats.cfg

rials		Outside MPI	MPI_Isend	MPI_Irecv	MPI_Wait	MPI_Waitall	MPI_Barrier	MPI_Reduce
a first question to answer when analyzing a narallel code is "how afficient does it	THREAD 1.19.	93.54 %	0.82 %	0.02 %	0.72 %	0.60 %	0.02 %	0.38 %
1?". The efficiency of a parallel program can be defined based on two aspects: the real-based on two aspects: the real-based on two aspects: the	THREAD 1.20.	90.08 %	0.75 %	0.03 %	0.78 %	0.58 %	0.02 %	0.79 %
gions. These two metrics would be the first checks on the proposed methodology.	THREAD 1.21.	88.97 %	0.71 %	0.02 %	1.39 %	0.54 %	0.01 %	0.85 %
• To measure the parallel efficiency load the configuration file	THREAD 1.22.	90.63 %	0.69 %	0.03 %	1.14 %	0.08 %	0.02 %	0.64 %
every thread spends in every MPI call. Look at the global statistics at the bottom of	THREAD 1.23.	96.40 %	0.74 %	0.04 %	0.91 %	0.12 %	0.06 %	0.16 %
efficiency, entry Avg/Max represents the global load balance and entry Maximum	THREAD 1.24.	93.46 %	0.64 %	0.03 %	2.22 %	0.01 %	0.02 %	0.37 %
represents the communication enciency. If any of mose values are lower than 85% is recommended to look at the corresponding metric in detail. Open the	THREAD 1.25.	96.07 %	0.39 %	0.02 %	1.83 %	0.01 %	0.01 %	0.00 %
control window to identify the phases and iterations of the code.	THREAD 1.26.	94.16 %	0.31 %	0.02 %	1.93 %	0.01 %	0.02 %	0.26 %
 To measure the computation time distribution load the configuration hile <u>cfgs/general/2dh_usefulduration.cfg</u> This configuration pops up a histogram of the duration for the computation regions. The computation regions are delimited by the exit from an MPI call and the entry to the next call if the 	THREAD 1.27.	74.70 %	0.25 %	0.01 %	12.43 %	0.01 %	0.01 %	1.26 %
histogram does not show vertical lines, it indicates the computation time may be not balanced. Open the control window to look at the time distribution and visually correlate hoth views.	Total	2,458.55 %	15.64 %	0.87 %	64.18 %	3.15 %	0.63 %	12.93 %
	Average	91.06 %	0.58 %	0.03 %	2.38 %	0.12 %	0.02 %	0.48 %
	Maximum	98.10 %	1.13 %	0.08 %	12.43 %	0.60 %	0.06 %	1.26 %
The computation regions	Minimum	74.70 %	0.05 %	0.01 %	0.72 %	0.01 %	0.00 %	0.00 %
distribution of the instruct	StDev	5.31 %	0.25 %	0.01 %	2.78 %	0.18 %	0.01 %	0.31 %
To measure the serial Comm efficiency	Avg/Max	0.93	0.51	0.42	0.19	0.19	0.40	0.38
					-			

MPI call profile @ lulesh2.0_27p.prv

Computation time and work distribution

■ Click on 2dh_usefulduration.cfg (2nd link) → Shows time computing

Tutorials
The first question to answer when analyzing a parallel code is "how efficient does it run?". The efficiency of a parallel program can be defined based on two aspects: the parallelization efficiency and the efficiency obtained in the execution of the serial regions. These two metrics would be the first checks on the proposed methodology.
 To measure the parallel efficiency load the configuration file <u>cfos*pi/api_stats.cfo</u> This configuration pops up a table with %time that every thread spends in every MPI call. Look at the global statistics at the bottom of the outside mpi column. Entry Average represents the application parallel efficiency, entry Avg/Max represents the global load balance and entry Maximum represents the communication efficiency. If any of those values are lower than 85% is recommended to look at the corresponding metric in detail. Open the control window to identify the phases and iterations of the code.
 To measure the computation time distribution load the configuration file <u>cfgs/general/2dh_usefulduration</u>.cfg This onfiguration pops up a histogram of the duration for the computation regions. The computation regions by the next call. If the histogram does not show vertical lines, it indicates the computation time may be not balanced. Open the control window to look at the time distribution and visually correlate both views.
 To measure the computational load (instructions) distribution load the configuration file <u>cfgs/pepi/2dh_useful_instructions.cfg</u> This configuration pops up a histogram of the instructions for the computation regions. The computation regions are delimited by the exit from an MPI call and the entry to the next call. If the histogram doesn't show vertical lines, it indicates the distribution of the instructions may be not balanced. Open the control window to look at the time distribution and correlate both views.
To measure the serial regions performance look at the IPC timeline loaded



Computation time and work distribution

■ ... and 2dh_useful_instructions.cfg (3rd link) → Shows amount of work

lutorials	
The first question to answer when analyzing a parallel code is "how efficient does it run?". The efficiency of a parallel program can be defined based on two aspects: the parallelization efficiency and the efficiency obtained in the execution of the serial regions. These two metrics would be the first checks on the proposed methodology.	
 To measure the parallel efficiency load the configuration file <u>cfgs/mpi/spi_stats.cfg</u> This configuration pops up a table with %time that every thread spends in every MPI call. Look at the global statistics at the bottom of the outside mpi column. Entry Average represents the application parallel efficiency, entry Avg/Max represents the global load balance and entry Maximum represents the communication efficiency. If any of those values are lower than 85% is recommended to look at the corresponding metric in detail. Open the control window to identify the phases and iterations of the code. 	
 To measure the computation time distribution load the configuration file cfgs/general/2dh, usefulduration.cfg This configuration pops up a histogram of the duration for the computation regions. The computation regions are delimited by the exit from an MPI call and the entry to the next call. If the histogram does not show vertical lines, it indicates the computation time may be not balanced. Open the control window to look at the time distribution and visually correlate both views. 	
 To meas re the computational load (instructions) distribution load the configura on file <u>cfgs/papi/2dh_useful_instructions</u>.cfg This configura ion pops up a histogram of the instructions for the computation regions. The computation regions are demined by the externon an merical and the entry to the next call. If the histogram doesn't show vertical lines, it indicates the distribution of the instructions may be not balanced. Open the control window to look at the time distribution and correlate both views. 	
To measure the serial regions performance look at the IPC timeline loaded	•



Where does this happen?

Go from tables to timelines



• |

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Where does this happen

Slow & **Fast** at the same time **→ Imbalance**



Where does this happen

• Hints \rightarrow Callers \rightarrow Caller function



Save CFG's (2 methods)



Save CFG's (2 methods)



CFG's distribution

Paraver comes with many more included CFG's

Paraver				
<u>File</u> Hints <u>H</u> elp				
Load <u>T</u> race				
Previous Traces	Load Configur	ration		
Unioad Traces	Loud Contrigui	ation		
Load <u>Configuration</u>	Look in:	cfas		
Save Configuration	LOOK IN.	cigs	<u>•</u> •	
Load Session CTRL+I	Discust made		🖻 as malia a	. Falalin a
Save Session CTRL+S	burst_mode	Java	sampling	+rolaing
Preferences	clustering	🚞 mpi	📄 scripts	
Ouit	Counters PA	PI CompSs	🚞 software	counters
2dh useful instructions				
Suseful instructions 2DZoom range [1.60154e+00]	CUDA	OpenCL	spectral	
	folding	DpenMP 📄		
	General	Dethread		
	General	paneda		
Files & Window Properties				
	·			
⊕ share	File name:		-	Open
	-	1		open
Telegram Telegram Telegram	Files of type:	Paraver configuration	file (* cfa)	Cancal
E. Wine	Thes of type.	Paraver configuration	me (*.crg)	Cancel
🖻 🗁 paraver				
Paraver files				

Hints: a good place to start!

Paraver suggests CFG's based on the information present in the trace





Cluster-based analysis



Install Clustering in your laptop

- Download a binary for your OS
 - https://tools.bsc.es/downloads

laptop> tar xf clusteringsuite-2.6.8-Linux_x86_64.tar.bz2

laptop> mv clusteringsuite-2.6.8-Linux_x86_64 clustering

Also available in JURECA

jrlXX> module load ClusteringSuite/2.6.8

Use clustering analysis

Run clustering

laptop> cd \$HOME/tools-material/clustering

laptop> \$HOME/clustering/bin/BurstClustering \

- -d cluster.xml \
- -i ../extrae/lulesh2.0_27p.prv \
- -o lulesh2.0_27p_clustered.prv

If you didn't get your own trace, use a prepared one from:

jrlXX> ls \$HOME/tools-material/traces/lulesh2.0_27p.prv

Cluster-based analysis

Check the resulting scatter plot

laptop> gnuplot lulesh2.0_27p_clustered.IPC.PAPI_TOT_INS.gnuplot

- Identify main computing trends
- Work (Y) vs. Speed (X)
- Look at the clusters shape
 - Variability in both axes indicate potential imbalances



Correlating scatter plot and time distribution

Open the clustered trace with Paraver and look at it

laptop> \$HOME/paraver/bin/wxparaver <path-to>/lulesh2.0_27p_clustered.prv

- Display the distribution of clusters over time
 - File → Load configuration → \$HOME/paraver/cfgs/clustering/clusterID_window.cfg





BSC Tools Hands-On

Lau Mercadal, Germán Llort (tools@bsc.es) Barcelona Supercomputing Center

