

BSC Tools Hands-On

Germán Llort, Judit Giménez Barcelona Supercomputing Center





Getting a trace with Extrae



Extrae features

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

No need to recompile / relink!

Extrae overheads

| | Average values | INTI (Haswell) |
|------------------------------|----------------|----------------|
| Event | 150-200 ns | 140 ns |
| Event + PAPI | 750 ns – 1 us | 816 ns |
| Event + callstack (1 level) | 600 ns | 719 ns |
| Event + callstack (6 levels) | 1.9 us | 1.8 us |

How does Extrae work?

- Symbol substitution through LD_PRELOAD <</p>
 - Specific libraries for each combination of runtimes
 - MPI
 - OpenMP
 - OpenMP+MPI
 - ...
- Dynamic instrumentation
 - Based on DynInst (developed by U.Wisconsin/U.Maryland)
 - Instrumentation in memory
 - Binary rewriting

```
Static link (i.e., PMPI, Extrae API)
```



How to use Extrae?

- 1. Adapt the job submission script
- 2. [Optional] Tune the Extrae XML configuration file
 - Examples distributed with Extrae at \$EXTRAE_HOME/share/example
- 3. Run with instrumentation

- For further reference check the **Extrae User Guide:**
 - Also distributed with Extrae at \$EXTRAE_HOME/share/doc

http://www.bsc.es/computer-sciences/performance-tools/documentation

Login to INTI and copy the examples

- > ssh -X <USER>@inti.ocre.cea.fr
- > module load datadir/formation
- > cp -r \$FORMATION_HOME/BSC_Tools/tools-material \$SCRATCHDIR
- > ls \$SCRATCHDIR/tools-material
 - ... apps/
 - ... clustering/
 - ... extrae/
 - ... traces/

Adapt the job script to load Extrae with LD_PRELOAD

> vi \$SCRATCHDIR/tools-material/extrae/job.mpi

```
#!/bin/bash
#MSUB -r job.mpi
#MSUB -n 27
#MSUB -T 600
#MSUB -o job %I.o
#MSUB -e job_%I.e
#MSUB -q haswell
set -x
cd ${BRIDGE MSUB PWD}
ccc mprun -n 27
../apps/lulesh/lulesh2.0 -i 10 -p -s 65
```

Adapt the job script to load Extrae with LD_PRELOAD

> vi \$SCRATCHDIR/tools-material/extrae/trace.sh #!/bin/bash #!/bin/bash #MSUB -r job.mpi export EXTRAE HOME=<installation-path> #MSUB -n 27 #MSUB -T 600 export EXTRAE CONFIG FILE=./extrae.xml #MSUB -o job %I.o #MSUB -e job %I.e #export LD PRELOAD=\$EXTRAE HOME/lib/libmpitracef.so #Fortran #MSUB -q haswell export LD_PRELOAD=\$EXTRAE HOME/lib(libmpitrace.so #C set -x \$@ cd \${BRIDGE MSUB PWD} ccc mprun -n 27 (./trace.sh ../apps/lulesh/lulesh2.0 1 10 -p -s 65 Select tracing library

LD_PRELOAD library selection

Choose depending on the 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 include suffix "f" in Fortran codes

Run with instrumentation

Submit your job

@ inti

- > cd \$SCRATCHDIR/tools-material/extrae
- > ccc msub job.mpi -E "--reservation=vihps"
- Once finished the trace will be in the same folder: lulesh2.0.{pcf,prv,row} (3 files)
- Any issue?
 - Already generated at \$SCRATCHDIR/tools-material/traces

Extrae XML configuration: extrae.xml

> vi \$SCRATCHDIR/tools-material/extrae/extrae.xml



Extrae XML configuration: extrae_config.xml (II)

```
<counters enabled="yes">
 <cpu enabled="yes" starting-set-distribution="cyclic">
     <set enabled="yes" changeat-time="500000us" domain="all">
                                                                     Define which HW
       PAPI TOT INS, PAPI TOT CYC, PAPI L2 DCM, PAPI L3 TCM
                                                                       counters are
     </set>
    <set enabled="yes" changeat-time="500000us" domain="all">
                                                                        measured
      PAPI TOT INS, PAPI TOT CYC, RESOURCE STALLS
    </set>
    <set ... /set>
 </cpu>
 <network enabled="no" />
 <resource-usage enabled="no" />
  <memory-usage enabled="no" />
</counters>
```

Extrae XML configuration: extrae_config.xml (III)

```
<buffer enabled="yes">
                                          Trace buffer size
  <size enabled="yes">5000000</size>
  <circular enabled="no" />
</buffer>
<sampling enabled="no" type="default" period="50m" variability="10m" />
<merge enabled="yes"
       synchronization="default"
                                        Merge intermediate
       tree-fan-out="16"
       max-memory="512"
                                          files into Paraver
       joint-states="yes"
                                                trace
       keep-mpits="yes"
       sort-addresses="yes"
       overwrite="yes"
/>
```



Installing Paraver & First analysis steps



Extrae

Research »

Documentation »

Publications

Download

Install Paraver in your laptop

- Download from <u>http://tools.bsc.es</u>
- Also available @inti
 - \$FORMATION_HOME/BSC_Tools/pkgs



BSC

Home

news@tools:~ > Paraver

Paraver »

Install Paraver in your laptop (II)

Uncompress into your home directory

@ your laptop

```
> tar xvfz wxparaver-4.6.4.rc1-linux-x86_64.tar.gz -C $HOME
```

> cd \$HOME

> ln -s wxparaver-4.6.4.rcl-linux-x86 64 paraver

- Download Paraver tutorials and uncompress into the Paraver directory
 - From website: <u>https://tools.bsc.es/sites/default/files/documentation/paraver-tutorials-20150526.tar.gz</u>
 - From @inti: \$FORMATION_HOME/BSC_Tools/pkgs/paraver-tutorials-20150526.tar.gz

@ your laptop

> tar xvfz paraver-tutorials-20150526.tar.gz -C \$HOME/paraver

Launch Paraver



> \$HOME/paraver/bin/wxparaver &

Check that tutorials are available





Launch Paraver (II)

Any issues installing in your laptop? Also available @ INTI

@ inti

> cd \$FORMATION_HOME/BSC_Tools/tools/wxparaver-4.6.4/bin

> ./wxparaver \$SCRATCHDIR/tools-material/extrae/lulesh2.0.prv &

First steps of analysis

Copy the trace to your desktop (All 3 files: *.prv, *.pcf, *.row)



Tutorials

Barcelona

Follow Tutorial #3

Load the trace

Introduction to Paraver and Dimemas methodology



Measure the parallel efficiency



Measure the computation time distribution





Cluster-based analysis



Use clustering analysis

Run clustering

@ inti

> cd \$SCRATCHDIR/tools-material/clustering
> ./clusterize.sh ../extrae/lulesh2.0.prv

Look at the results

@ inti

> gnuplot lulesh2.0_clustered.IPC.PAPI_TOT_INS.gnuplot

Cluster-based analysis

Check the resulting scatter plot

> gnuplot lulesh2.0_clustered.IPC.PAPI_TOT_INS.gnuplot



Identify main computing trends with respect to work & performance

Correlating scatter plot and time distribution

• Copy the clustered trace to your laptop and look at it

@ your laptop

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

■ File → Load configuration → \$HOME/paraver/cfgs/clustering/clusterID_window.cfg

