# BSCTools
# Hands-On

Judit Giménez, Germán Llort

Barcelona Supercomputing Center

# Getting a trace with Extrae

# Extrae features

- Parallel programming models
  - MPI, OpenMP, pthreads, OmpSs, CUDA, OpenCL, Intel MIC…

- 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!**

# How does Extrae work?

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

**Recommended**

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

- Alternatives
  - 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 it!

- 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

# Log in to Hornet and copy the example files to your home directory

```
> ssh <user>@hornet.hww.de

> cp -r ~xhpgl/paraver-tutorial $HOME

> cd $HOME/paraver-tutorial
```

# Example: Adapt the job script to load Extrae with LD_PRELOAD

## run-27procs.pbs

```
#!/bin/bash
#PBS -l nodes=2:ppn=24
#PBS -l walltime=00:10:00
#PBS -N LULESH
#PBS -e $PBS_JOBID.err
#PBS -o $PBS_JOBID.out
#PBS –q R_vihps04

cd $PBS_O_WORKDIR
module switch PrgEnv-cray PrgEnv-gnu

time aprun -n 27 ./lulesh
```

# Example: Adapt the job script to load Extrae with LD_PRELOAD

### run-and-trace-27procs-24+3.pbs

```
#!/bin/bash
#PBS -l nodes=2:ppn=24
#PBS -l walltime=00:10:00
#PBS -N LULESH
#PBS -e $PBS_JOBID.err
#PBS -o $PBS_JOBID.out
#PBS -q R_vihps04

cd $PBS_O_WORKDIR
module switch PrgEnv-cray PrgEnv-gnu

export TRACE_NAME=lulesh.27procs.24+3.prv

time aprun -n 27 ./trace.sh ./lulesh
```

# Example: Adapt the job script to load Extrae with LD_PRELOAD

**run-and-trace-27procs-24+3.pbs**

**trace.sh**

```
#!/bin/bash
#PBS -l nodes=2:ppn=24
#PBS -l walltime=00:10:00
#PBS -N LULESH
#PBS -e $PBS_JOBID.err
#PBS -o $PBS_JOBID.out
#PBS -q R_vihps04

cd $PBS_O_WORKDIR
module switch PrgEnv-cray PrgEnv-gnu

export TRACE_NAME=lulesh.27procs.24+3.prv

time aprun -n 27 ./trace.sh ./lulesh
```

```
#!/bin/bash

rm -rf TRACE*.sym TRACE*.mpits TRACE*.spawn set-*

# Configure Extrae
export EXTRAE_HOME=~xhpgl/tools/extrae/default
source $EXTRAE_HOME/etc/extrae.sh
export EXTRAE_CONFIG_FILE=./extrae-config.xml

# Load the tracing library (C/Fortran)
export LD_PRELOAD=${EXTRAE_HOME}/lib/libmpitrace.so
#export LD_PRELOAD=${EXTRAE_HOME}/lib/libmpitracef.so

# Run the program
$*
```

# LD_PRELOAD library selection

- Choose depending on the application type

| Library | Serial | MPI | OpenMP | pthread | CUDA |
|---|---|---|---|---|---|
| libseqtrace | ✓ | | | | |
| libmpitrace[f][1] | | ✓ | | | |
| libomptrace | | | ✓ | | |
| libpttrace | | | | ✓ | |
| libcudatrace | | | | | ✓ |
| libompitrace[f] [1] | | ✓ | ✓ | | |
| libptmpitrace[f] [1] | | ✓ | | ✓ | |
| libcudampitrace[f] [1] | | ✓ | | | ✓ |

**[1] include suffix "f" in Fortran codes**

# Extrae XML configuration: extrae_config.xml

```xml
<mpi enabled="yes">
  <counters enabled="yes" />
</mpi>

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

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

<callers enabled="yes">
  <mpi enabled="yes">1-3</mpi>
  <sampling enabled="no">1-5</sampling>
</callers>
```
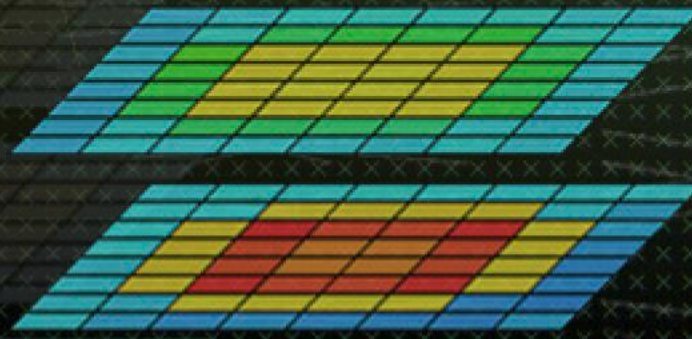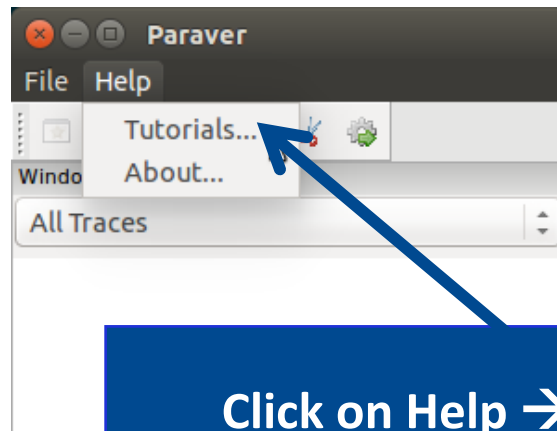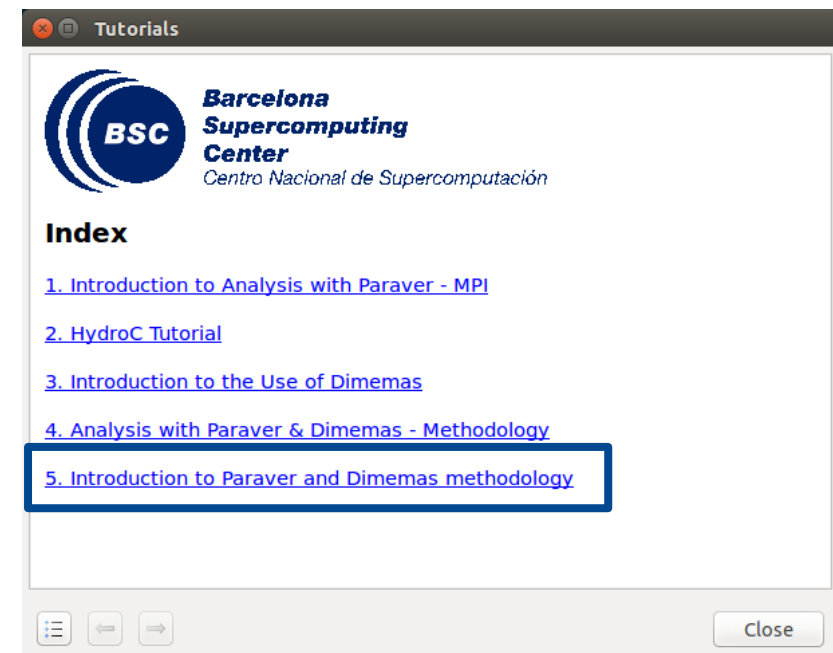
**Trace MPI calls + HW counters**

**Trace call-stack events @ MPI calls**

# Extrae XML configuration: extrae_config.xml (II)

```xml
<counters enabled="yes">
  <cpu enabled="yes" starting-set-distribution="cyclic">
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_L1_DCM,PAPI_L2_DCM
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_L3_TCM,RESOURCE_STALLS:SB,RESOURCE_STALLS:ROB
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_SR_INS,PAPI_BR_CN,PAPI_BR_UCN
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_BR_MSP,PAPI_LD_INS
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS,PAPI_TOT_CYC,RESOURCE_STALLS,RESOURCE_STALLS:RS
    </set>
  </cpu>

  <network enabled="no" />

  <resource-usage enabled="no" />

  <memory-usage enabled="no" />
</counters>
```

**Define which HW counters are measured**

# Extrae XML configuration: extrae_config.xml (III)

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


<sampling enabled="no" type="default" period="50m" variability="10m" />


<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 buffer size**

**Enable sampling**

**Merge intermediate files into Paraver trace**

# Run it!

- Submit your job

```
> cd $HOME/paraver-tutorial
> qsub run-and-trace-27procs-24+3.pbs
```

- Copy the resulting trace to your laptop

```
> gzip lulesh.27procs.24+3.prv
> scp <user>@hornet.hww.de:paraver-tutorial/lulesh.27procs.24+3.\* $HOME
```

- Or click to download from: ( 3 files, lulesh.27procs.24+3.[ prv.gz | row | pcf ] )
https://www.dropbox.com/sh/42teru7vunc4yot/AACy1ckEr2sUKgp-3re3IK8Aa?dl=0

# Installing Paraver

# Install Paraver in your laptop

- Download the Paraver binaries

| Linux 64 bits | Click to download Paraver for Linux 64 bits from: https://www.dropbox.com/s/hu6u2pfoa8zx0ha/wxparaver-4.5.6-linux-x86_64.tar.gz?dl=0 |
|---|---|

| Linux 32 bits | Click to download Paraver for Linux 32 bits from: https://www.dropbox.com/s/jtkik0mrazu0ue2/wxparaver-4.5.6-linux-x86_32.tar.gz?dl=0 |
|---|---|

| Mac | Click to download Paraver for Mac from: https://www.dropbox.com/s/7vq6ereyo5incb6/wxparaver_mac.zip?dl=0 |
|---|---|

| Windows | Click to download Paraver for Windows from: https://www.dropbox.com/s/pwcruoezjuyoqfc/wxparaver-4.5.6-win.zip?dl=0 |
|---|---|

# Install Paraver in your laptop (II)

- Uncompress into your home directory

```
> tar xvfz wxparaver-4.5.6-linux-x86_64.tar.gz -C $HOME

> cd $HOME

> ln -s wxparaver-4.5.6-linux-x86_64 paraver
```

- Download Paraver tutorials and uncompress into the Paraver directory

```
              Click to download Paraver tutorials from:
    https://www.dropbox.com/s/ece15xzqxa4j2qo/paraver-tutorials-20141016.tar.gz?dl=0

> tar xvfz paraver-tutorials-20141016.tar.gz -C $HOME/paraver --strip=1
```

# First steps of analysis

- Load the trace with Paraver

```
> $HOME/paraver/bin/wxparaver lulesh.27procs.24+3.prv
```

- Open Tutorial #5
  - Help → Tutorials



**Click on Help → Tutorials**

# Measure the parallel efficiency

▪ Click on the "mpi_stats.cfg"

# Measure the computation time distribution

- Click on the "2dh_usefulduration.cfg"

# Compare with other configurations

- Nodes more balanced (14 processes in node 1, 13 in node 2)
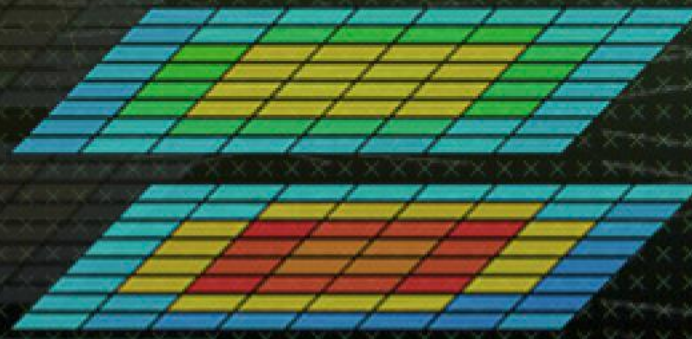  - If you want to get the trace…

```
> cd $HOME/paraver-tutorial
> qsub run-and-trace-27procs-14+13.pbs
```

- CPU sockets more balanced (max 7 processes per socket)
  - If you want to get the trace…

```
> cd $HOME/paraver-tutorial
> qsub run-and-trace-27procs-14+13-7PExSOCKET.pbs
```

- <u>Or click to download from:</u>
  https://www.dropbox.com/sh/42teru7vunc4yot/AACy1ckEr2sUKgp-3re3IK8Aa?dl=0

# Thank you!