# Extrae and Paraver Hands-on

# Install Paraver

- Download from https://tools.bsc.es/downloads



Pick your version

wxparaver-4.7.2-win.zip

wxparaver-4.7.2-mac.zip

wxparaver-4.7.2-Linux_i686.tar.gz (32-bits)
wxparaver-4.7.2-Linux_x86_64.tar.gz (64-bits)

~gimenez1/BSC-paraver

# Install Paraver (II)

- Download tutorials:
  - Documentation
    - Paraver tutorials



**Download links**

`~gimenez1/BSC-paraver`

# Uncompress, rename & move

- Paraver

```
> tar xf wxparaver-4.8.1-linux-x86_64.tar.gz

> mv wxparaver-4.8.1-linux-x86_64 paraver
```

- Tutorials

```
> tar xf paraver-tutorials-20150526.tar.gz

> mv paraver-tutorials-20150526 paraver/tutorials
```

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

# Check that everything works

- Start Paraver

```
> paraver/bin/wxparaver
```

- Check that tutorials are available

Click on Help → Tutorials

- Trouble installing locally? Remote open from kabre

```
> ssh –Y <USER>@cluster.cenat.ac.cr
> ln –s ~gimenez1/tools/wxparaver64/bin/wxparaver
  wxparaver
```

# Log in to kabré

```
> ssh -Y <USER>@cluster.cenat.ac.cr
```

- Copy the examples to your home folder:

@ cluster.cenat.ac.cr

```
> cp -r ~jgimenez/BSC-handson/ $HOME

> ls -l $HOME/BSC-handson/
    … apps
    … traces
    … slides
```

# OpenMP example: matrix

- Comparing the location of the parallel for pragma

```
#pragma omp parallel for
for(int i=0; i<N; i++){
        for(int k=0; k<N; k++){
                for(int j=0; j<N; j++)
```

```
for(int i=0; i<N; i++){
#pragma omp parallel for
        for(int k=0; k<N; k++){
                for(int j=0; j<N; j++){
```

matrix.l1                                            matrix.l2

- Check the script content
- Submit the job script

@ cluster.cenat.ac.cr

```
> cd $HOME/BSC-handson/apps/matrix
> qsub matrix.pbs
```

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

# OpenMP example: matrix

- Copy the traces to your laptop

```
> scp <USER>@cluster.cenat.ac.cr: \
    BSC_handson/apps/matrix/matrix.l?.* $HOME
```

- Load the traces with Paraver

```
> wxparaver matrix.l?.prv*
```

- Trouble getting in the queues? Already available at ~/BSC-handson/traces/matrix/matrix.l?.*

- Compare the two executions with paraver

# MPI example: jacobi

- Submit the job script

```
> cd $HOME/BSC-handson/apps/jacobi
> qsub job.kabre
```

- Copy the traces to your laptop and load the trace with Paraver

@ your laptop

```
> scp <USER>@cluster.cenat.ac.cr: \
    BSC_handson/apps/jacobi/jacobi.???* $HOME
> wxparaver jacobi.prv*
```

- Trouble getting in the queues? Already available at ~/BSC-handson/traces/jacobi/jacobi.???*

- Look at the execution with paraver

# Extrae features

- Platforms
  - Intel, Cray, BlueGene, MIC, ARM, Android, Fujitsu Sparc...

- Parallel programming models
  - 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 (Dyninst)

- Periodic sampling

- User events (Extrae API)

**No need to recompile / relink!**

# Extrae overheads

| Average values | Kábre (login) | Kabré (Zárate) |
|---|---|---|
| Event | 150-200 ns | 180 ns | 600 ns |
| Event + PAPI | 750 ns – 1 us | 580 ns | 4.7 us |
| Event + callstack (1 level) | 600 ns | 750 ns | 3.4 us |
| Event + callstack (6 levels) | 2 us | 1.7us | 8.2 us |

# 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
  - Compiler instrumentation (-finstrument-functions – Intel, GNU)
  - Static link (i.e., PMPI, Extrae API)

# Using Extrae in 3 steps

1. **Adapt** your job submission scripts

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:**

  - https://tools.bsc.es/sites/default/files/documentation/html/extrae/index.html

  - Also distributed with Extrae at $EXTRAE_HOME/share/doc

# Step 1: Adapt the job script to load Extrae

```
> vi $HOME/BSC-handson/apps/lulesh/job.kabre
```

**job.kabre**

```
#PBS -N extrae
#PBS -q phi-n1h72
#PBS -l nodes=1:ppn=27
#PBS -l walltime=00:20:00

cd $PBS_O_WORKDIR

module load mpich/3.2.1




mpirun  ./lulesh2.0 -i 20 -s 64 -p
```

# Step 1: Adapt the job script to load Extrae

```
> vi $HOME/BSC-handson/apps/lulesh/job.kabre
```

**job.kabre**

```
#PBS -N extrae
#PBS -q phi-n1h72
#PBS -l nodes=1:ppn=27
#PBS -l walltime=00:20:00

cd $PBS_O_WORKDIR

module load mpich/3.2.1


export TRACE_NAME=lulesh2_27p.prv


mpirun ./extrae/trace.sh  ./lulesh2.0 -i 20 -s
64 -p
```

# Step 1: Adapt the job script to load Extrae

@ cluster.cenat.ac.cr

```
> vi $HOME/BSC-handson/apps/lulesh/extrae/trace.sh
```

**job.kabre**

**Select "what to trace"**

**trace.sh**

```
#PBS -N extrae
#PBS -q phi-n1h72
#PBS -l nodes=1:ppn=27
#PBS -l walltime=00:20:00

cd $PBS_O_WORKDIR

module load mpich/3.2.1


export TRACE_NAME=lulesh2_27p.prv



mpirun  /extrae/trace.sh   /lulesh2.0 -i 20 -s
64 -p
```

```
#!/bin/bash

export EXTRAE_HOME=/home/jgimenez/tools/extrae-3.6.1
export EXTRAE_CONFIG_FILE=extrae/extrae.xml
export LD_PRELOAD=${EXTRAE_HOME}/lib/libmpitrace.so #
For C apps
#export LD_PRELOAD=${EXTRAE_HOME}/lib/libmpitracef.so
# For Fortran apps

## Run the desired program
$*
```

**Select your type of application**

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

17

# Step 1: Which tracing library?

- 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**

# Step 3: Run it!

- Submit your job

@ cluster.cenat.ac.cr

```
> cd $HOME/BSC-handson/apps/lulesh

> qsub job.kabre
```

- Easy! ☺

# Step 2: Extrae XML configuration

@ cluster.cenat.ac.cr

```
> vi $HOME/BSC-handson/apps/lulesh/extrae/extrae.xml
```

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

<openmp enabled="no">
  <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 the MPI calls
(What's the program doing?)**

**Trace the call-stack
(Where in my code?)**

Barcelona
**Supercomputing
Center**
Centro Nacional de Supercomputación

# Step 2: Extrae XML configuration (II)

@ cluster.cenat.ac.cr

```
> vi $HOME/BSC-handson/apps/lulesh/extrae/extrae.xml
```

```xml
<counters enabled="yes">
  <cpu enabled="yes" starting-set-distribution="1">
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_L1_DCM, PAPI_L2_DCM, PAPI_BR_MSP, RESOURCE_STALLS
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_L3_TCM, PAPI_LD_INS, PAPI_SR_INS
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_VEC_DP
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_VEC_SP, PAPI_FP_INS
    </set>
  </cpu>
  <network enabled="no" />
  <resource-usage enabled="no" />
  <memory-usage enabled="no" />
</counters>
```

**Select which HW counters are measured**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Step 2: Extrae XML configuration (III)

@ cluster.cenat.ac.cr

```
> vi $HOME/BSC-handson/apps/lulesh/extrae/extrae.xml
```

```xml
<buffer enabled="yes">
  <size enabled="yes">5000000</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 (Flush/memory trade-off)**

**Enable sampling (Want more details?)**

**Automatic post-processing to generate the trace**

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

22

# All done! Check your resulting trace

- Once finished (check with `squeue`) you will have the trace (3 files):

```
➤ ls -l $HOME/BSC-handson/apps/lulesh
  ...
  lulesh2_27p.pcf
  lulesh2_27p.prv
  lulesh2_27p.row
```

- Trouble getting in the queues? Already available at ~/BSC-handson/traces/ lulesh/lulesh2_27p.*

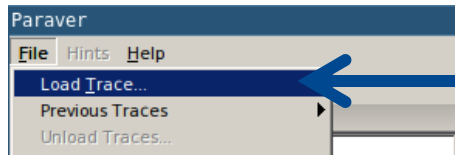- Now let's look into it !

# First steps of analysis
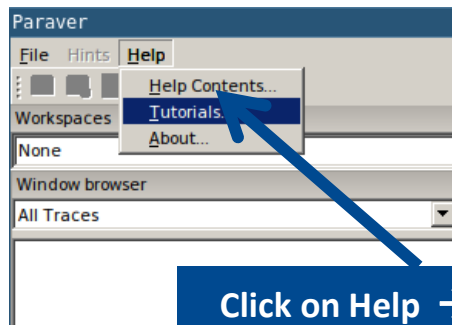
- Copy the trace to your laptop

@ your laptop

```
> scp <USER>@cluster.cenat.ac.cr: \
     BSC_handson/apps/lulesh/lulesh2_27p.* $HOME
```
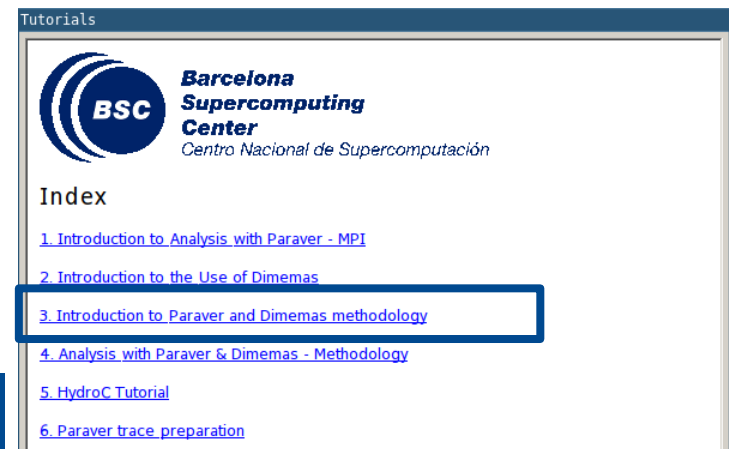
- Load the trace with Paraver

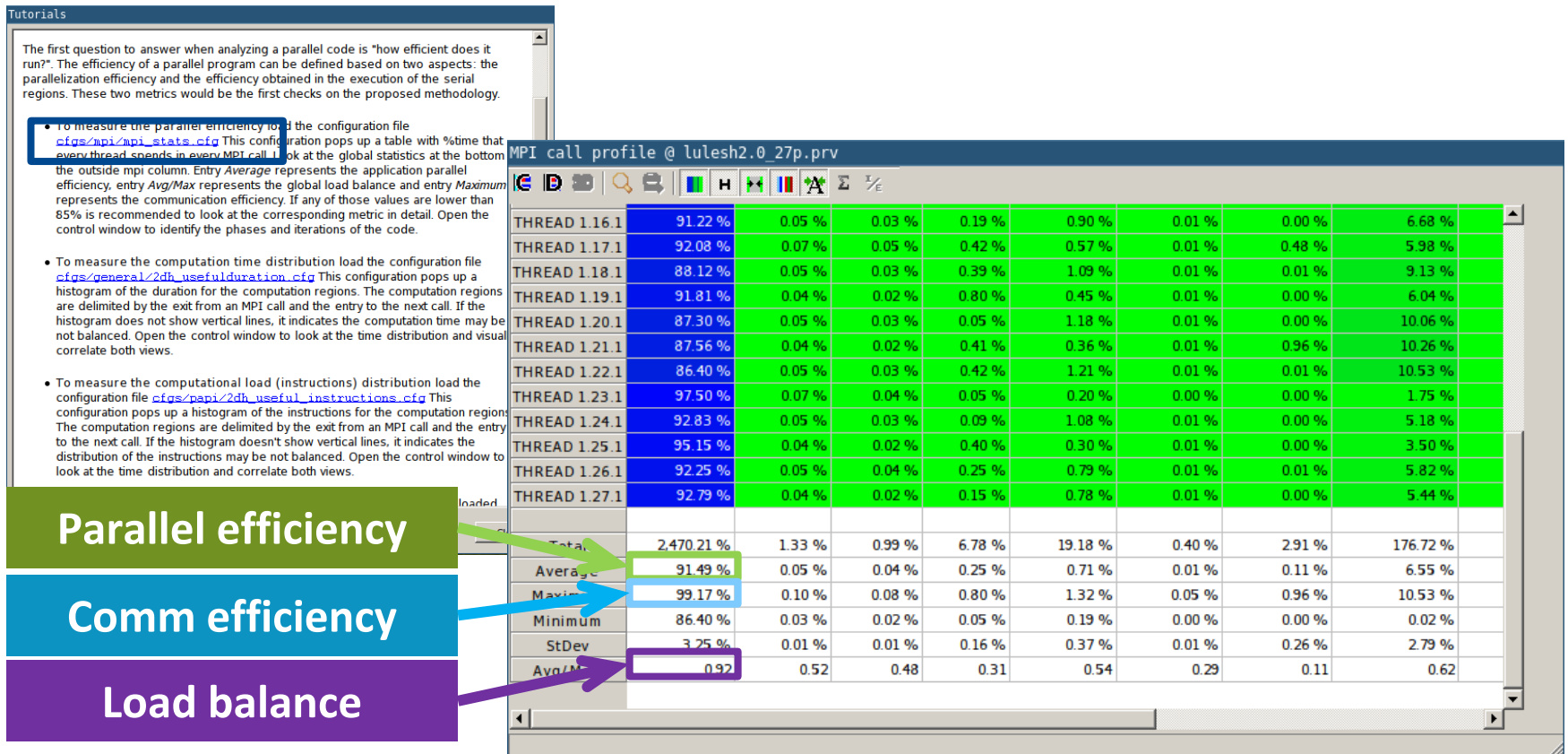**Click on File → Load Trace → Browse to "lulesh2.0_27p.prv"**

- Follow Tutorial #3

**Click on Help → Tutorials**

Tutorials

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

Index

1. Introduction to Analysis with Paraver - MPI

2. Introduction to the Use of Dimemas

3. Introduction to Paraver and Dimemas methodology

4. Analysis with Paraver & Dimemas - Methodology
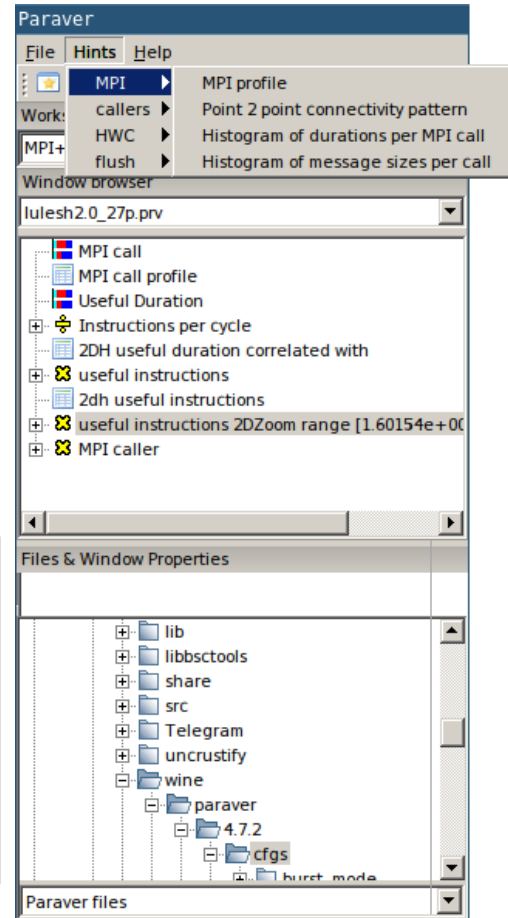
5. HydroC Tutorial

6. Paraver trace preparation

# Measure the parallel efficiency

- Click on "mpi_stats.cfg"
  - Check the **Average** for the column labeled "**Outside MPI**"

# Hints: a good place to start!

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

# CFG's distribution

- Paraver comes with many more included CFG's