**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

EXCELENCIA SEVERO OCHOA

# Extrae Hands-On

✉ [tools@bsc.es](mailto:tools@bsc.es)

DEEP-EST Early Access Programme

# 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 or 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)

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# Extrae on DEEP-EST (I)

- Log-in to DEEP-EST:

```
laptop$ ssh -Y <USER>@deep-fz.juelich.de
```

- Extrae is available via modules for 4 toolchains... choose yours!

  1. GCC compiler with ParaStation MPI

```
deep$ module use $OTHERSTAGES
deep$ ml Stages/Devel-2019a
deep$ ml GCC/8.3.0
deep$ ml ParaStationMPI/5.4.6-1
deep$ ml Extrae
```

  2. Intel compiler with Intel MPI

```
deep$ module use $OTHERSTAGES
deep$ ml Stages/Devel-2019a
deep$ ml Intel/2019.3.199-GCC-8.3.0
deep$ ml IntelMPI/2018.5.288
deep$ ml Extrae
```

# Extrae on DEEP-EST (II)

3. Intel compiler with ParaStationMPI

```
deep$ module use $OTHERSTAGES
deep$ ml Stages/Devel-2019a
deep$ ml Intel/2019.5.281-GCC-8.3.0
deep$ ml ParaStationMPI/5.4.6-1
deep$ ml Extrae
```

4. Intel compiler with ParaStationMPI (MPI multithreaded support)

```
deep$ module use $OTHERSTAGES
deep$ ml Stages/Devel-2019a
deep$ ml Intel/2019.5.281-GCC-8.3.0
deep$ ml ParaStationMPI/5.4.6-1-mt
deep$ ml Extrae
```

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

# Getting your first trace

- Provided package **lulesh-example.tar.xz** contains:
  - README
  - Application (lulesh2.0)
  - Jobscripts to execute and trace (job.slurm, trace.sh)
  - Configuration of the tracing tool (extrae.xml)
  - Already generated trace (trace/*.{pcf,prv,row})

- Copy this package to DEEP-EST cluster and uncompress into **your own /work/<project>/<user> folder**

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# Using Extrae in 3 steps

1. **Adapt** your job submission scripts

2. **Configure** what to trace

   - XML configuration file
   - Example configurations at `$EXTRAE_HOME/share/example`

3. **Run** it!

- For further reference check the **Extrae User Guide:**

  - https://tools.bsc.es/doc/html/extrae
  - Also distributed with Extrae at `$EXTRAE_HOME/share/doc`

# Step 1: Adapt the job script to load Extrae

- Example of a standard jobscript (without tracing)

```
#!/bin/bash
#SBATCH --job-name=lulesh2.0_27p
#SBATCH --output=%x_%j.out
#SBATCH --error=%x_%j.err
#SBATCH --ntasks=27
#SBATCH --nodes=2
#SBATCH --cpus-per-task=1
#SBATCH --exclusive
#SBATCH --time=00:10:00
#SBATCH --partition=dp-cn
```

**Request resources**

```
srun ./lulesh2.0 -i 10 -s 65
```

**Run the program**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Step 1: Adapt the job script to load Extrae

- Jobscript modified to load Extrae

```
#!/bin/bash
#SBATCH --job-name=lulesh2.0_27p
#SBATCH --output=%x_%j.out
#SBATCH --error=%x_%j.err
#SBATCH --ntasks=27
#SBATCH --nodes=2
#SBATCH --cpus-per-task=1
#SBATCH --exclusive
#SBATCH --time=00:10:00
#SBATCH --partition=dp-cn

module use $OTHERSTAGES
ml Stages/Devel-2019a
ml GCC/8.3.0
ml ParaStationMPI/5.4.6-1
ml Extrae

export TRACE_NAME=lulesh2.0_27p.prv

srun ./trace.sh ./lulesh2.0 -i 10 -s 65
```

**Load Extrae (choose proper toolchain)**

**Optionally specify name of output trace**

**Run with Extrae**

# Step 1: Adapt the job script to load Extrae

- Tracing launcher helper script (trace.sh)

```bash
#!/bin/bash
#SBATCH --job-name=lulesh2.0_27p
#SBATCH --output=%x_%j.out
#SBATCH --error=%x_%j.err
#SBATCH --ntasks=27
#SBATCH --nodes=2
#SBATCH --cpus-per-task=1
#SBATCH --exclusive
#SBATCH --time=00:10:00
#SBATCH --partition=dp-cn

module use $OTHERSTAGES
ml Stages/Devel-2019a
ml GCC/8.3.0
ml ParaStationMPI/5.4.6-1
ml Extrae

export TRACE_NAME=lulesh2.0_27p.prv

srun ./trace.sh ./lulesh2.0 -i 10 -s 65
```

```bash
#!/usr/bin/env bash

export EXTRAE_ENFORCE_FS_SYNC=1

# Configure Extrae
export EXTRAE_CONFIG_FILE=./extrae.xml

# Load the tracing library (choose C/Fortran)
export LD_PRELOAD=$EBROOTEXTRAE/lib/libmpitrace.so
#export LD_PRELOAD=$EBROOTEXTRAE/lib/libmpitracef.so

# Run the program
$*
```

**What to trace?**

**Choose a tracing library depending on the app type (see next slide)**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

10

# 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**

**Barcelona Supercomputing Center**
Centro Nacional de Supercomputación

# Step 2: Extrae XML configuration

```
deep$ vi extrae.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 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)

```
deep$ vi 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
    </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?)**

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Step 2: Extrae XML configuration (III)

```
deep$ vi 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)

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

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

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

# Step 3: Run it!

- Submit your job as usual

```
deep$ sbatch job.slurm
```

- **REMEMBER!** Run job from your /work folder (NOT IN HOME!)

# All done! Check your resulting trace

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

```
deep$ ls -l
...
lulesh2.0_27p.pcf
lulesh2.0_27p.prv
lulesh2.0_27p.row
```

- Any trouble? There's a trace already generated under the "trace" folder

- Now let's look into it !

# First steps of analysis

- Copy the trace to your computer

- Load the trace with Paraver

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

# First steps of analysis

- Follow Tutorial #3
    - Introduction to Paraver and Dimemas methodology

# Measure the parallel efficiency

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



Parallel efficiency (Avg)

Comm efficiency (Max)

Load balance (Avg/Max)

# Focus on the iterative part



Click on Open Control Window

# Focus on the iterative part

# Recalculate efficiency of iterative region

# Recalculate efficiency of iterative region



Right click → Paste → Time

# Efficiency of iterative region

- 3 numbers to quickly describe the efficiency of your code
  - Parallel efficiency → % of time my program is computing (100% is perfect)
  - Comm efficiency → At least 1 process can finish all communications in 100 - Maximum % of the program's time (100% is perfect)
  - Load balance → Ratio of slow/fast processes (1 is perfectly balanced)
  - Any value below 85% (0.85)? Pay attention there…

**Parallel efficiency**

**Comm efficiency**

**Load balance**

MPI call profile @ lulesh2.0_27p.prv

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| THREAD 1.17.1 | 89.00 % | 0.18 % | 0.14 % | 0.14 % | 1.77 % | 8.71 % | 0.06 |
| THREAD 1.18.1 | 82.79 % | 0.13 % | 0.10 % | 0.09 % | 2.83 % | 14.01 % | 0.06 |
| THREAD 1.19.1 | 90.45 % | 0.08 % | 0.07 % | 0.10 % | 0.74 % | 8.49 % | 0.06 |
| THREAD 1.20.1 | 82.55 % | 0.12 % | 0.09 % | 0.24 % | 1.69 % | 15.25 % | 0.06 |
| THREAD 1.21.1 | 82.74 % | 0.09 % | 0.06 % | 0.33 % | 0.42 % | 16.29 % | 0.06 |
| THREAD 1.22.1 | 80.45 % | 0.13 % | 0.09 % | 0.25 % | 3.20 % | 15.82 % | 0.06 |
| THREAD 1.23.1 | 97.57 % | 0.19 % | 0.15 % | 0.12 % | 1.01 % | 0.90 % | 0.06 |
| THREAD 1.24.1 | 90.57 % | 0.14 % | 0.08 % | 0.08 % | 2.01 % | 7.05 % | 0.06 |
| THREAD 1.25.1 | 94.00 % | 0.10 % | 0.06 % | 0.30 % | 0.32 % | 5.16 % | 0.06 |
| THREAD 1.26.1 | 88.89 % | 0.14 % | 0.08 % | 0.08 % | 1.82 % | 8.94 % | 0.06 |
| THREAD 1.27.1 | 90.81 % | 0.10 % | 0.05 % | 0.10 % | 0.63 % | 8.25 % | 0.06 |
| | | | | | | | |
| Total | 2,399.39 % | 3.40 % | 2.87 % | 6.69 % | 32.95 % | 253.10 % | 1.61 |
| Average | 88.87 % | 0.13 % | 0.11 % | 0.25 % | 1.22 % | 9.37 % | 0.06 |
| Maximum | 98.78 % | 0.26 % | 0.23 % | 0.57 % | 3.23 % | 16.62 % | 0.06 |
| Minimum | 80.45 % | 0.06 % | 0.05 % | 0.08 % | 0.25 % | 0.02 % | 0.06 |
| StDev | 4.76 % | 0.04 % | 0.04 % | 0.14 % | 0.83 % | 4.54 % | 0.00 |
| Avg/Max | 0.90 | 0.49 | 0.47 | 0.44 | 0.38 | 0.56 | 0. |

Barcelona Supercomputing Center
Centro Nacional de Supercomputación

# Computation time distribution

- Click on "2dh_usefulduration.cfg" (2nd link) → Shows **time computing**
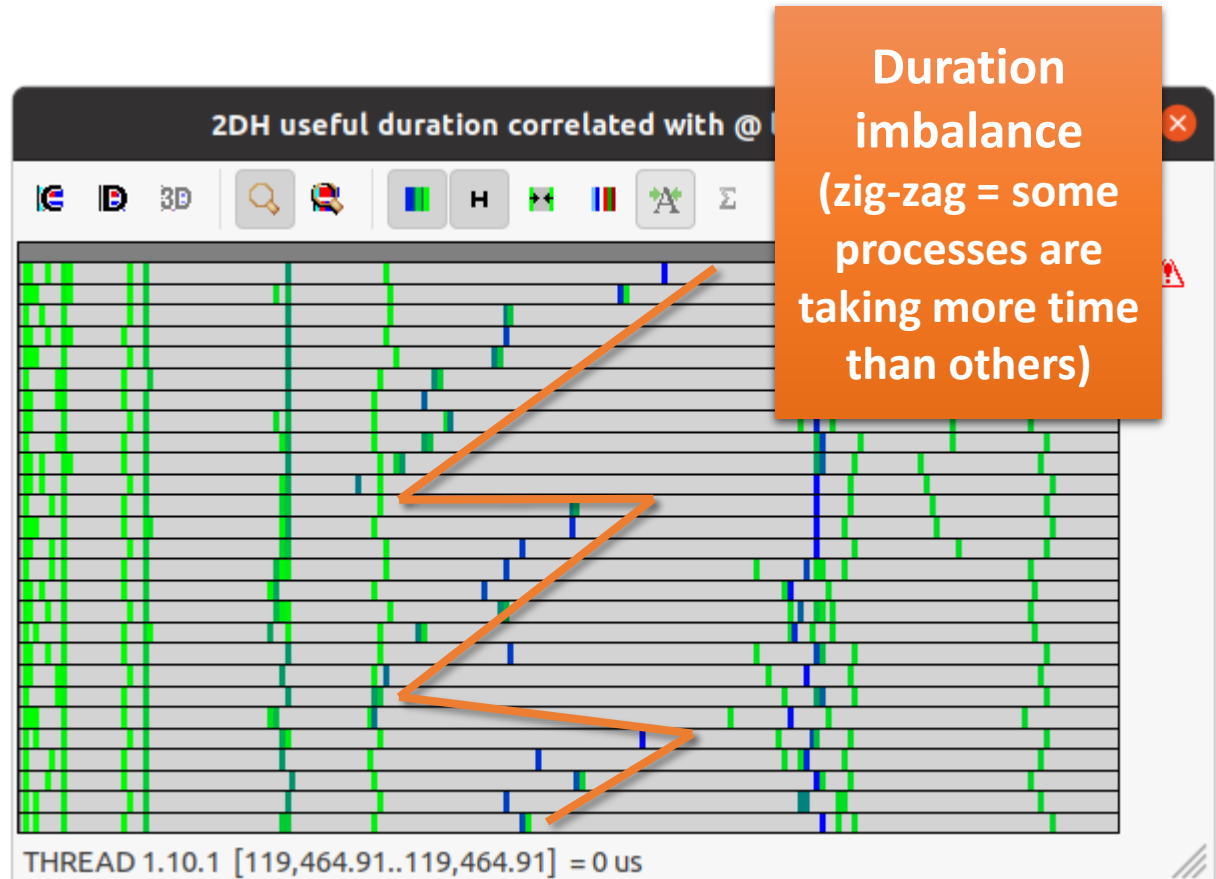
# Focus on the iterative part

- Click on "2dh_usefulduration.cfg" (2nd link) → Shows **time computing**



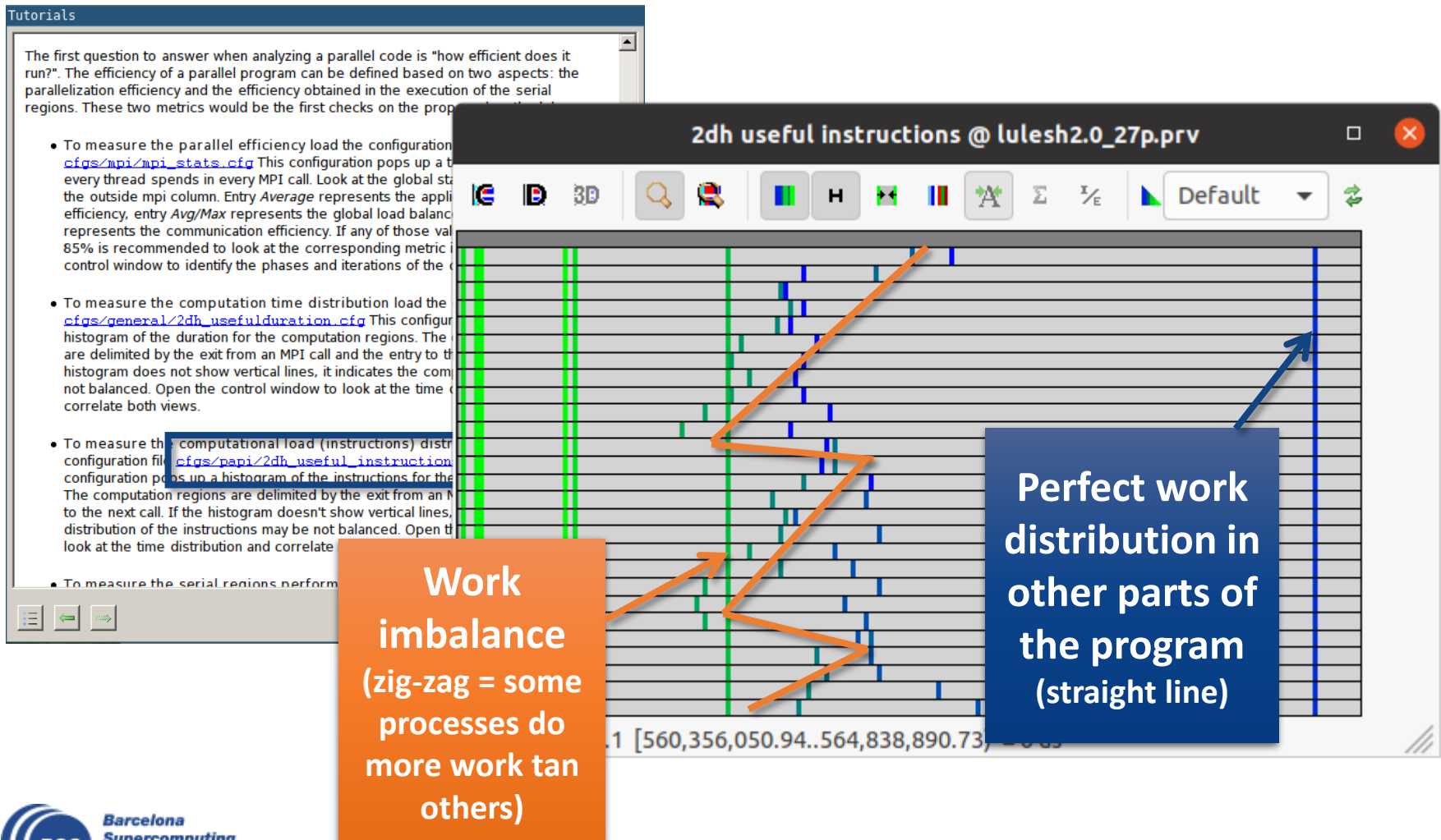**Drag & drop on this area to skip the initializations**

# Computation time distribution

- Click on "2dh_usefulduration.cfg" (2nd link) → Shows **time computing**
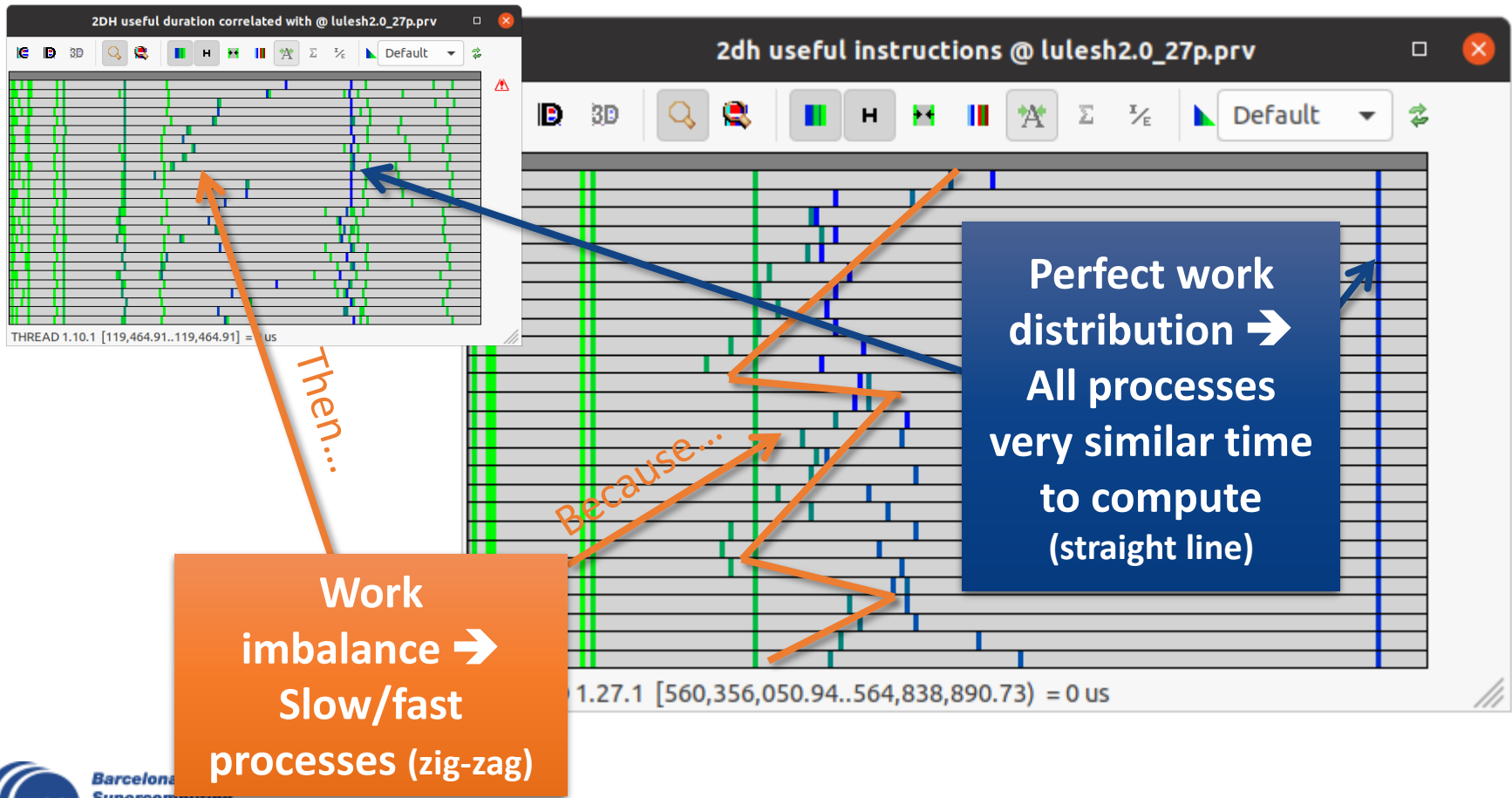


Duration imbalance (zig-zag = some processes are taking more time than others)

# Computation load distribution

- Click on "2dh_useful_instructions.cfg" (3rd link) → Shows **amount of work**



**Work imbalance** (zig-zag = some processes do more work tan others)

**Perfect work distribution in other parts of the program** (straight line)

# Computation load distribution

- Comparing the two histograms → **Similar shapes** → Work distribution determines time computing



**Then...**

**Because...**

**Work imbalance →  Slow/fast processes (zig-zag)**

**Perfect work distribution →  All processes very similar time to compute (straight line)**

# Where does this happen?

- Go from the table to the timeline



**1. Click on "Open Filtered Control Window"**

**2. Select this area (drag-and-drop)**

# Where does this happen?

- Go from the table to the timeline



**Clicking here always rescales. Same as…
Right click →
Fit Semantic Scale →
Fit Both**

# Where does this happen?

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



**Zoom into
1 of the iterations
(by drag-and-dropping)**

# Where does this happen?

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



**1. Right click → Copy**

useful instructions 2DZoom range [1.47399e+08,4.25349e+08) @ lulesh2....

THREAD 1.19.1
THREAD 1.27.1
2,481,808 us          2,739,887 us

- **Hints → Call stack references → Caller function**



**2. Right click → Paste → Time**

MPI caller @ lulesh2.0_27p_binding.prv

THREAD 1.13.1
THREAD 1.17.1
THREAD 1.21.1
THREAD 1.25.1
THREAD 1.27.1
0 us          4,915,946 us

# Where does this happen?

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



- **Hints → Call stack references → Caller function**

# Save CFG's (method 1)



useful instructions 2DZoom range [1.72355e+08,4.97913e+08).c1 (u

THREAD 1.1.1
THREAD 1.10.1
THREAD 1.19.1
THREAD 1.27.1
2,571,619 us                                    2,912,845 us

**Right click on timeline**

| Copy | Ctrl+C |
| Paste | ▸ |
| Clone | |
| Undo Zoom | Ctrl+U |
| Redo Zoom | Ctrl+R |
| Fit Time Scale | |
| Fit Semantic Scale | ▸ |
| Fit Objects | |
| Select Objects... | |
| View | ▸ |
| Paint As | ▸ |
| Drawmode | ▸ |
| Pixel Size | ▸ |
| Object Labels | ▸ |
| Object Axis | ▸ |
| Run | ▸ |
| Synchronize | |
| Remove all sync | |
| **Save** | ▸ |
| Info Panel | |

| Configuration... |
| Image... |
| Image Legend... |
| Text... |

Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

35

# Save CFG's (method 2)

# CFG's distribution

- Paraver comes with many included CFG's ➔ Apply any CFG to any trace!

# Hints: a good place to start!

- Paraver suggests CFG's based on the contents of the trace

# Do it on your code!

- Follow guidelines from slides 7-16 to your own code to get a trace
  - There are more examples of tracing scripts for different programming models under $EBROOTEXTRAE/share/examples

- Follow guidelines from slides 17-34 to conduct an initial analysis
  - The usual suspects:
    - Parallel Efficiency is low? Load balance issues?
    - Imbalances in the durations of computations?
    - Are these caused by work imbalance?