

BSC Tools Hands-On

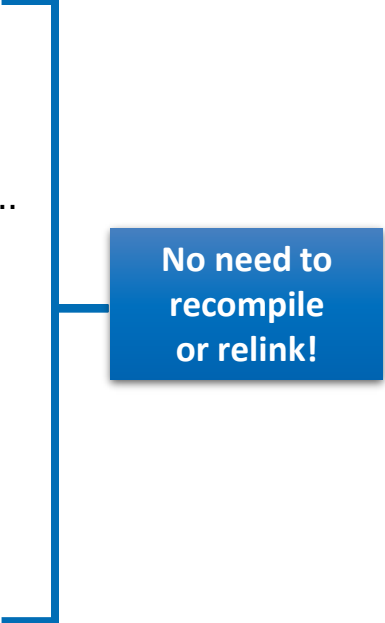
Lau Mercadal
(tools@bsc.es)
Barcelona Supercomputing Center

Getting a trace with Extrae

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



Extræe on IVYMUC (I)

- Log into IVYMUC:

```
laptop$ ssh -Y <USER>@lxlogin10.lrz.de
```

- Extræe is available via modules...

```
ivymuc$ module use /lrz/sys/courses/vihps/modulefiles  
ivymuc$ module load extræe/3.8.3
```

- ... as are the other BSCTOOLS:

```
ivymuc$ module load wxparaver  
ivymuc$ module load clustering
```

Getting your first trace

- Provided folder `bsctools` in `/lrz/sys/courses/vihps/material` contains:
 - Application compiled for the Intel and GNU toolchains (`lulesh2.0-intel`, `lulesh2.0-gnu`)
 - Jobscripts to execute and trace (`job-intel.sh`, `job-gnu.sh`, `trace.sh`)
 - Configuration of the tracing tool (`extrae.xml`)
 - Already generated tracefiles (`traces/*.{pcf,prv,row}`)
 - Clustering analysis configuration file (`cluster.xml`)
- Copy this folder to your `$HOME` and you are ready to follow this hands-on tutorial

```
ivymuc$ cp -R /lrz/sys/courses/vihps/material/bsctools $HOME
```

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)

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

```
#SBATCH --job-name=lulesh2.0-intel_27p  
#SBATCH --output=%x.%j.out  
#SBATCH --error=%x.%j.err  
#SBATCH --clusters=ivymuc  
#SBATCH --nodes=2  
#SBATCH --ntasks=27  
#SBATCH --get-user-env  
#SBATCH --time=00:10:00
```

Request resources

```
module load slurm_setup
```

```
export OMP_NUM_THREADS=1
```

```
mpirun ../bin/lulesh2.0-intel -i 10 -s 65 -p
```

Run the program

Step 1: Adapt the job script to load Extrae

- Jobscript modified to load Extrae (extrae/job-intel.sh)

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

#SBATCH --job-name=lulesh2.0-intel_27p
#SBATCH --output=%x.%j.out
#SBATCH --error=%x.%j.err
#SBATCH --clusters=ivymuc
#SBATCH --nodes=2
#SBATCH --ntasks=27
#SBATCH --get-user-env
#SBATCH --time=00:10:00

module load slurm_setup

export OMP_NUM_THREADS=1
export TRACE_NAME=lulesh2.0-intel_27p.prv

mpirun ./trace.sh ../bin/lulesh2.0-intel -i 10 -s 65 -p
```

Run with Extrae

Step 1: Adapt the job script to load Extrae

- Tracing launcher helper script (`extrae/trace.sh`)

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

#SBATCH --job-name=lulesh2.0-intel_27p
#SBATCH --output=%x.%j.out
#SBATCH --error=%x.%j.err
#SBATCH --clusters=ivymuc
#SBATCH --nodes=2
#SBATCH --ntasks=27
#SBATCH --get-user-env
#SBATCH --time=00:10:00

module load slurm_setup

export OMP_NUM_THREADS=1
export TRACE_NAME=lulesh2.0-intel_27p.prv

mpirun ./trace.sh ./bin/lulesh2.0-intel -i 10 -s 65 -p
```

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

module use /lrz/sys/courses/vihps/modulefiles
module load extrae

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

## Run the desired program
$*
```

What to trace?

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

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 add suffix "f" in Fortran codes

Step 2: Extrae XML configuration

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

Instrument the MPI calls
(What's the program doing?)

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

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

Instrument the call-stack
(Where in my code?)

```
<callers enabled="yes">  
  <mpi enabled="yes">1-3</mpi>  
  <sampling enabled="no">1-5</sampling>  
</callers>
```

Step 2: Extrae XML configuration (II)

```
<counters enabled="yes">
  <cpu enabled="yes" starting-set-distribution="1">
    <set enabled="yes" domain="all" changeat-time="0">
      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?)

```
<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>
```

Extrae buffer size
(Flush/memory trade-off)

Additional sampling
(Want more details?)

Automatic
post-processing
to generate the
Paraver trace

```
ivymuc$ papi_best_set \
omnipresent <list-of-counters> \
<list-of-counters>
```

Step 3: Run it!

- Submit your job as usual

```
ivymuc$ sbatch --reservation=hhps1s21_workshop job-intel.sh
```

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

```
ivymuc$ ls -l
...
lulesh2.0-intel_27p.pcf
lulesh2.0-intel_27p.prv
lulesh2.0-intel_27p.row
```

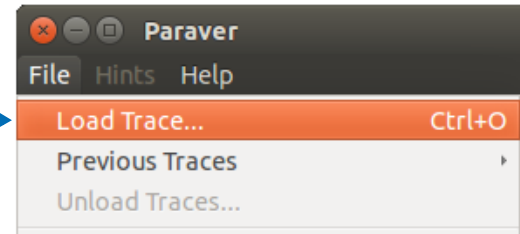
- Any trouble? There's a trace already generated under the "traces" folder
- Now let's look into it!

Analysing a trace with Extrae

First steps of analysis

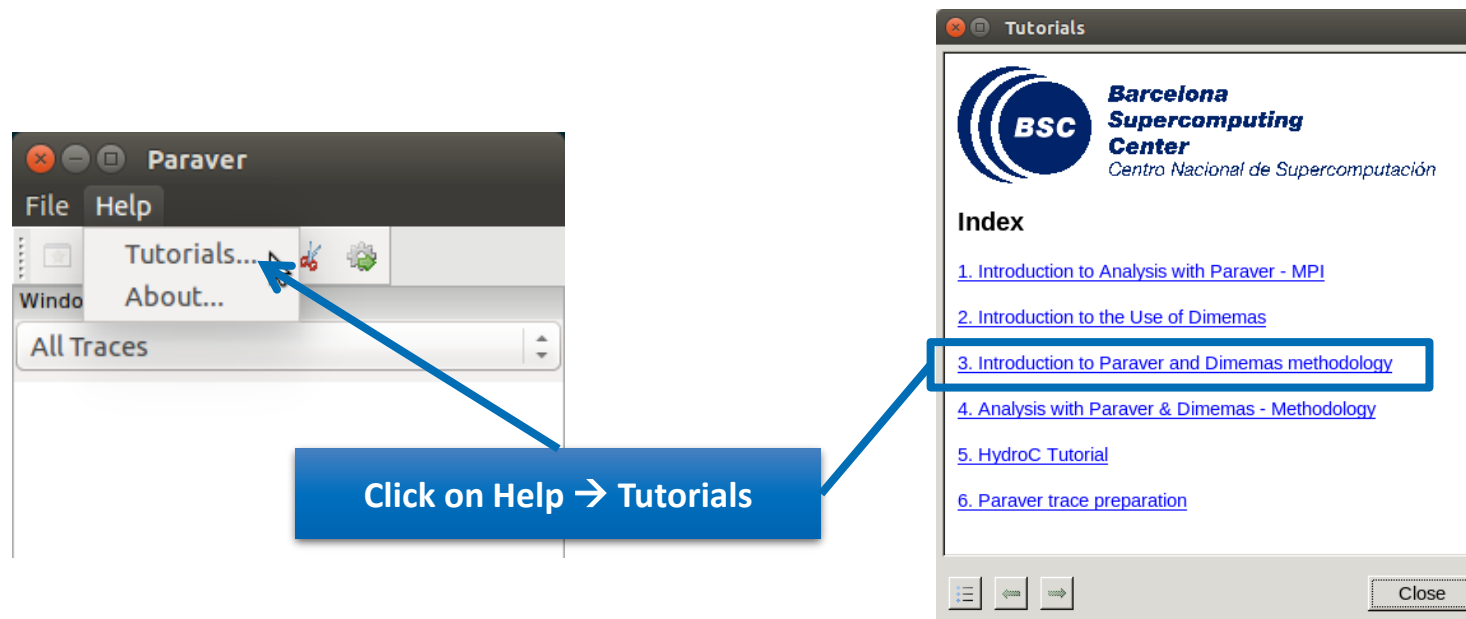
- Copy the trace to your computer
- Load the trace with Paraver

Click on File → Load Trace
→ Browse to “lulesh2.0-intel_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 labelled “**Outside MPI**”

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 `cfgs/mpi/mpi_stats.cfg`. 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 measure the computational load (instructions) distribution load the configuration file `cfgs/papi/2dh_useful_instructions.cfg`. 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.

Close

Parallel efficiency (Avg)

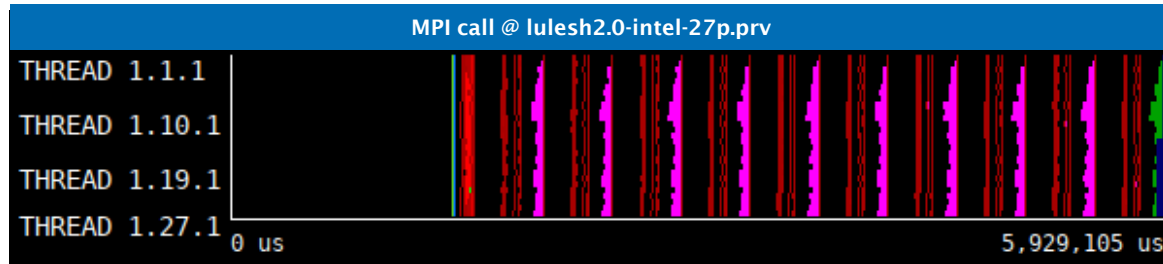
Comm efficiency (Max)

Load balance (Avg/Max)

MPI call profile @ lulesh2.0-intel-27p.prv

	Outside MPI	MPI_Isend	MPI_Irecv	MPI_Wait	MPI_Waitall	MPI_Barrier	MPI_Reduce
THREAD 1.15.1	90.65 %	0.93 %	0.12 %	0.35 %	0.14 %	0.05 %	0.15 %
THREAD 1.16.1	92.19 %	0.84 %	0.08 %	0.37 %	0.12 %	0.03 %	0.00 %
THREAD 1.17.1	91.62 %	0.89 %	0.11 %	0.43 %	0.19 %	0.06 %	0.03 %
THREAD 1.18.1	89.04 %	0.72 %	0.08 %	1.28 %	0.07 %	0.03 %	0.23 %
THREAD 1.19.1	92.20 %	0.88 %	0.05 %	0.21 %	0.10 %	0.03 %	0.03 %
THREAD 1.20.1	88.59 %	0.88 %	0.07 %	0.47 %	0.10 %	0.03 %	0.37 %
THREAD 1.21.1	86.85 %	0.81 %	0.05 %	1.56 %	0.08 %	0.03 %	0.43 %
THREAD 1.22.1	88.02 %	0.82 %	0.07 %	0.55 %	0.09 %	0.03 %	0.42 %
THREAD 1.23.1	89.31 %	1.01 %	0.10 %	2.84 %	0.35 %	0.10 %	0.00 %
THREAD 1.24.1	85.73 %	0.80 %	0.07 %	3.96 %	0.01 %	0.03 %	0.31 %
THREAD 1.25.1	87.79 %	0.53 %	0.05 %	3.68 %	0.01 %	0.03 %	0.16 %
THREAD 1.26.1	86.03 %	0.36 %	0.06 %	4.02 %	0.01 %	0.03 %	0.32 %
THREAD 1.27.1	85.52 %	0.46 %	0.04 %	3.27 %	0.01 %	0.03 %	0.44 %
Total	2,414.52 %	19.35 %	2.34 %	47.45 %	1.69 %	1.15 %	13.08 %
Average	89.43 %	0.72 %	0.09 %	1.76 %	0.06 %	0.04 %	0.48 %
Maximum	98.21 %	1.20 %	0.18 %	4.02 %	0.35 %	0.10 %	1.38 %
Minimum	82.46 %	0.09 %	0.04 %	0.21 %	0.01 %	0.01 %	0.00 %
StDev	3.78 %	0.27 %	0.03 %	1.19 %	0.07 %	0.02 %	0.39 %
Avg/Max	0.91	0.60	0.48	0.44	0.18	0.41	0.35

Focus on the iterative part

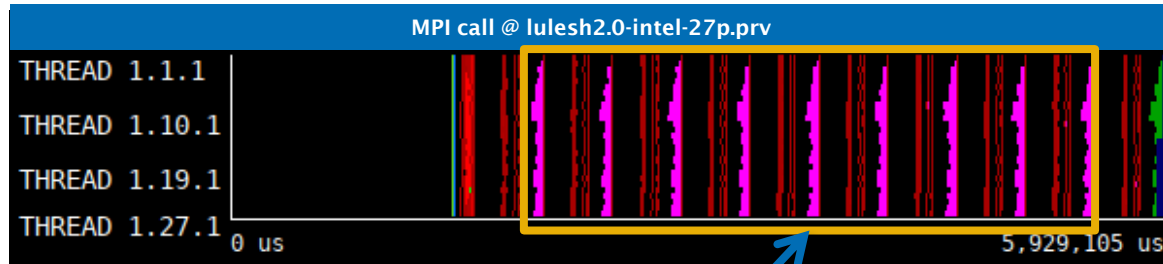


Click on
Open Control Window

MPI call profile @ lulesh2.0-intel-27p.prv

	Outside MPI	MPI_Isend	MPI_Irecv	MPI_Wait	MPI_Waitall	MPI_Barrier	MPI_Reduce
THREAD 1.15.1	90.65 %	0.93 %	0.12 %	0.35 %	0.14 %	0.05 %	0.15 %
THREAD 1.16.1	92.19 %	0.84 %	0.08 %	0.37 %	0.12 %	0.03 %	0.00 %
THREAD 1.17.1	91.62 %	0.89 %	0.11 %	0.43 %	0.19 %	0.06 %	0.03 %
THREAD 1.18.1	89.04 %	0.72 %	0.08 %	1.28 %	0.07 %	0.03 %	0.23 %
THREAD 1.19.1	92.20 %	0.88 %	0.05 %	0.21 %	0.10 %	0.03 %	0.03 %
THREAD 1.20.1	88.59 %	0.88 %	0.07 %	0.47 %	0.10 %	0.03 %	0.37 %
THREAD 1.21.1	86.85 %	0.81 %	0.05 %	1.56 %	0.08 %	0.03 %	0.43 %
THREAD 1.22.1	88.02 %	0.82 %	0.07 %	0.55 %	0.09 %	0.03 %	0.42 %
THREAD 1.23.1	89.31 %	1.01 %	0.10 %	2.84 %	0.35 %	0.10 %	0.00 %
THREAD 1.24.1	85.73 %	0.80 %	0.07 %	3.96 %	0.01 %	0.03 %	0.31 %
THREAD 1.25.1	87.79 %	0.53 %	0.05 %	3.68 %	0.01 %	0.03 %	0.16 %
THREAD 1.26.1	86.03 %	0.36 %	0.06 %	4.02 %	0.01 %	0.03 %	0.32 %
THREAD 1.27.1	85.52 %	0.46 %	0.04 %	3.27 %	0.01 %	0.03 %	0.44 %
Total	2,414.52 %	19.35 %	2.34 %	47.45 %	1.69 %	1.15 %	13.08 %
Average	89.43 %	0.72 %	0.09 %	1.76 %	0.06 %	0.04 %	0.48 %
Maximum	98.21 %	1.20 %	0.18 %	4.02 %	0.35 %	0.10 %	1.38 %
Minimum	82.46 %	0.09 %	0.04 %	0.21 %	0.01 %	0.01 %	0.00 %
StDev	3.78 %	0.27 %	0.03 %	1.19 %	0.07 %	0.02 %	0.39 %
Avg/Max	0.91	0.60	0.48	0.44	0.18	0.41	0.35

Focus on the iterative part

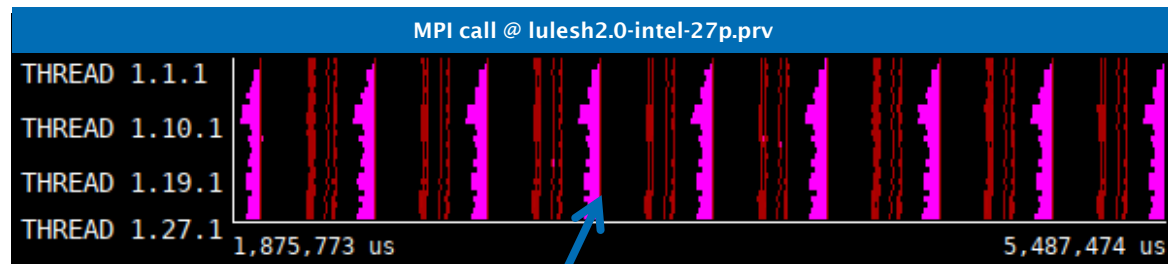


Drag & drop on this area
to zoom on the iterative region

MPI call profile @ lulesh2.0-gnu-27p.prv

	Outside MPI	MPI_Isend	MPI_Irecv	MPI_Wait	MPI_Waitall	MPI_Barrier	MPI_Reduce
THREAD 1.15.1	90.65 %	0.93 %	0.12 %	0.35 %	0.14 %	0.05 %	0.15 %
THREAD 1.16.1	92.19 %	0.84 %	0.08 %	0.37 %	0.12 %	0.03 %	0.00 %
THREAD 1.17.1	91.62 %	0.89 %	0.11 %	0.43 %	0.19 %	0.06 %	0.03 %
THREAD 1.18.1	89.04 %	0.72 %	0.08 %	1.28 %	0.07 %	0.03 %	0.23 %
THREAD 1.19.1	92.20 %	0.88 %	0.05 %	0.21 %	0.10 %	0.03 %	0.03 %
THREAD 1.20.1	88.59 %	0.88 %	0.07 %	0.47 %	0.10 %	0.03 %	0.37 %
THREAD 1.21.1	86.85 %	0.81 %	0.05 %	1.56 %	0.08 %	0.03 %	0.43 %
THREAD 1.22.1	88.02 %	0.82 %	0.07 %	0.55 %	0.09 %	0.03 %	0.42 %
THREAD 1.23.1	89.31 %	1.01 %	0.10 %	2.84 %	0.35 %	0.10 %	0.00 %
THREAD 1.24.1	85.73 %	0.80 %	0.07 %	3.96 %	0.01 %	0.03 %	0.31 %
THREAD 1.25.1	87.79 %	0.53 %	0.05 %	3.68 %	0.01 %	0.03 %	0.16 %
THREAD 1.26.1	86.03 %	0.36 %	0.06 %	4.02 %	0.01 %	0.03 %	0.32 %
THREAD 1.27.1	85.52 %	0.46 %	0.04 %	3.27 %	0.01 %	0.03 %	0.44 %
Total	2,414.52 %	19.35 %	2.34 %	47.45 %	1.69 %	1.15 %	13.08 %
Average	89.43 %	0.72 %	0.09 %	1.76 %	0.06 %	0.04 %	0.48 %
Maximum	98.21 %	1.20 %	0.18 %	4.02 %	0.35 %	0.10 %	1.38 %
Minimum	82.46 %	0.09 %	0.04 %	0.21 %	0.01 %	0.01 %	0.00 %
StDev	3.78 %	0.27 %	0.03 %	1.19 %	0.07 %	0.02 %	0.39 %
Avg/Max	0.91	0.60	0.48	0.44	0.18	0.41	0.35

Recalculate efficiency of iterative region

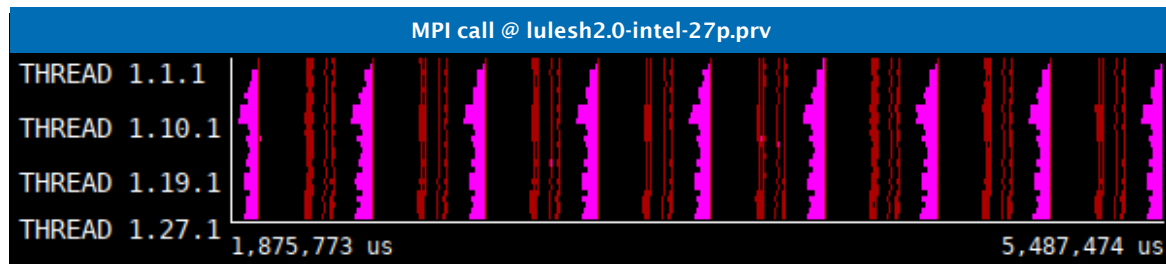


Right click →
Copy

MPI call profile @ lulesh2.0-intel-27p.prv

	Outside MPI	MPI_Isend	MPI_Irecv	MPI_Wait	MPI_Waitall	MPI_Barrier	MPI_Reduce
THREAD 1.15.1	90.65 %	0.93 %	0.12 %	0.35 %	0.14 %	0.05 %	0.15 %
THREAD 1.16.1	92.19 %	0.84 %	0.08 %	0.37 %	0.12 %	0.03 %	0.00 %
THREAD 1.17.1	91.62 %	0.89 %	0.11 %	0.43 %	0.19 %	0.06 %	0.03 %
THREAD 1.18.1	89.04 %	0.72 %	0.08 %	1.28 %	0.07 %	0.03 %	0.23 %
THREAD 1.19.1	92.20 %	0.88 %	0.05 %	0.21 %	0.10 %	0.03 %	0.03 %
THREAD 1.20.1	88.59 %	0.88 %	0.07 %	0.47 %	0.10 %	0.03 %	0.37 %
THREAD 1.21.1	86.85 %	0.81 %	0.05 %	1.56 %	0.08 %	0.03 %	0.43 %
THREAD 1.22.1	88.02 %	0.82 %	0.07 %	0.55 %	0.09 %	0.03 %	0.42 %
THREAD 1.23.1	89.31 %	1.01 %	0.10 %	2.84 %	0.35 %	0.10 %	0.00 %
THREAD 1.24.1	85.73 %	0.80 %	0.07 %	3.96 %	0.01 %	0.03 %	0.31 %
THREAD 1.25.1	87.79 %	0.53 %	0.05 %	3.68 %	0.01 %	0.03 %	0.16 %
THREAD 1.26.1	86.03 %	0.36 %	0.06 %	4.02 %	0.01 %	0.03 %	0.32 %
THREAD 1.27.1	85.52 %	0.46 %	0.04 %	3.27 %	0.01 %	0.03 %	0.44 %
Total	2,414.52 %	19.35 %	2.34 %	47.45 %	1.69 %	1.15 %	13.08 %
Average	89.43 %	0.72 %	0.09 %	1.76 %	0.06 %	0.04 %	0.48 %
Maximum	98.21 %	1.20 %	0.18 %	4.02 %	0.35 %	0.10 %	1.38 %
Minimum	82.46 %	0.09 %	0.04 %	0.21 %	0.01 %	0.01 %	0.00 %
StDev	3.78 %	0.27 %	0.03 %	1.19 %	0.07 %	0.02 %	0.39 %
Avg/Max	0.91	0.60	0.48	0.44	0.18	0.41	0.35

Recalculate efficiency of iterative region



Right click →
 Paste → Time

MPI call profile @ lulesh2.0-intel-27p.prv

	Outside MPI	MPI_Isend	MPI_Irecv	MPI_Wait	MPI_Waitall	MPI_Allreduce	MPI_Comm
THREAD 1.15.1	87.92 %	0.27 %	0.16 %	0.20 %	0.16 %	11.21 %	
THREAD 1.16.1	90.05 %	0.20 %	0.10 %	0.29 %	0.11 %	9.18 %	
THREAD 1.17.1	89.28 %	0.29 %	0.15 %	0.40 %	0.12 %	9.68 %	
THREAD 1.18.1	85.52 %	0.25 %	0.10 %	1.21 %	0.08 %	12.75 %	
THREAD 1.19.1	90.15 %	0.13 %	0.07 %	0.09 %	0.09 %	9.40 %	
THREAD 1.20.1	84.83 %	0.21 %	0.10 %	0.37 %	0.09 %	14.33 %	
THREAD 1.21.1	82.19 %	0.20 %	0.07 %	1.94 %	0.08 %	15.43 %	
THREAD 1.22.1	84.00 %	0.22 %	0.09 %	0.38 %	0.09 %	15.15 %	
THREAD 1.23.1	86.42 %	0.35 %	0.14 %	3.52 %	0.46 %	9.05 %	
THREAD 1.24.1	81.07 %	0.28 %	0.09 %	4.82 %	0.02 %	13.64 %	
THREAD 1.25.1	83.80 %	0.21 %	0.06 %	4.38 %	0.02 %	11.46 %	
THREAD 1.26.1	81.47 %	0.26 %	0.08 %	4.28 %	0.02 %	13.81 %	
THREAD 1.27.1	80.69 %	0.21 %	0.05 %	3.40 %	0.02 %	15.57 %	
Total	2,328.06 %	6.03 %	3.06 %	49.86 %	1.83 %	309.06 %	
Average	86.22 %	0.22 %	0.11 %	1.85 %	0.07 %	11.45 %	
Maximum	98.82 %	0.40 %	0.24 %	4.82 %	0.46 %	20.17 %	
Minimum	76.18 %	0.11 %	0.05 %	0.09 %	0.02 %	0.03 %	
StDev	5.41 %	0.06 %	0.04 %	1.51 %	0.09 %	4.52 %	
Avg/Max	0.87	0.56	0.47	0.38	0.15	0.57	

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 (Avg)

Comm efficiency (Max)

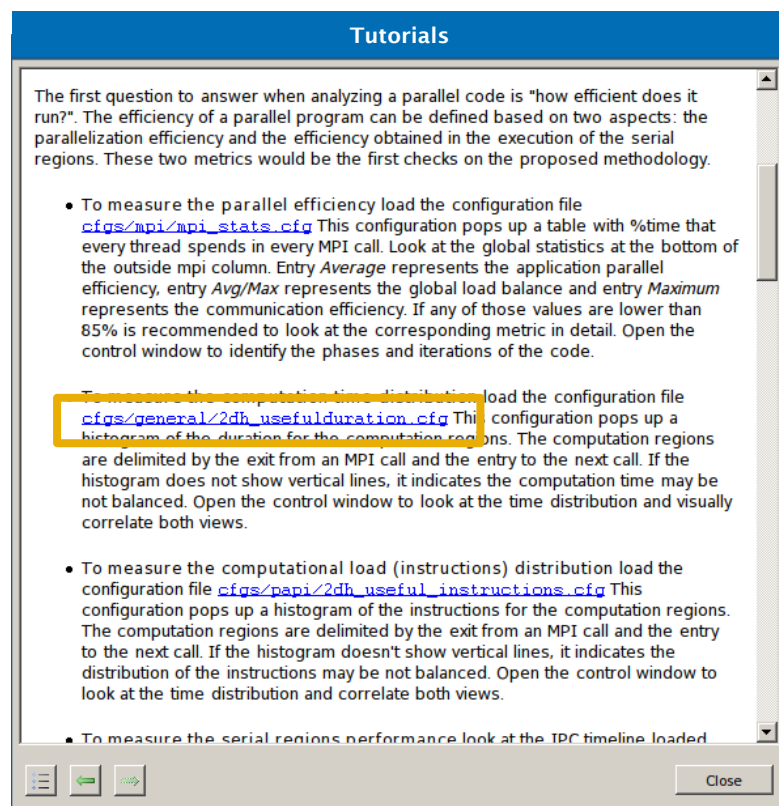
Load balance (Avg/Max)

MPI call profile @ lulesh2.0-intel-27p.prv

	Outside MPI	MPI_Isend	MPI_Irecv	MPI_Wait	MPI_Waitall	MPI_Allreduce	MPI_Comm
THREAD 1.15.1	87.92 %	0.27 %	0.16 %	0.20 %	0.16 %	11.21 %	
THREAD 1.16.1	90.05 %	0.20 %	0.10 %	0.29 %	0.11 %	9.18 %	
THREAD 1.17.1	89.28 %	0.29 %	0.15 %	0.40 %	0.12 %	9.68 %	
THREAD 1.18.1	85.52 %	0.25 %	0.10 %	1.21 %	0.08 %	12.75 %	
THREAD 1.19.1	90.15 %	0.13 %	0.07 %	0.09 %	0.09 %	9.40 %	
THREAD 1.20.1	84.83 %	0.21 %	0.10 %	0.37 %	0.09 %	14.33 %	
THREAD 1.21.1	82.19 %	0.20 %	0.07 %	1.94 %	0.08 %	15.43 %	
THREAD 1.22.1	84.00 %	0.22 %	0.09 %	0.38 %	0.09 %	15.15 %	
THREAD 1.23.1	86.42 %	0.35 %	0.14 %	3.52 %	0.46 %	9.05 %	
THREAD 1.24.1	81.07 %	0.28 %	0.09 %	4.82 %	0.02 %	13.64 %	
THREAD 1.25.1	83.80 %	0.21 %	0.06 %	4.38 %	0.02 %	11.46 %	
THREAD 1.26.1	81.47 %	0.26 %	0.08 %	4.28 %	0.02 %	13.81 %	
THREAD 1.27.1	80.69 %	0.21 %	0.05 %	3.40 %	0.02 %	15.57 %	
Total	2,328.06 %	6.03 %	3.06 %	49.86 %	1.83 %	309.06 %	
Average	86.22 %	0.22 %	0.11 %	1.85 %	0.07 %	11.45 %	
Maximum	98.82 %	0.40 %	0.24 %	4.82 %	0.46 %	20.17 %	
Minimum	76.18 %	0.11 %	0.05 %	0.09 %	0.02 %	0.03 %	
StDev	5.41 %	0.06 %	0.04 %	1.51 %	0.09 %	4.52 %	
Avg/Max	0.87	0.56	0.47	0.38	0.15	0.57	

Computation time distribution

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

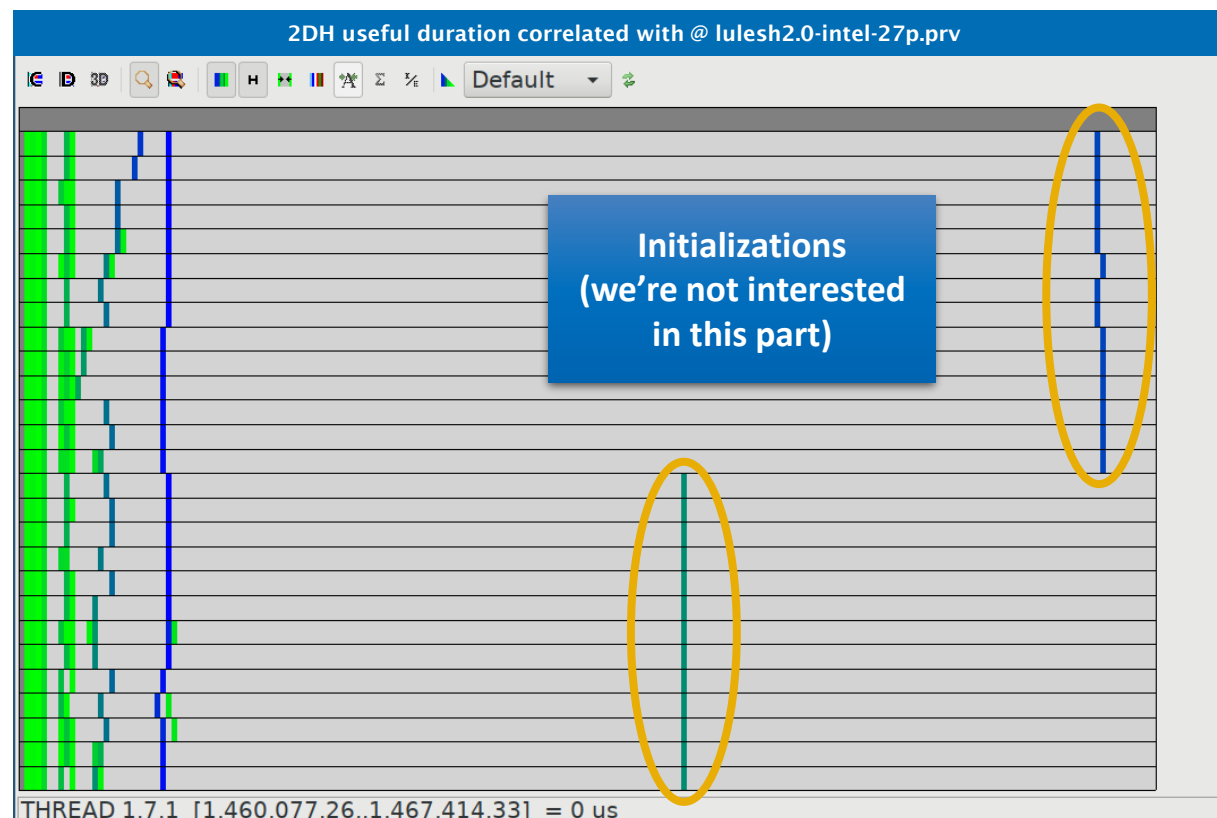


Tutorials

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 [cfigs/mpi/mpi_stats.cfg](#) 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 [cfigs/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 measure the computational load (instructions) distribution load the configuration file [cfigs/papi/2dh_useful_instructions.cfg](#) 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.

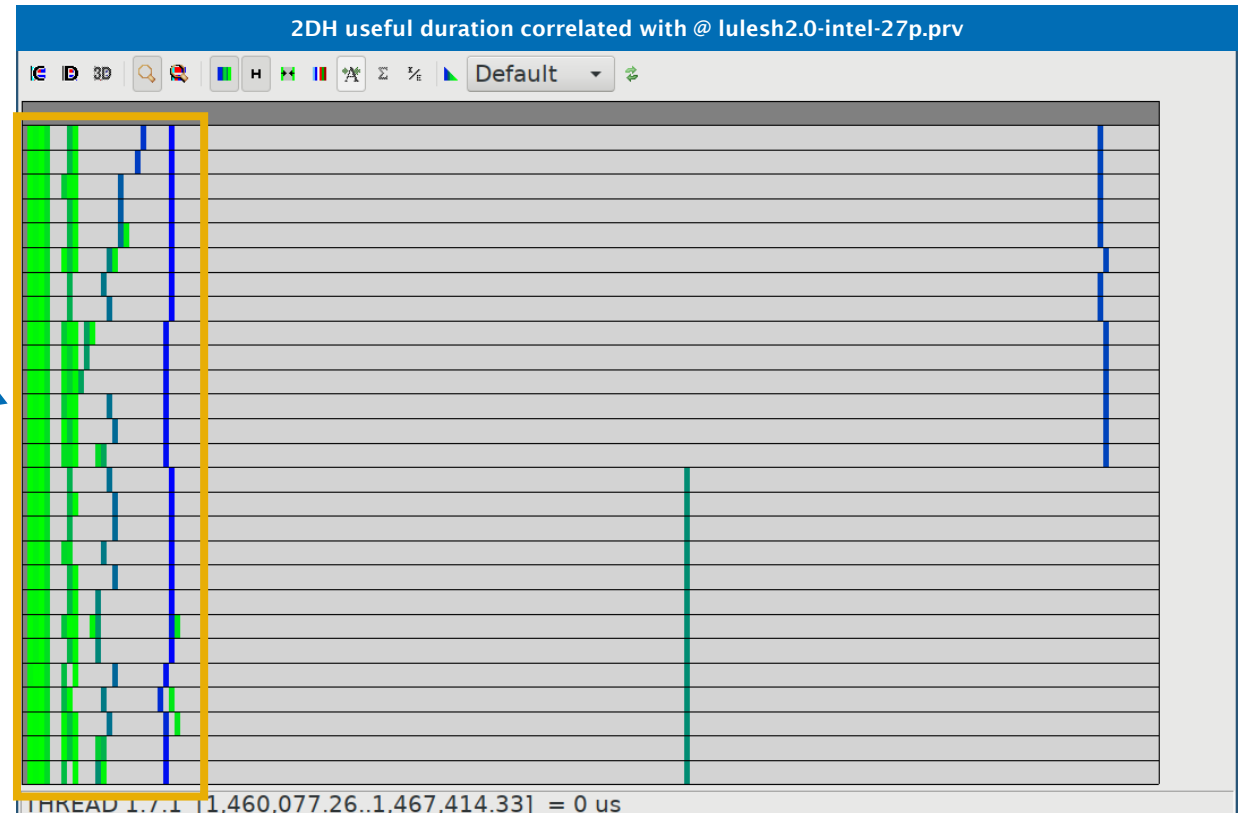
Close



Focus on the iterative part

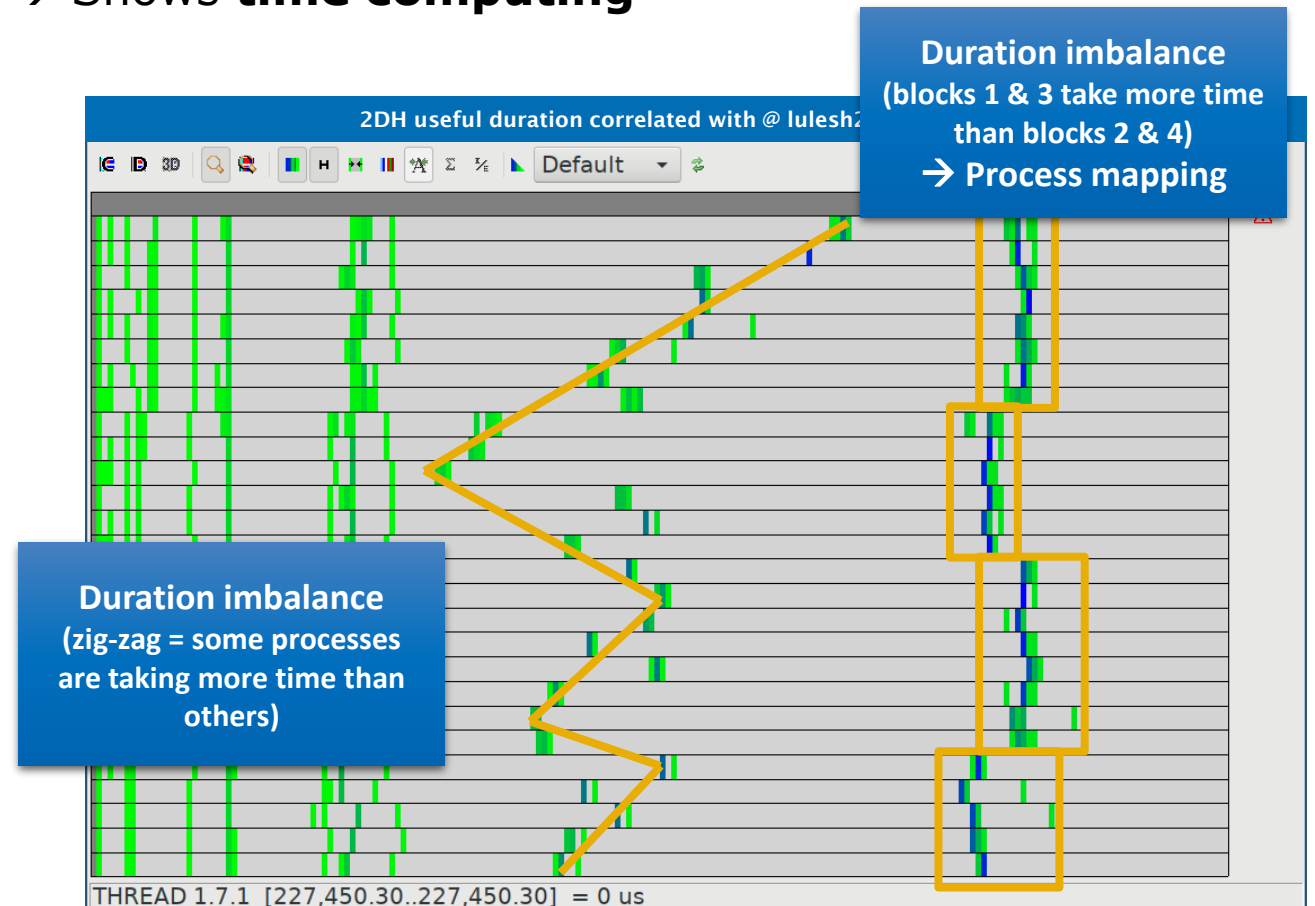
- Click on “2dh_usefulduration.cfg” (2nd link) → Shows **time computing**

Drag & drop on this area
to skip the initializations



Computation time distribution

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



Computation load distribution

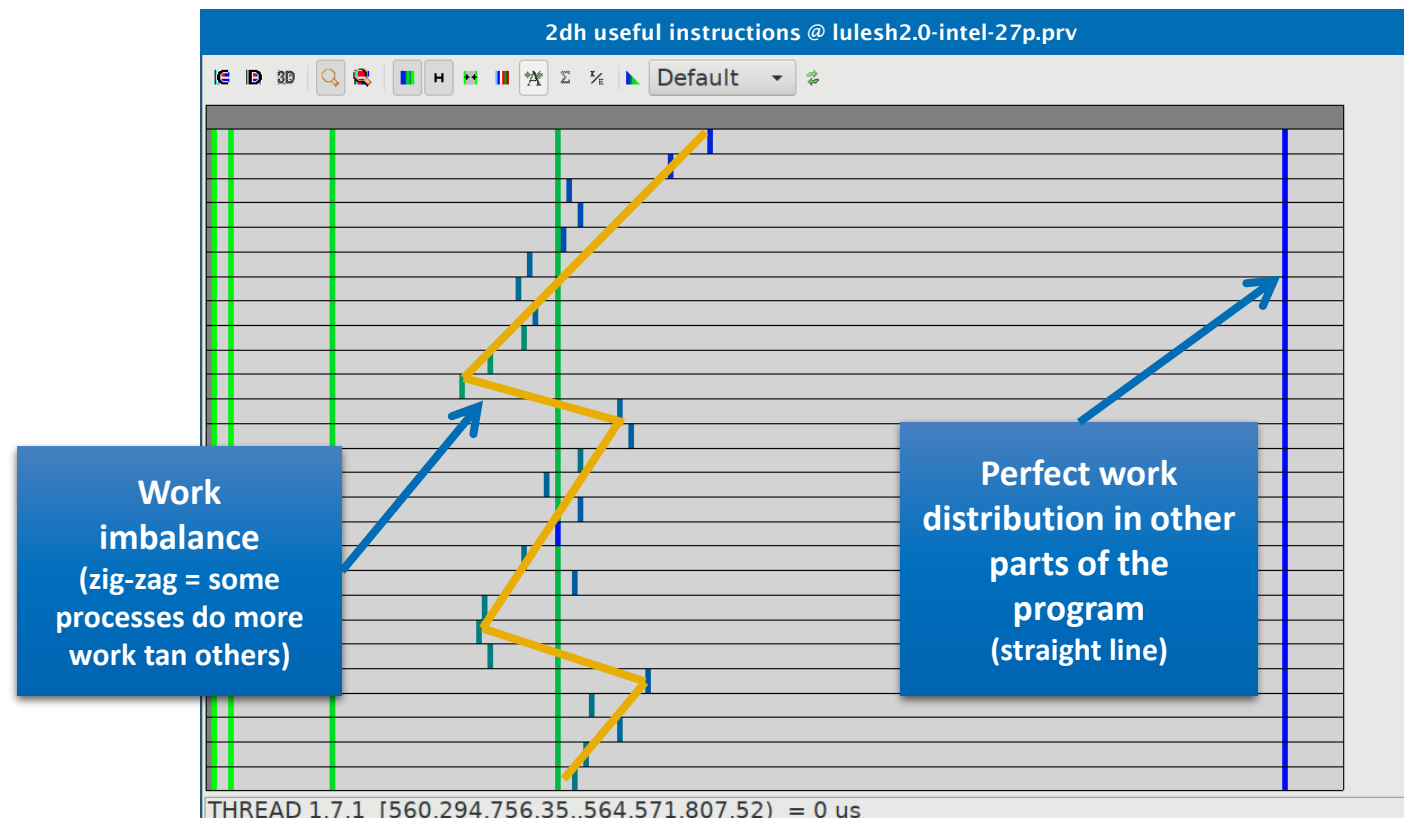
- Click on “2dh_useful_instructions.cfg” (3rd link) → Shows **amount of work**

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.

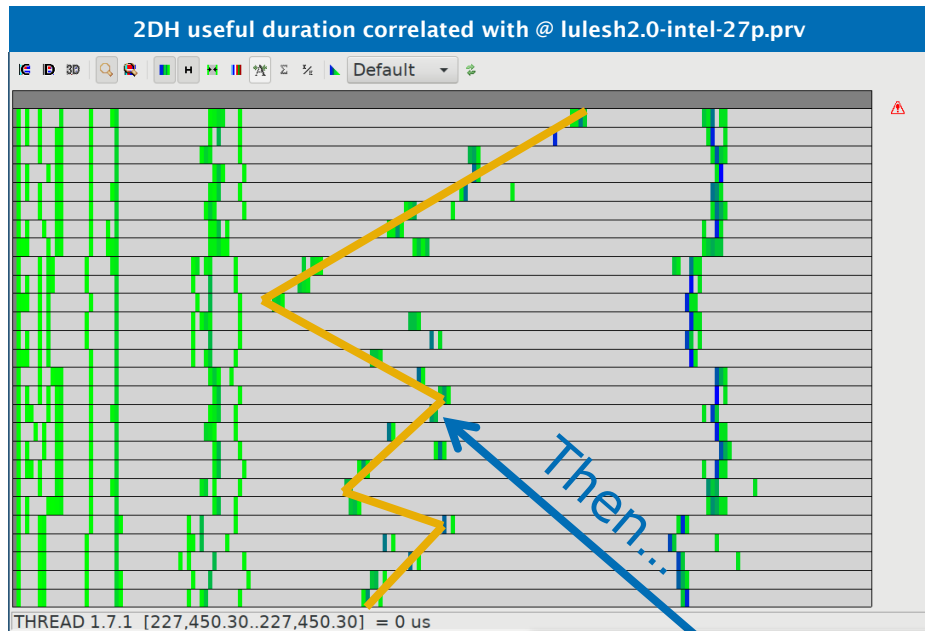
- To measure the parallel efficiency load the configuration file [cfigs/mpi/mpi_stats.cfg](#). This configuration pops up a table with %time that every thread spends in every MPI call. Look at the global statistics at the bottom of the outside mpi column. Entry *Average* represents the application parallel 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 [cfigs/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 measure the **computational load (instructions) distribution** load the configuration file [cfigs/papi/2dh_useful_instructions.cfg](#). 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.

Close



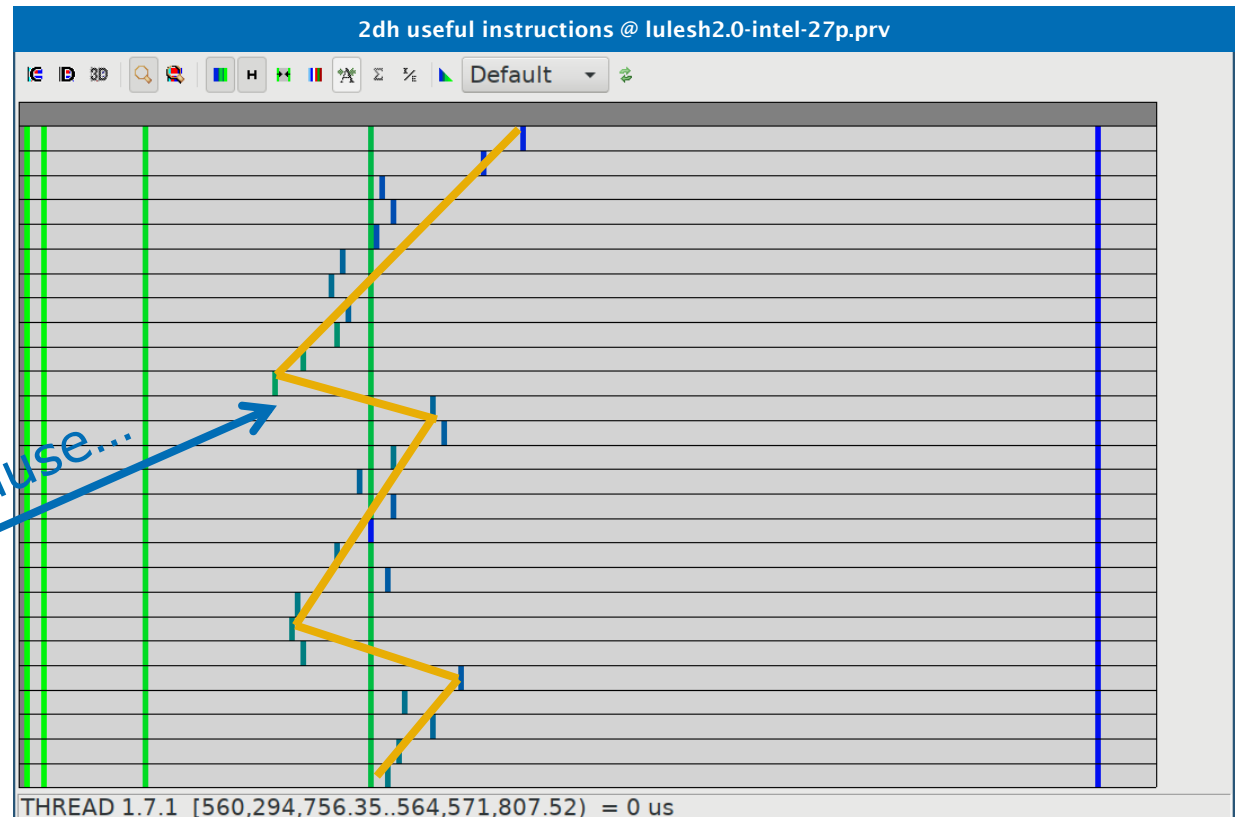
Computation load distribution

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



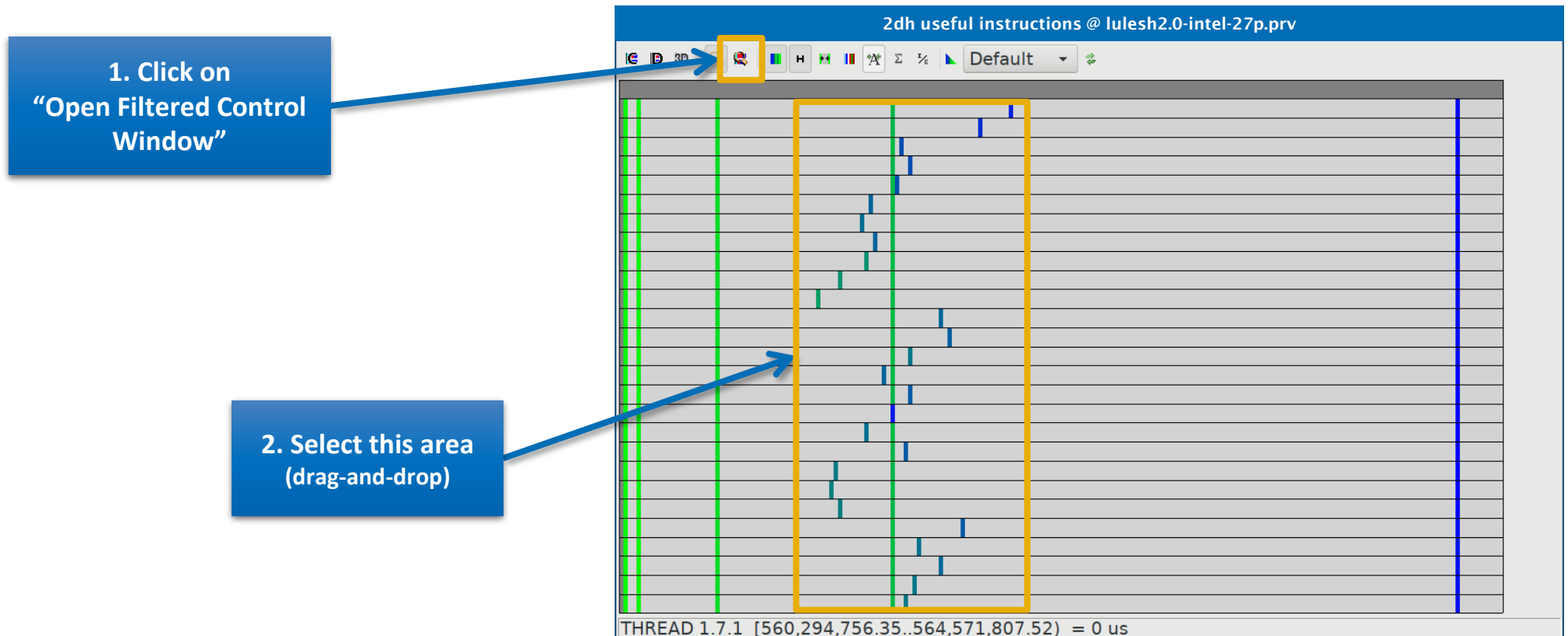
Work
imbalance → Slow/fast
processes
(zig-zag)

Because...



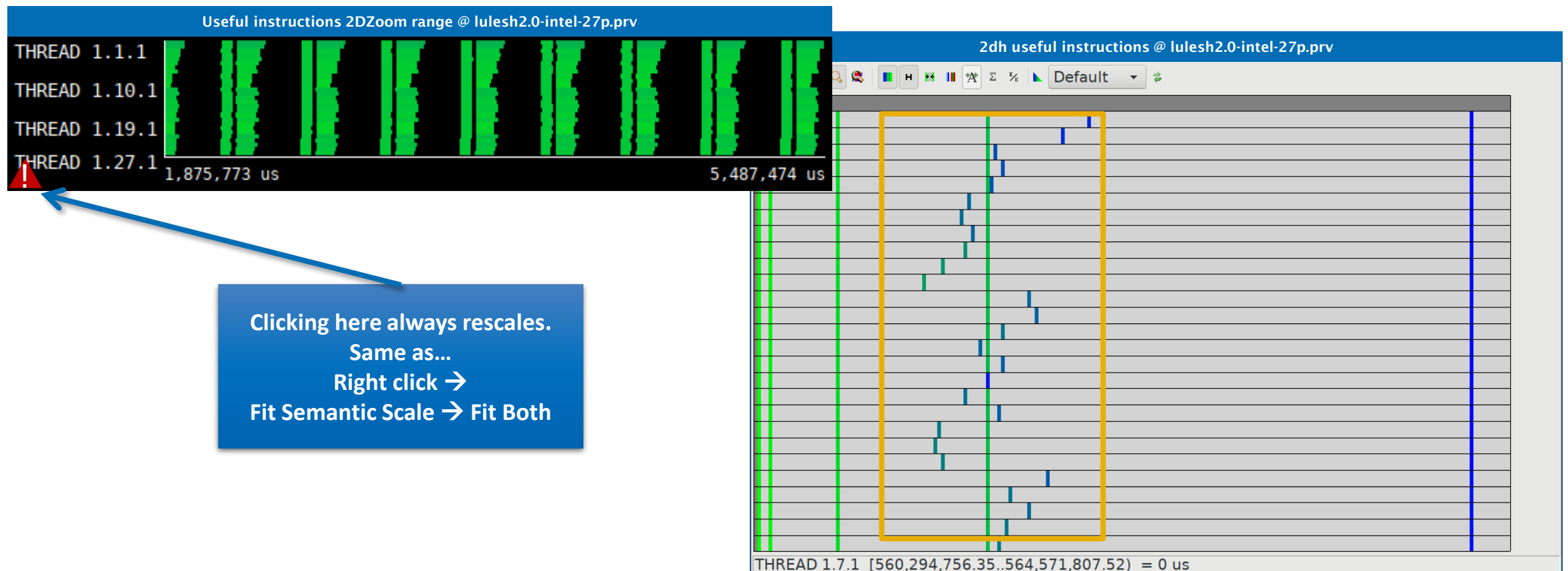
Where does this happen?

- Go from the table to the timeline



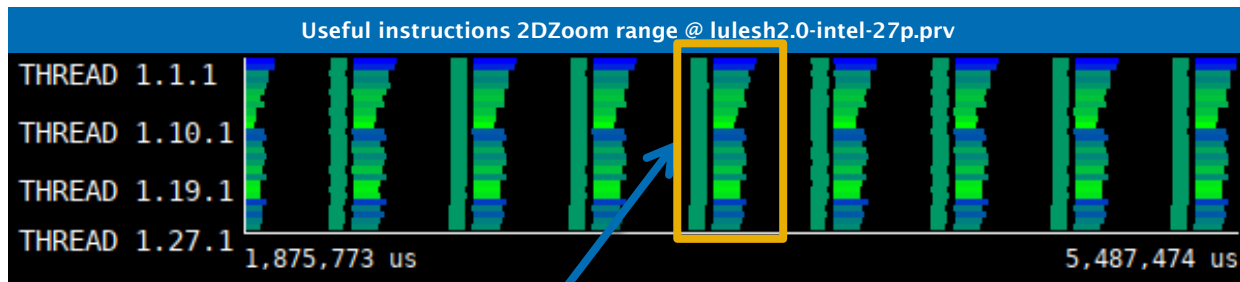
Where does this happen?

- Go from the table to the timeline



Where does this happen?

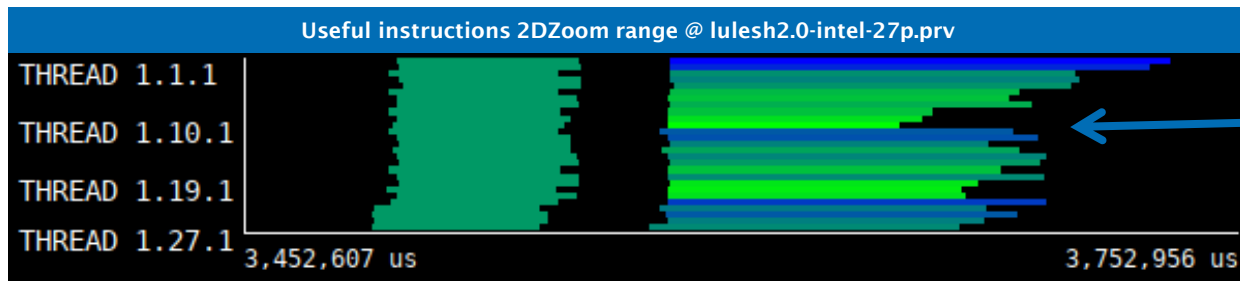
- **Slow** & **Fast** at the same time? → Imbalance



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

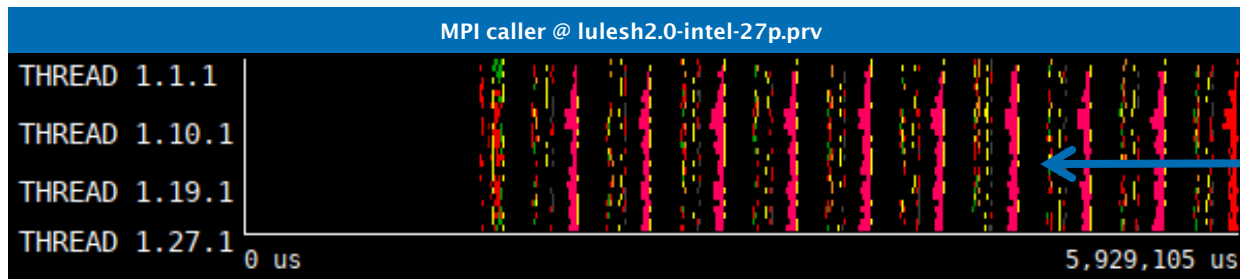
Where does this happen?

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



1. Right click → Copy

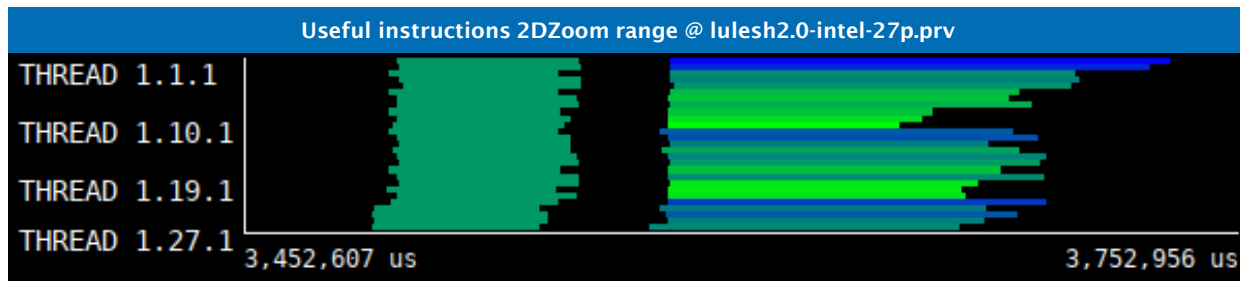
- **Hints** → Call stack references → Caller function



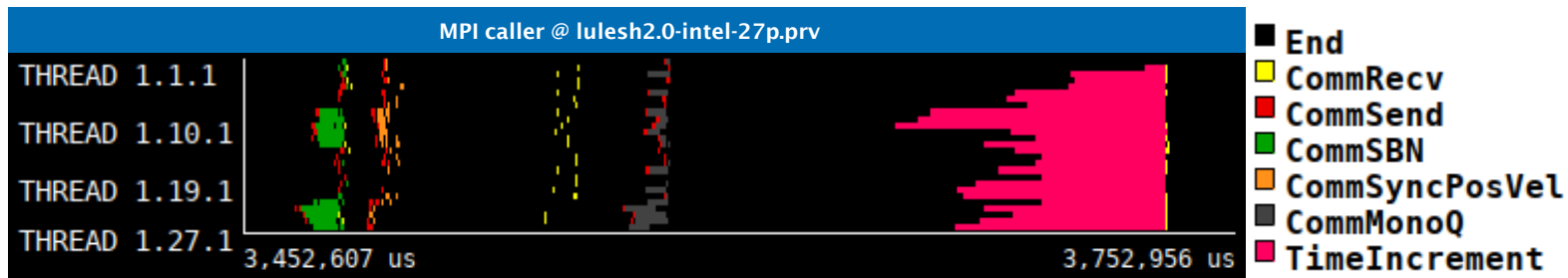
2. Right click →
Paste → Time

Where does this happen?

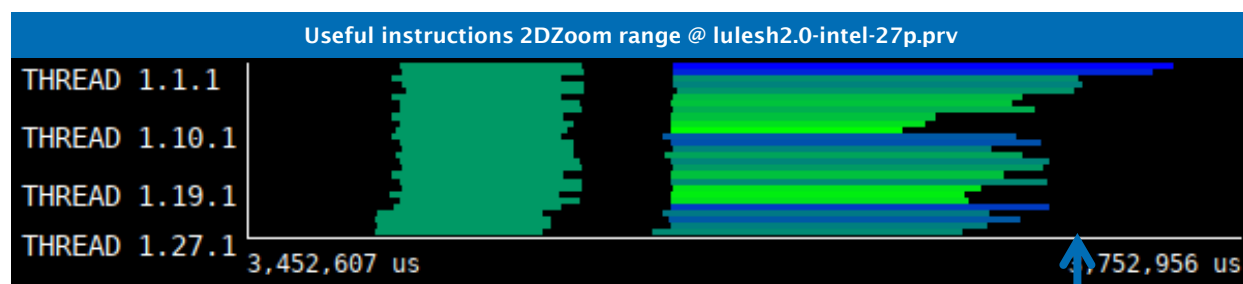
- **Slow** & **Fast** at the same time? → Imbalance



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



Save CFG's (method 1)

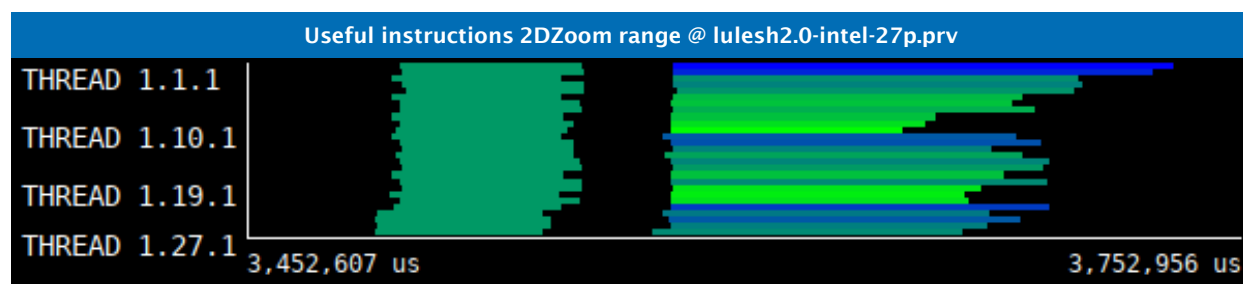


Right click on
timeline

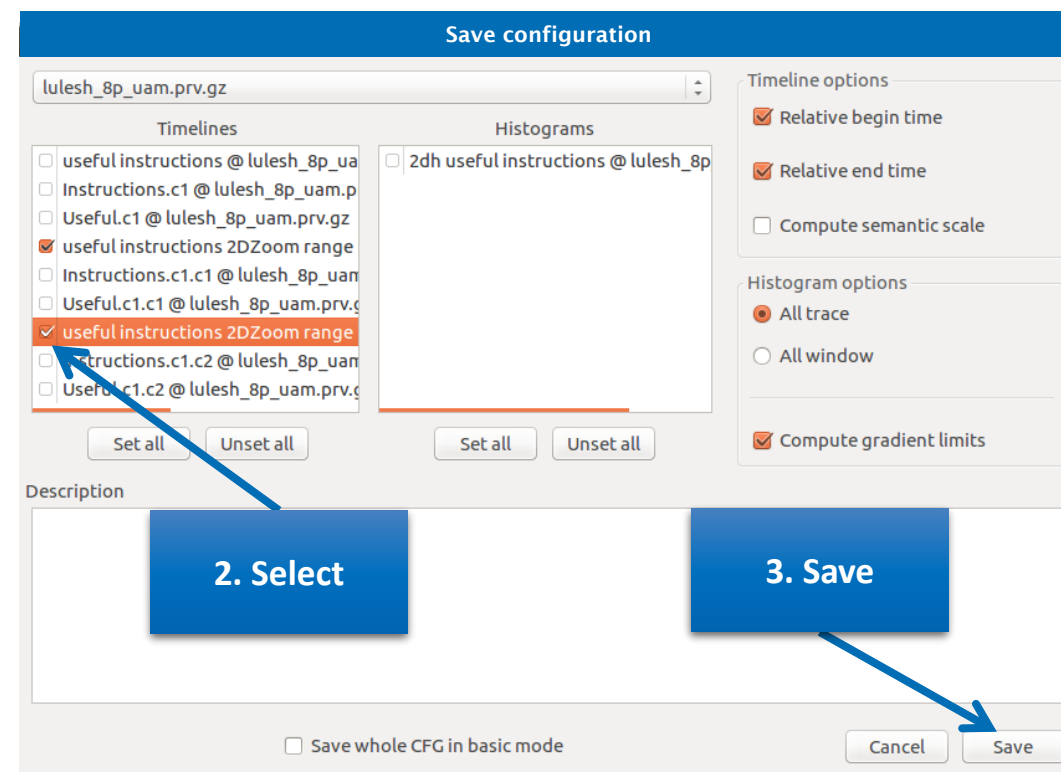
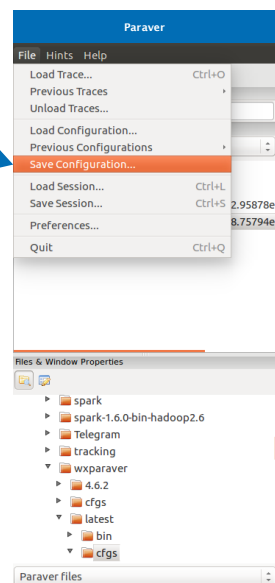
- Copy Ctrl+C
- Paste
- Clone
- Undo Zoom Ctrl+U
- Redo Zoom Ctrl+R
- Fit Time Scale
- Fit Semantic Scale
- Fit Objects
- Select Objects...
- View
- Paint As
- Drawmode
- Pixel Size
- Object Labels
- Object Axis
- Run
- Synchronize
- Remove all sync
- Save**
- Info Panel

- Configuration...
- Image...
- Image Legend...
- Text...

Save CFG's (method 2)

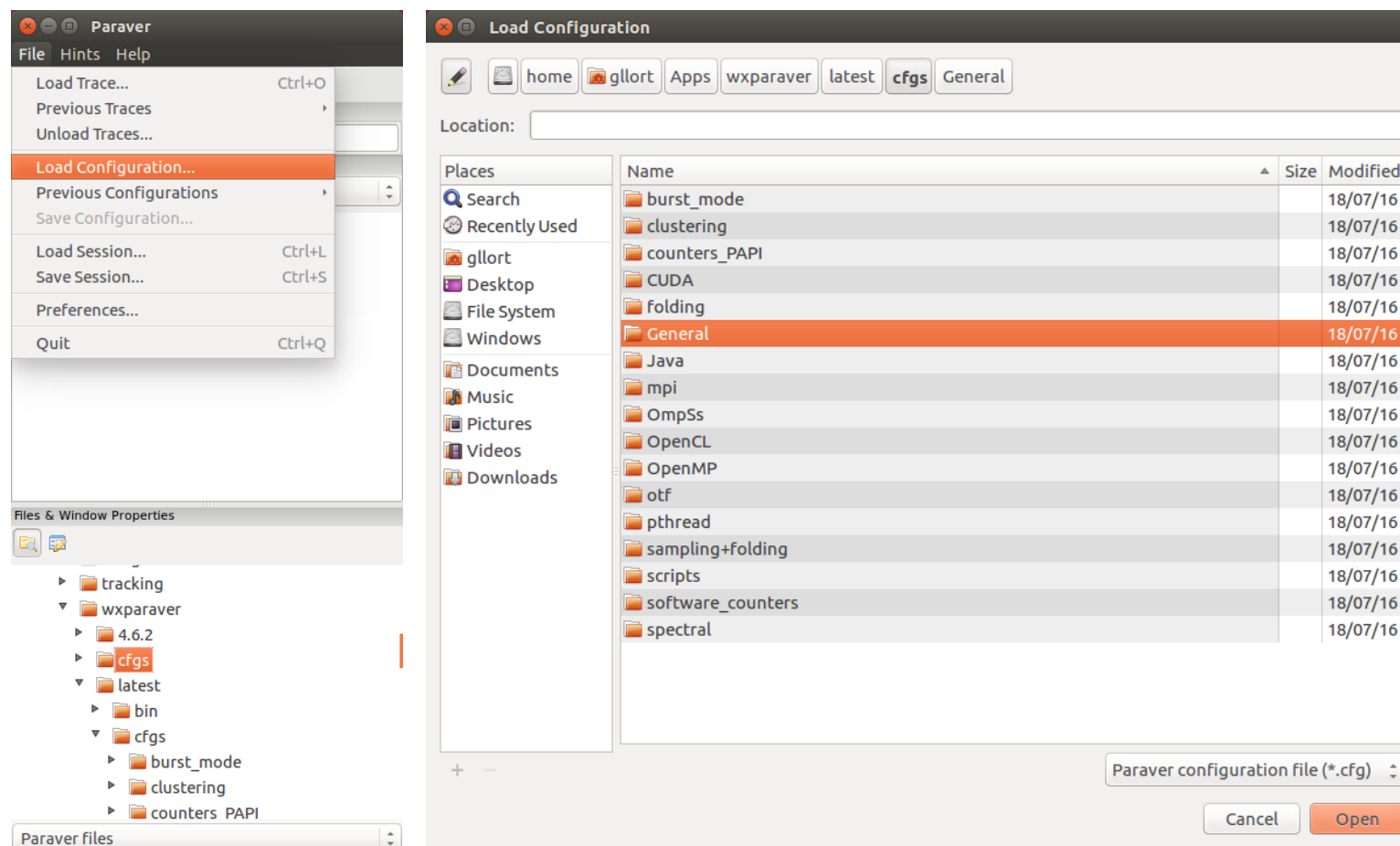


1. Main Paraver window



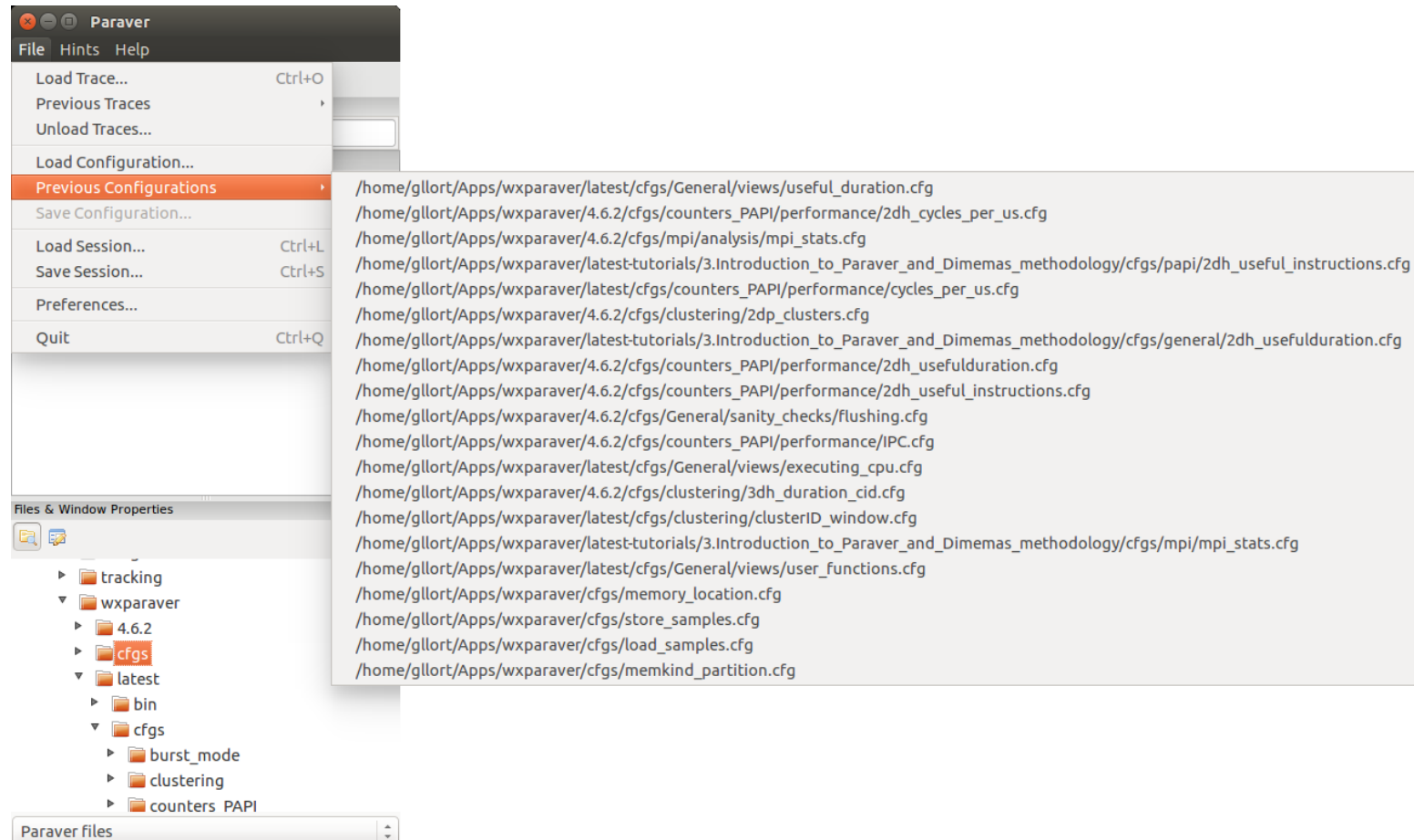
CFG's distribution

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



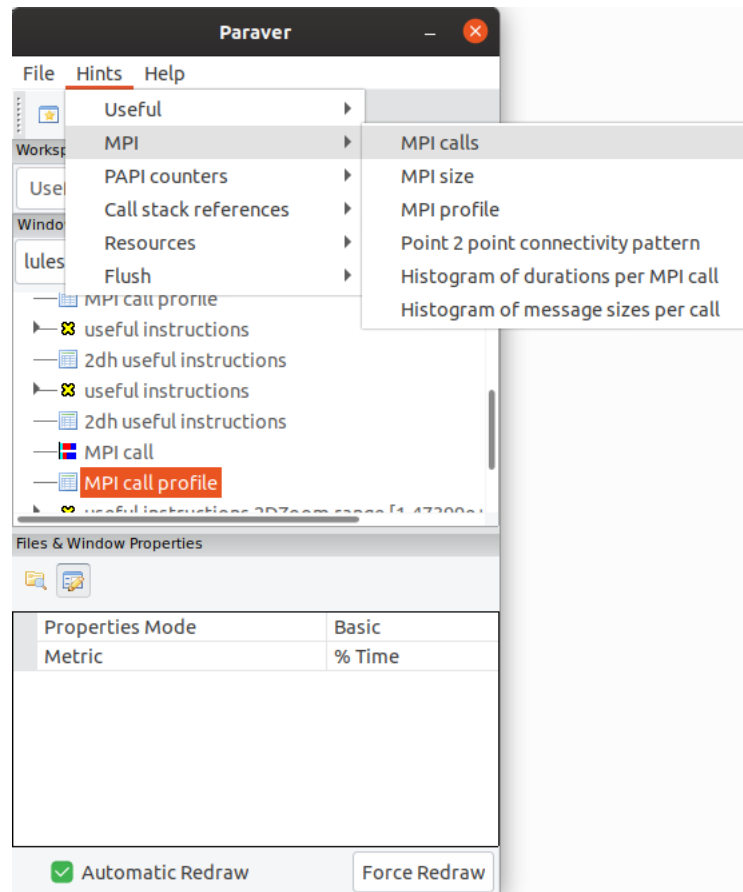
CFG's distribution

- 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 `$EXTRAE_HOME/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?

Cluster-based analysis

Use clustering analysis

- Run clustering

```
ivymuc> cd $HOME/bsctools/clustering
ivymuc> module load clustering
ivymuc> BurstClustering -d cluster.xml \
-i ../extrae/lulesh2.0-intel_27p.prv \
-o lulesh2.0-intel_27p_clustered.prv
```

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

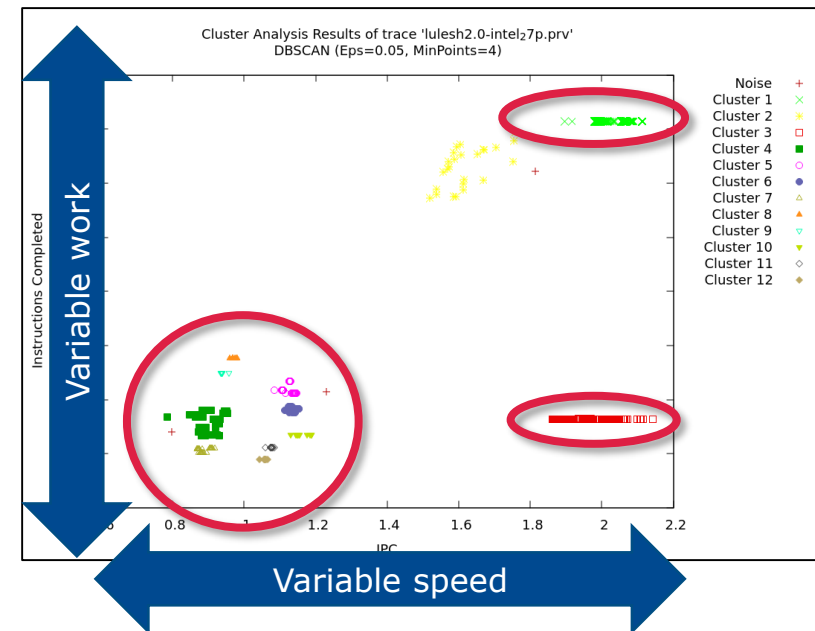
```
ivymuc> ls $HOME/bsctools/traces/lulesh2.0-intel_27p.prv
```

Cluster-based analysis

- Check the resulting scatter plot

```
laptop> gnuplot lulesh2.0-intel_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**

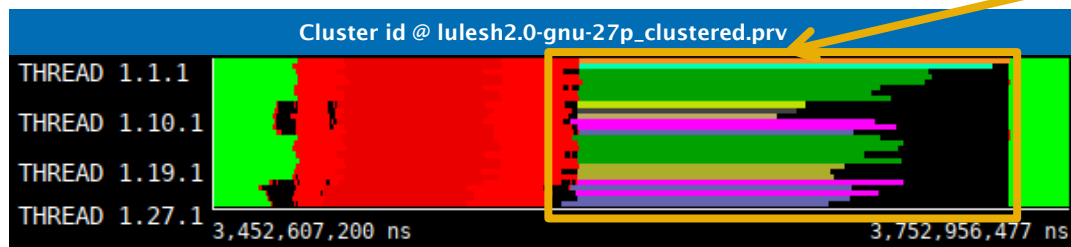
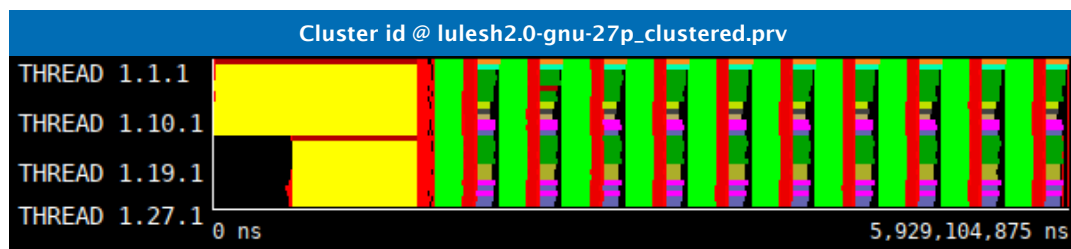


Correlating scatter plot and time distribution

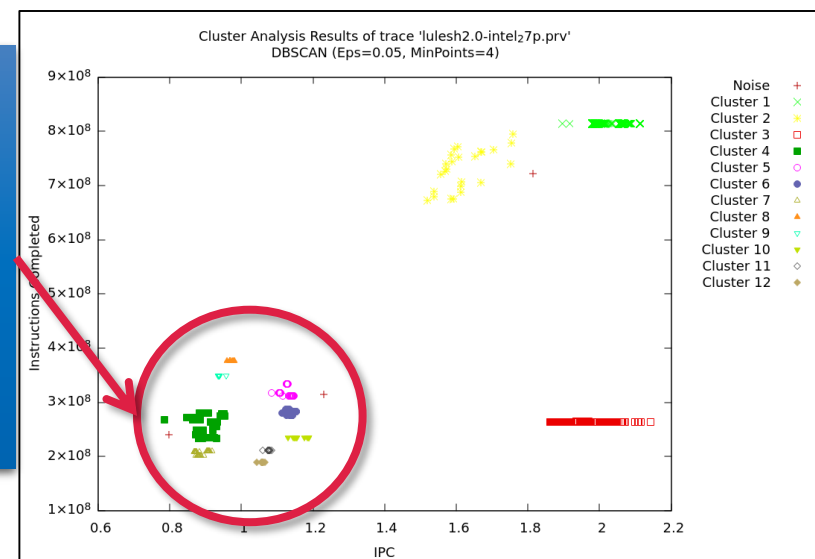
- Open the clustered trace with Paraver and look at it

```
laptop> $HOME/paraver/bin/wxparaver <path-to>/lulesh2.0-intel_27p_clustered.prv
```

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



Variable work /
speed
+
Simultaneously
@ different
processes
=
Imbalances



BSC Tools Hands-On

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