



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Extrae & Paraver

Judit Giménez, Lau Mercadal, Germán Llort

✉ tools@bsc.es

2024-09-05

Performance Analysis and Tools

Exrae features

- Platforms
 - Intel, ARM, RISC-V, POWER, Cray, BlueGene, Android, Fujitsu Sparc ...
- Parallel programming models
 - MPI, OpenMP, pthreads, OmpSs, CUDA, OpenACC, 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 anywhere in your program (Exrae API)

No need
to
recompile
nor relink!

Extrae overhead

Karolina	
Punctual event	175 ns
Event + PAPI counters	319 ns
Event + 1-level callstack	1.471 us
Event + 6-levels callstack	3.461 us

How does Extrae work?

- Symbol substitution through LD_PRELOAD

```
export LD_PRELOAD=$EXTRAE_HOME/lib/libmpitrace.so
```

- Specific libraries for each runtime and combinations
 - MPI
 - OpenMP
 - OpenMP+MPI
 - CUDA
 - ...
- Dynamic instrumentation
 - Based on Dyninst (developed by U.Wisconsin / U.Maryland)
 - Instrumentation in memory
 - Binary rewriting
- Static link (i.e., PMPI, Extrae API)

Recommended

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 Exrae

```
hpc$ vi job.slurm
```

```
#!/bin/bash

#SBATCH --job-name=lulesh2.0_s65_27p
#SBATCH --output=%x_%j.out
#SBATCH --error=%x_%j.err
#SBATCH --ntasks=27
#SBATCH --time=00:10:00

export TRACE_NAME=lulesh2.0_s65_27p.prv
srun ./trace.sh ./lulesh2.0 -i 10 -s 65 -p
```

Run with Exrae

Step 1: Adapt the job script to load Extrace

hpc\$ vi trace.sh

```
#!/bin/bash

#SBATCH --job-name=lulesh2.0_s65_27p
#SBATCH --output=%x_%j.out
#SBATCH --error=%x_%j.err
#SBATCH --ntasks=27
#SBATCH --time=00:10:00

export TRACE_NAME=lulesh2.0_s65_27p.prv
srun ./trace.sh ./lulesh2.0 -i 10 -s 65 -p
```

```
#!/bin/bash

# Configure Extrace
export EXTRAE_HOME=<installation_base_path>
export EXTRAE_CONFIG_FILE=./extrae.xml

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

# Run the program
$*
```

What to trace?

Type of application

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

```
hpc$ sbatch job.slurm
```

- Easy!

Step 2: Extrace XML configuration

```
hpc$ vi extrae.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?)

Step 2: Extrae XML configuration (II)

```
hpc$ vi extrae.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_L3_TCM, PAPI_BR_INS, PAPI_BR_MSP, RESOURCE_STALLS
    </set>
    <set enabled="yes" domain="all" changeat-time="500000us">
      PAPI_TOT_INS, PAPI_TOT_CYC, PAPI_VEC_SP, PAPI_SR_INS, PAPI_LD_INS
    </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?)

Step 2: Extrace XML configuration (III)

```
hpc$ vi extrae.xml
```

```
<buffer enabled="yes">
  <size enabled="yes">5000000</size>
  <circular enabled="no" />
</buffer>
```

Trace buffer size
(Flush/memory trade-off)

```
<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>
```

Enable sampling
(Want more details?)

Automatic
post-processing
to generate the
Paraver trace

All done! Check your resulting trace

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

```
hpc$ ls -l
...
lulesh2.0_s65_27p.pcf
lulesh2.0_s65_27p.prv
lulesh2.0_s65_27p.row
```

- Now let's look into it !



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Paraver Installation

Judit Giménez, Lau Mercadal, Germán Llort

✉ tools@bsc.es

2024-09-05

PATC: Performance Analysis and Tools

Install Paraver

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

Pick your version

CORE TOOLS

wxparaver-4.11.4-win.zip

wxparaver-4.11.4-mac.zip

wxparaver-4.11.4-Linux_x86_64.tar.gz

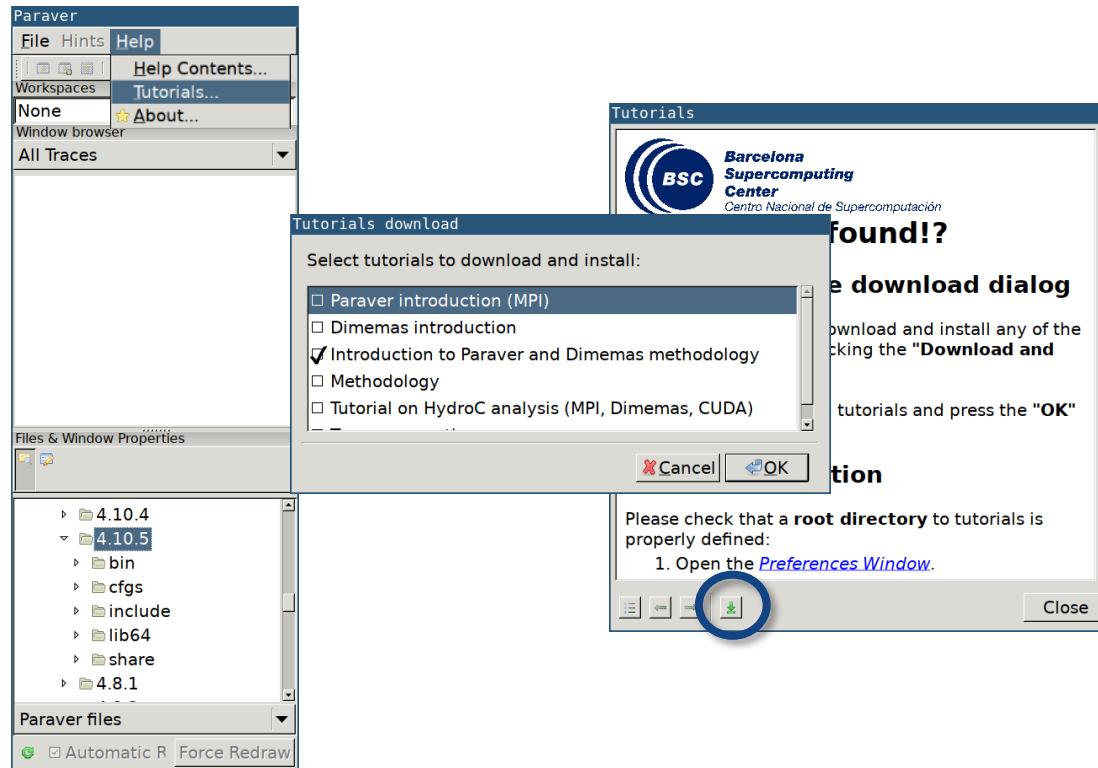
laptop\$ tar xf wxparaver-4.11.4-linux-x86_64.tar.gz

laptop\$ mv wxparaver-4.11.4-linux-x86_64 paraver

Install Tutorials (Automatic)

- Start Paraver

```
laptop$ paraver/bin/wxparaver
```



Install Tutorials (Manual)

- Download tutorials:
 - Documentation → Paraver Tutorials

The screenshot shows a web browser displaying the BSC (Barcelona Supercomputing Center) website. The URL in the address bar is `news@tools:~ > Paraver 4.7.2 available`. The page title is "Paraver tutorials". The navigation menu includes Home, Paraver, Dimemas, Extrae, Research, Documentation (which is currently selected), Downloads, and Publications. Below the menu, a breadcrumb trail shows "Home > Documentation > Paraver tutorials". The main content area contains a paragraph about seven tutorials available for Paraver 4.7.2. A list of tutorials is provided, with the third item, "Introduction to Paraver and Dimemas methodology", circled in red. A large blue callout box with the text "Download links" is positioned over this circled item. The list of tutorials includes:

- Paraver introduction (MPI)
- Dimemas introduction
- Introduction to Paraver and Dimemas methodology
- Methodology
- Tutorial on HydroC analysis (MPI, Dimemas, CUDA)
- Trace preparation
- Trace alignment tutorial.

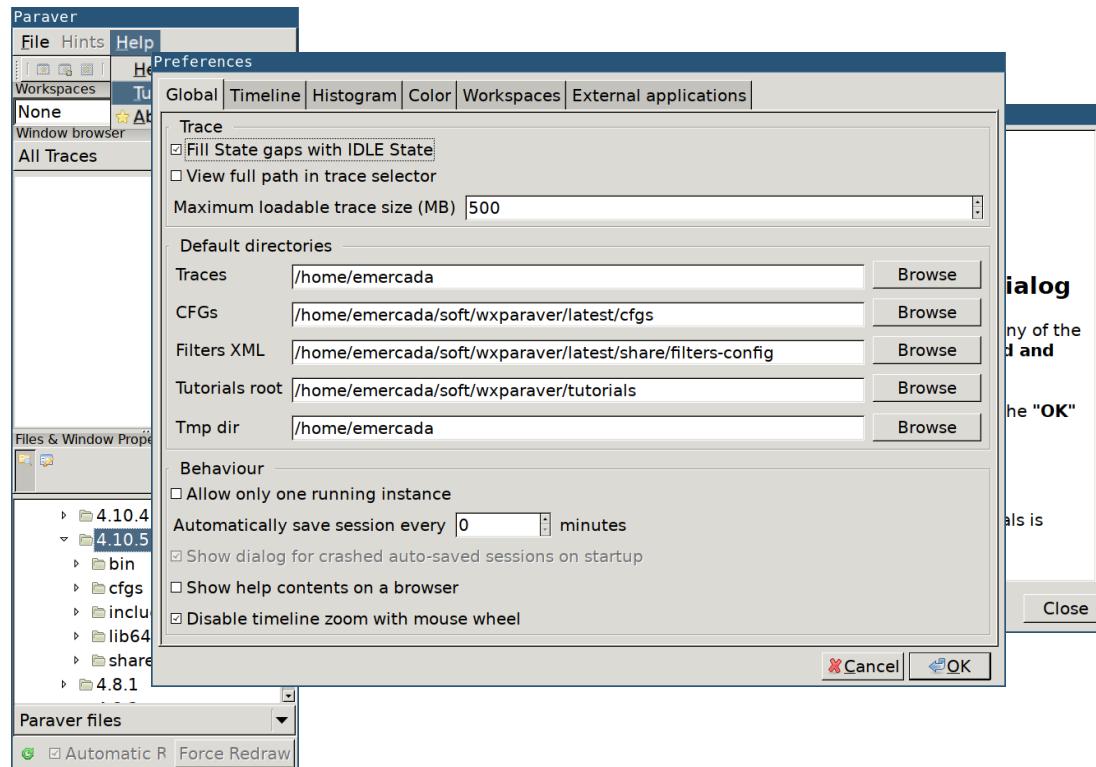
```
laptop$ mkdir paraver/tutorials
```

```
laptop$ tar xf 3.introduction* -C paraver/tutorials/
```

Install Tutorials (Manual)

- Start Paraver

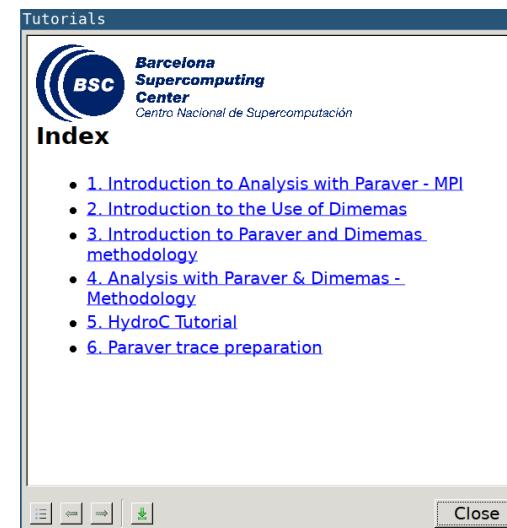
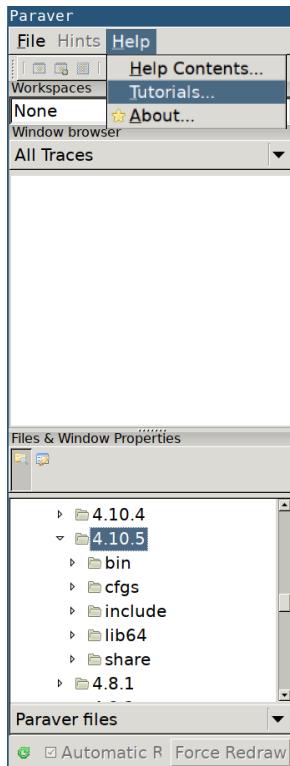
```
laptop$ paraver/bin/wxparaver
```



Install Paraver

- Start Paraver

```
laptop$ paraver/bin/wxparaver
```





**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Paraver Hands-On

Judit Giménez, Lau Mercadal, Germán Llort

✉ tools@bsc.es

2024-09-05

PATC: Performance Analysis and Tools

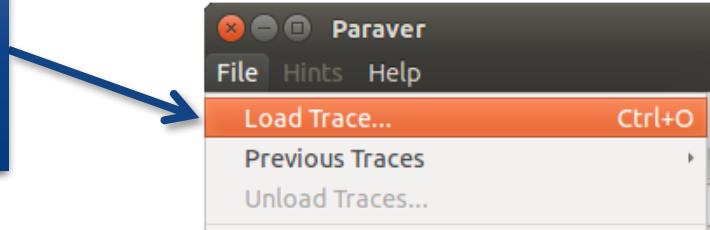
First steps of analysis

- Copy the trace to your laptop

```
laptop$ scp <USER>@hpc:lulesh2.0_s65_27p.* ./
```

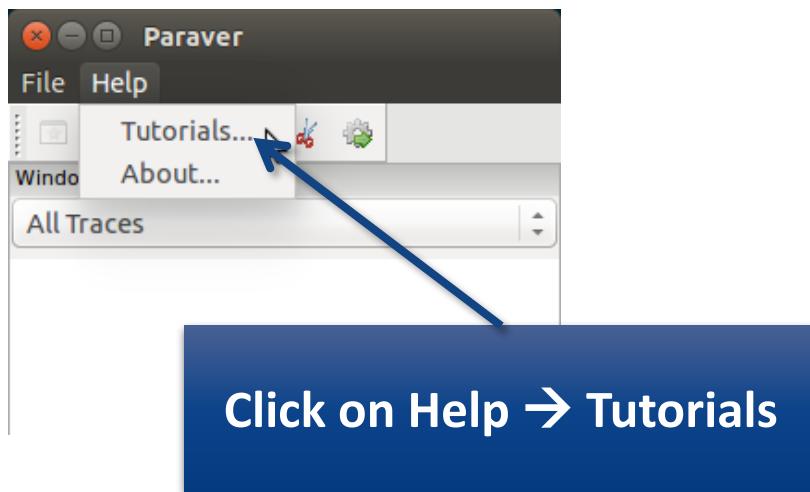
- Load the trace with Paraver

Click on File → Load Trace
→ Browse to ...
“lulesh2.0_s65_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”

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 every thread spends in every MPI call. Look at the global statistics at the outside MPI column. Entry Average represents the application efficiency, entry Avg/Max represents the global load balance and entry Avg/Min represents the communication efficiency. If any of those values are below 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 populates 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. The histogram does not show vertical lines; it indicates the computation is not balanced. Open the control window to look at the time distribution and correlate both views.

Parallel efficiency

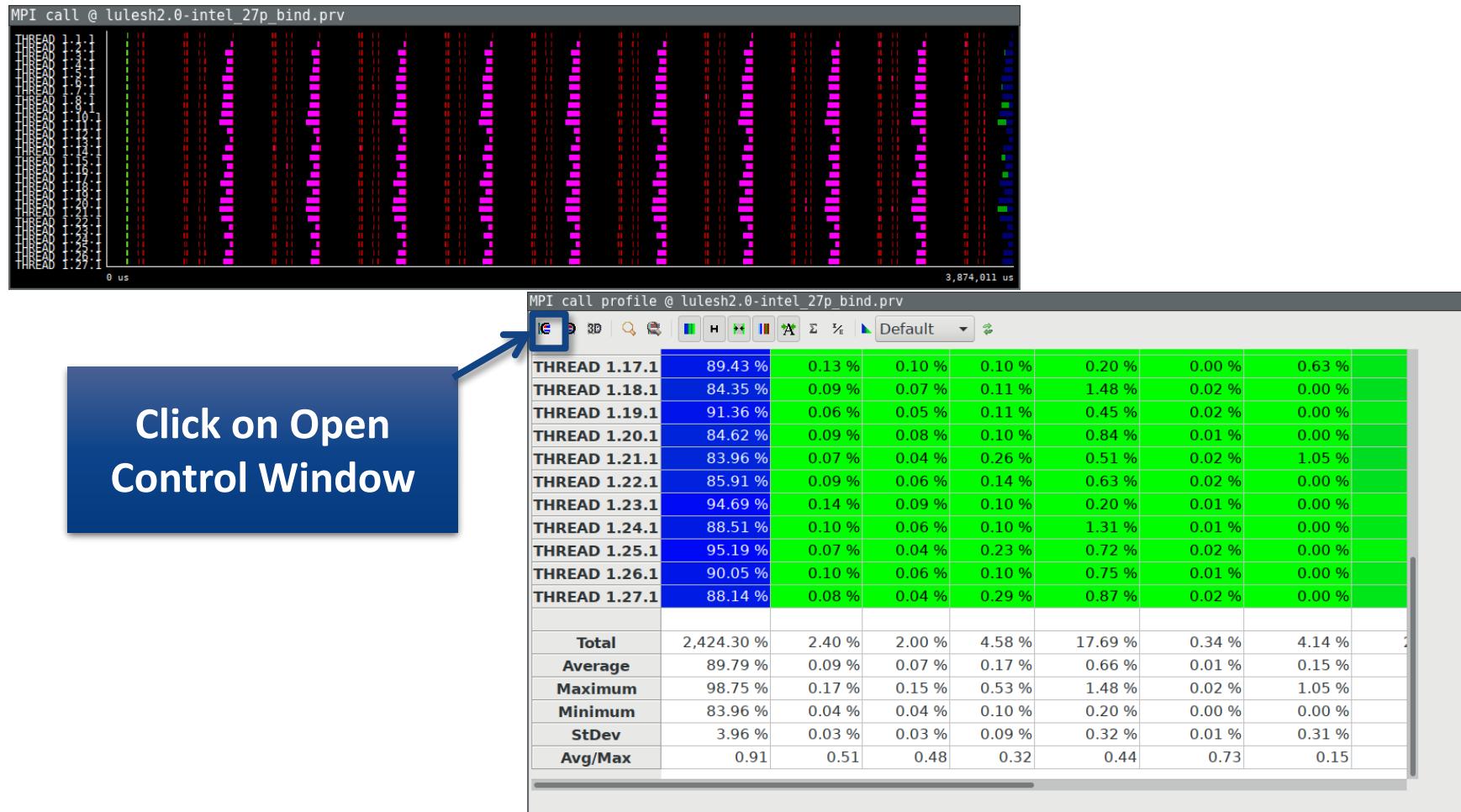
Comm efficiency

Load balance

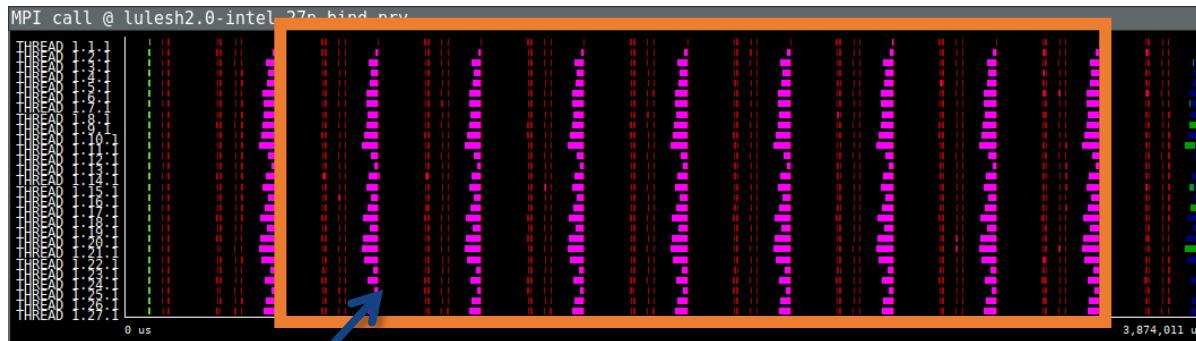
MPI call profile @ lulesh2.0-intel_27p_bind.prv

	Total	Avg	Max	Min	StDev	Avg/Max	Avg/Min
THREAD 1.17.1	89.43 %	0.13 %	0.10 %	0.10 %	0.20 %	0.00 %	0.63 %
THREAD 1.18.1	84.35 %	0.09 %	0.07 %	0.11 %	1.48 %	0.02 %	0.00 %
THREAD 1.19.1	91.36 %	0.06 %	0.05 %	0.11 %	0.45 %	0.02 %	0.00 %
THREAD 1.20.1	84.62 %	0.09 %	0.08 %	0.10 %	0.84 %	0.01 %	0.00 %
THREAD 1.21.1	83.96 %	0.07 %	0.04 %	0.26 %	0.51 %	0.02 %	1.05 %
THREAD 1.22.1	85.91 %	0.09 %	0.06 %	0.14 %	0.63 %	0.02 %	0.00 %
THREAD 1.23.1	94.69 %	0.14 %	0.09 %	0.10 %	0.20 %	0.01 %	0.00 %
THREAD 1.24.1	88.51 %	0.10 %	0.06 %	0.10 %	1.31 %	0.01 %	0.00 %
THREAD 1.25.1	95.19 %	0.07 %	0.04 %	0.23 %	0.72 %	0.02 %	0.00 %
THREAD 1.26.1	90.05 %	0.10 %	0.06 %	0.10 %	0.75 %	0.01 %	0.00 %
THREAD 1.27.1	88.14 %	0.08 %	0.04 %	0.29 %	0.87 %	0.02 %	0.00 %
Total	2,424.30 %	2.40 %	2.00 %	4.58 %	17.69 %	0.34 %	4.14 %
Avg	89.79 %	0.09 %	0.07 %	0.17 %	0.66 %	0.01 %	0.15 %
Max	98.75 %	0.17 %	0.15 %	0.53 %	1.48 %	0.02 %	1.05 %
Min	83.90 %	0.04 %	0.04 %	0.10 %	0.20 %	0.00 %	0.00 %
StDev	3.96 %	0.03 %	0.03 %	0.09 %	0.32 %	0.01 %	0.31 %
Avg/Max	0.91	0.51	0.48	0.32	0.44	0.73	0.15

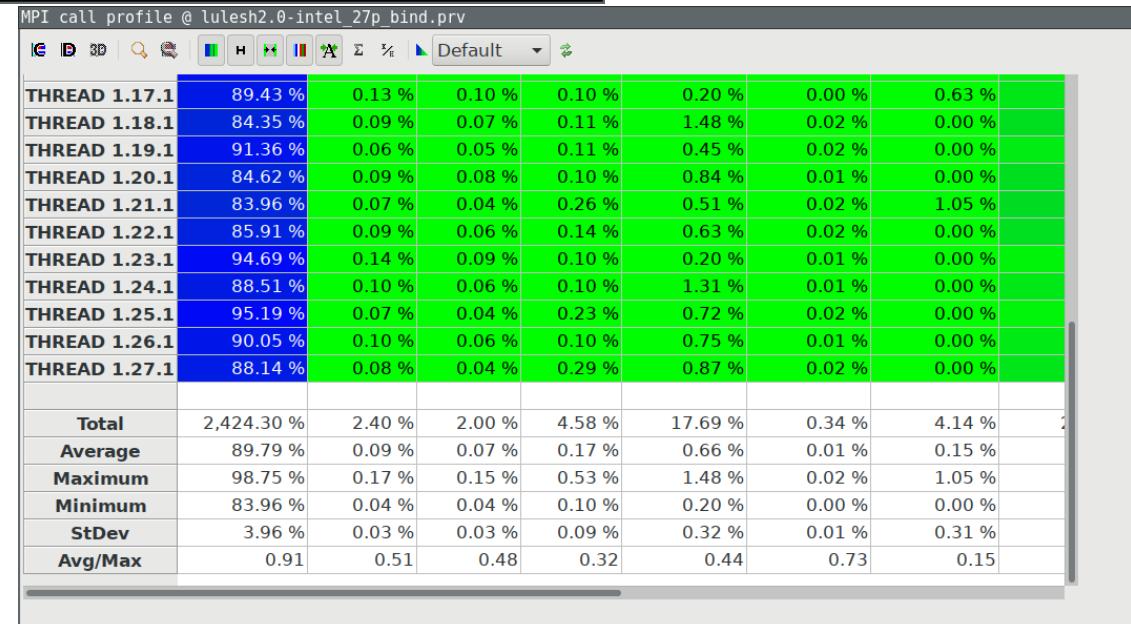
Focus on the iterative part



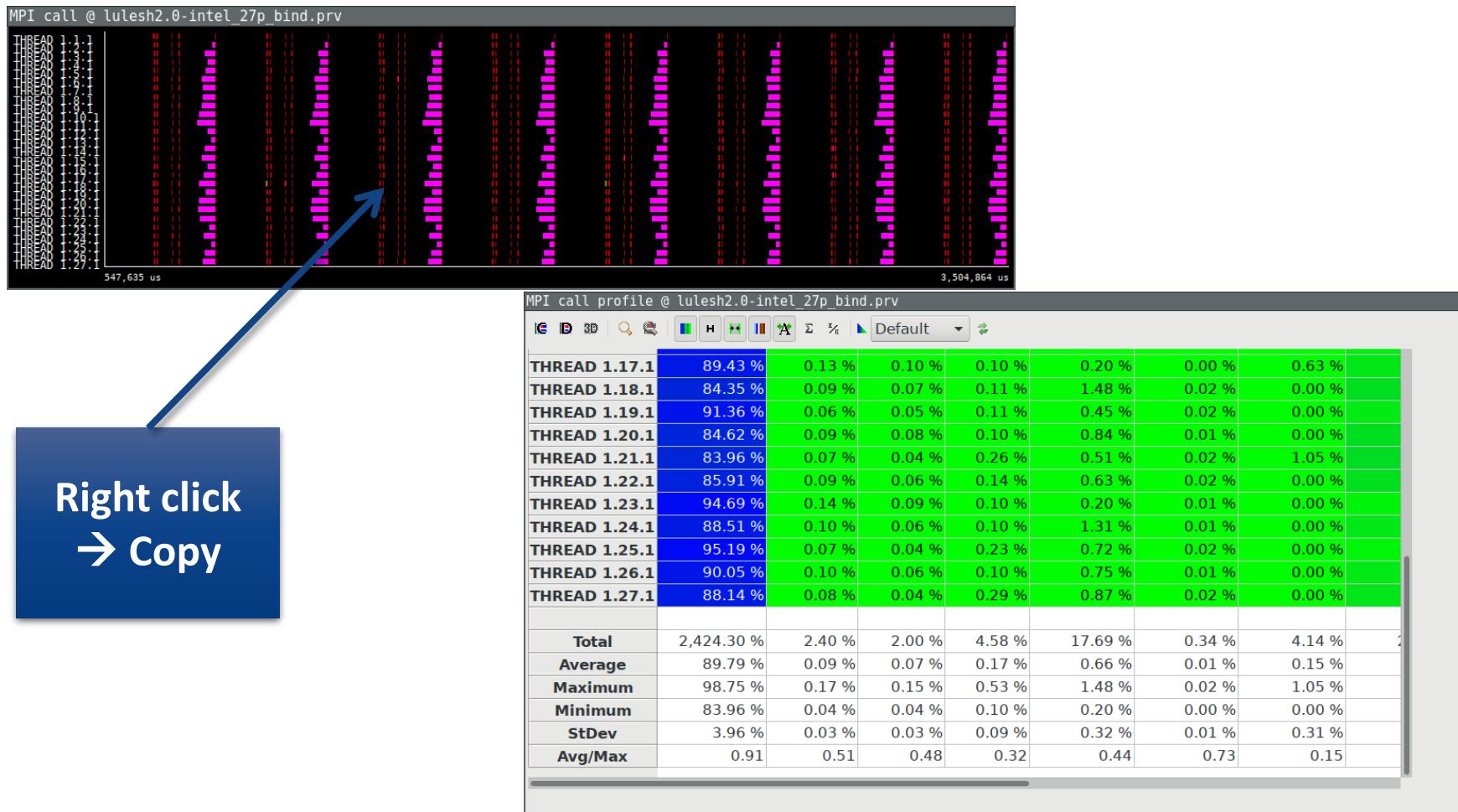
Focus on the iterative part



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



Recalculate efficiency of iterative region



Recalculate efficiency of iterative region



MPI call profile @ lulesh2.0-intel_27p_bind.prv

	89.03 %	0.12 %	0.09 %	0.10 %	0.20 %	10.41 %	0.05 %
THREAD 1.17.1	89.03 %	0.12 %	0.09 %	0.10 %	0.20 %	10.41 %	0.05 %
THREAD 1.18.1	83.71 %	0.09 %	0.06 %	0.06 %	1.58 %	14.44 %	0.05 %
THREAD 1.19.1	91.09 %	0.06 %	0.04 %	0.09 %	0.46 %	8.21 %	0.05 %
THREAD 1.20.1	84.05 %	0.09 %	0.08 %	0.06 %	0.88 %	14.79 %	0.05 %
THREAD 1.21.1	83.33 %	0.07 %	0.04 %	0.25 %	0.52 %	15.73 %	0.05 %
THREAD 1.22.1	85.40 %	0.09 %	0.06 %	0.10 %	0.65 %	13.66 %	0.05 %
THREAD 1.23.1	94.58 %	0.14 %	0.09 %	0.09 %	0.19 %	4.87 %	0.05 %
THREAD 1.24.1	88.12 %	0.10 %	0.06 %	0.06 %	1.41 %	10.20 %	0.05 %
THREAD 1.25.1	95.10 %	0.07 %	0.04 %	0.23 %	0.76 %	3.76 %	0.05 %
THREAD 1.26.1	99.73 %	0.10 %	0.05 %	0.05 %	0.80 %	9.23 %	0.05 %
THREAD 1.27.1	87.66 %	0.17 %	0.04 %	0.31 %	0.92 %	10.95 %	0.06 %
Total	2,414.73 %	2.29 %	1.96 %	4.09 %	18.55 %	257.03 %	1.35 %
Average	89.43 %	0.08 %	0.07 %	0.15 %	0.69 %	9.52 %	0.05 %
Maximum	98.76 %	0.17 %	0.15 %	0.55 %	1.58 %	16.08 %	0.06 %
Minimum	83.33 %	0.04 %	0.04 %	0.05 %	0.19 %	0.01 %	0.05 %
StDev	4.15 %	0.03 %	0.03 %	0.11 %	0.35 %	4.09 %	0.00 %
Avg/Max	0.91	0.51	0.48	0.28	0.44	0.59	0.87

Right click
→ Paste →
Time

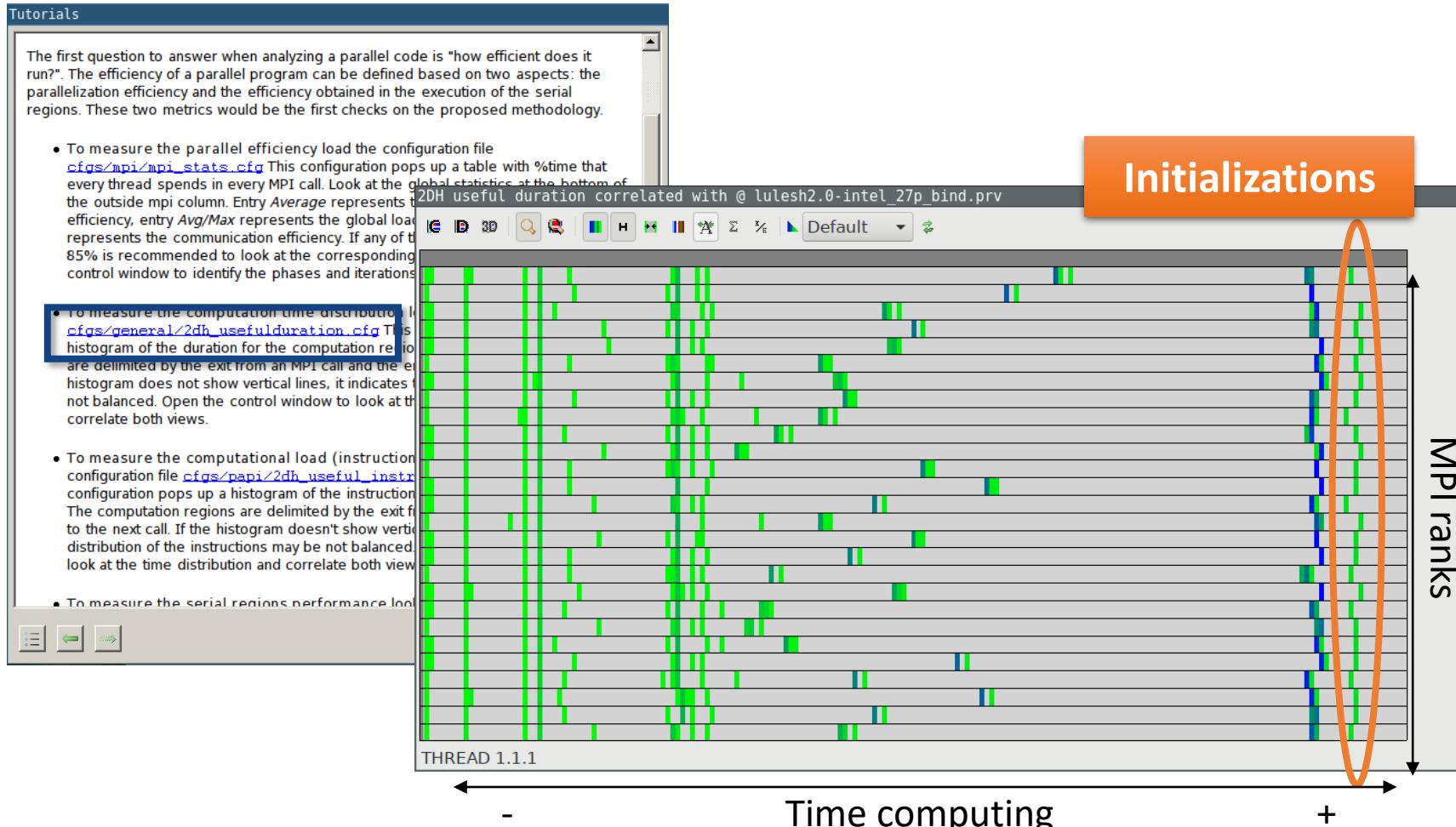
Efficiency of iterative region



MPI call profile @ lulesh2.0-intel_27p_bind.prv							
	Default						
THREAD 1.17.1	89.03 %	0.12 %	0.09 %	0.10 %	0.20 %	10.41 %	0.05 %
THREAD 1.18.1	83.71 %	0.09 %	0.06 %	0.06 %	1.58 %	14.44 %	0.05 %
THREAD 1.19.1	91.09 %	0.06 %	0.04 %	0.09 %	0.46 %	8.21 %	0.05 %
THREAD 1.20.1	84.05 %	0.09 %	0.08 %	0.06 %	0.88 %	14.79 %	0.05 %
THREAD 1.21.1	83.33 %	0.07 %	0.04 %	0.25 %	0.52 %	15.73 %	0.05 %
THREAD 1.22.1	85.40 %	0.09 %	0.06 %	0.10 %	0.65 %	13.66 %	0.05 %
THREAD 1.23.1	94.58 %	0.14 %	0.09 %	0.09 %	0.19 %	4.87 %	0.05 %
THREAD 1.24.1	88.12 %	0.10 %	0.06 %	0.06 %	1.41 %	10.20 %	0.05 %
THREAD 1.25.1	95.10 %	0.07 %	0.04 %	0.23 %	0.76 %	3.76 %	0.05 %
THREAD 1.26.1	89.73 %	0.10 %	0.05 %	0.05 %	0.80 %	9.23 %	0.05 %
THREAD 1.27.1	87.66 %	0.07 %	0.04 %	0.31 %	0.92 %	10.95 %	0.06 %
Total	2,414.73 %	2.29 %	1.96 %	4.09 %	18.55 %	257.03 %	1.35 %
Average	89.43 %	0.08 %	0.07 %	0.15 %	0.69 %	9.52 %	0.05 %
Maximum	98.76 %	0.17 %	0.15 %	0.55 %	1.58 %	16.08 %	0.06 %
Minimum	83.33 %	0.04 %	0.04 %	0.05 %	0.19 %	0.01 %	0.05 %
StDev	4.15 %	0.03 %	0.03 %	0.11 %	0.35 %	4.09 %	0.00 %
Avg/Rt	0.91	0.51	0.48	0.28	0.44	0.59	0.87

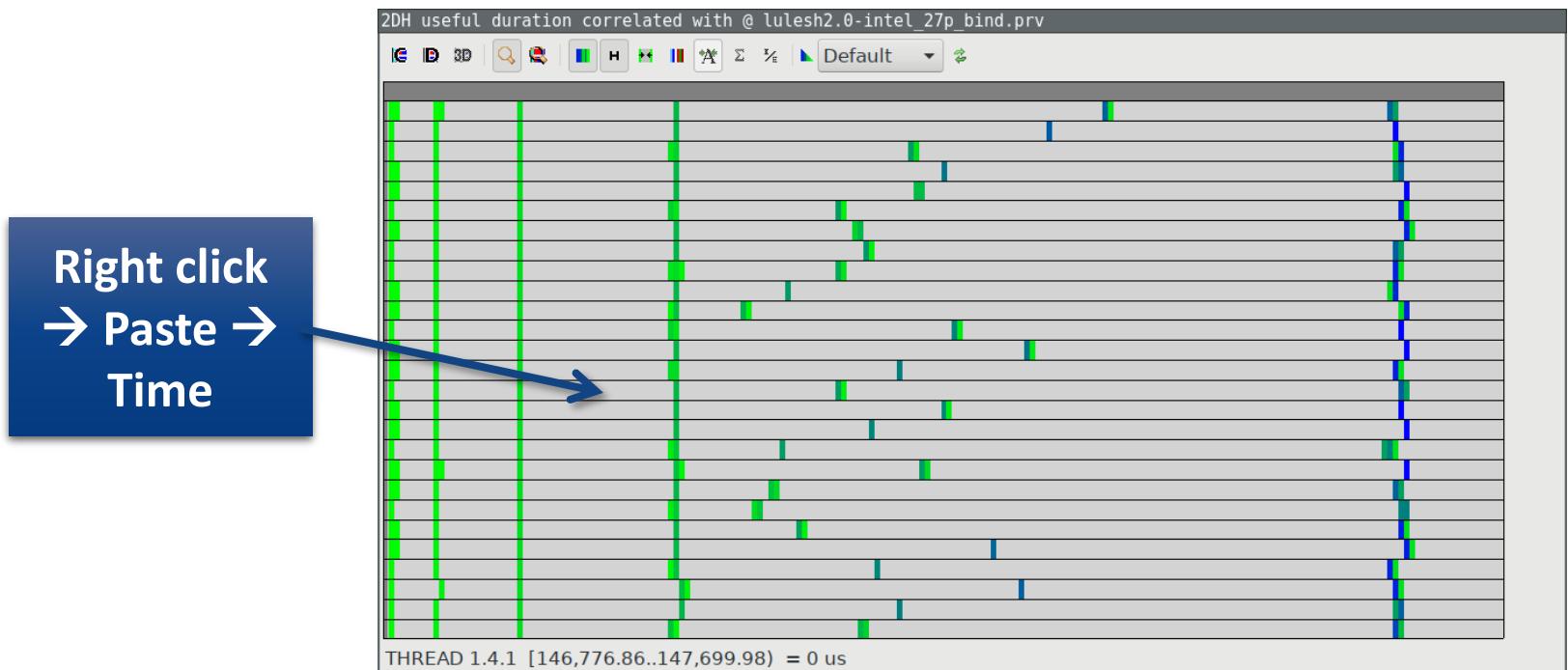
Computation time distribution

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



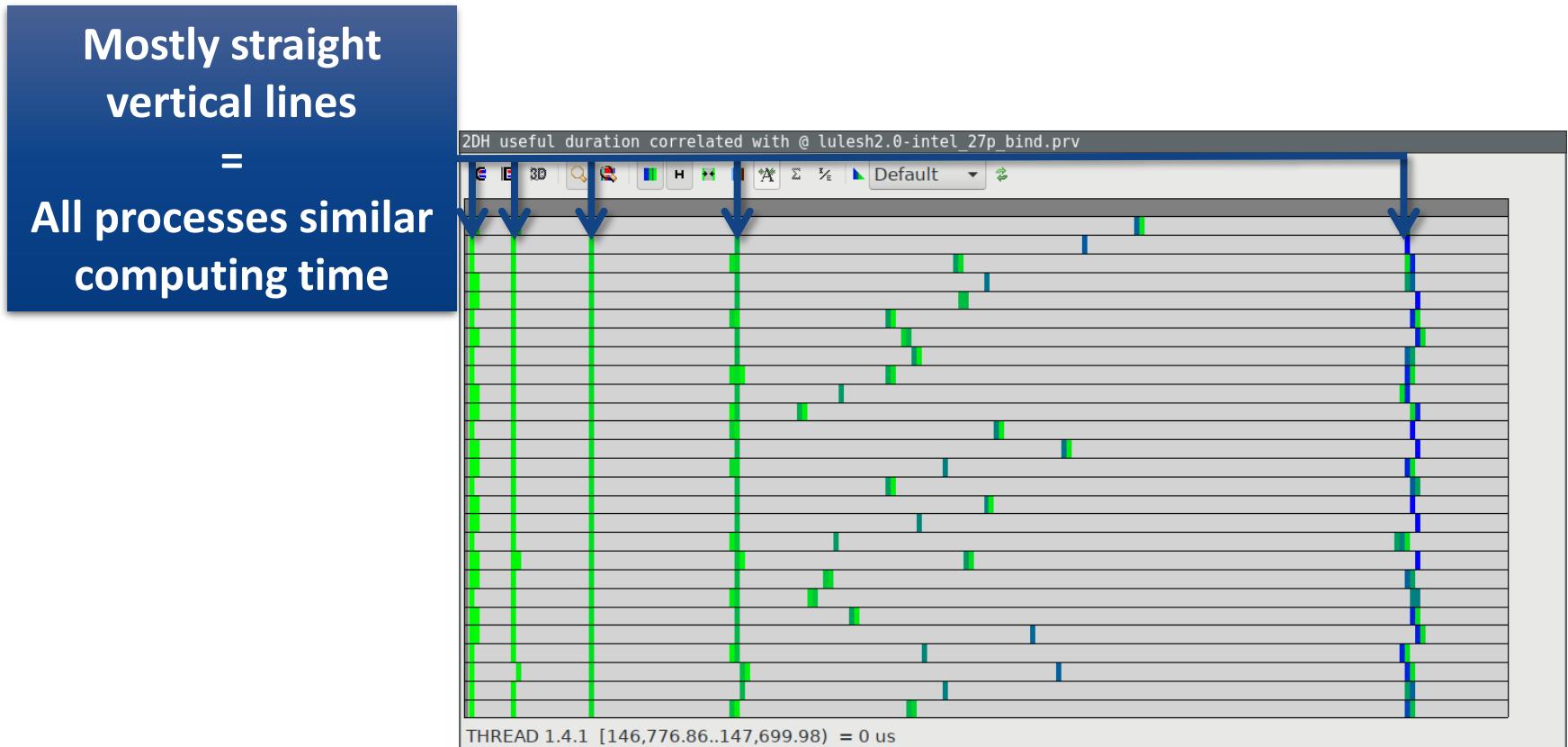
Focus on the iterative part

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



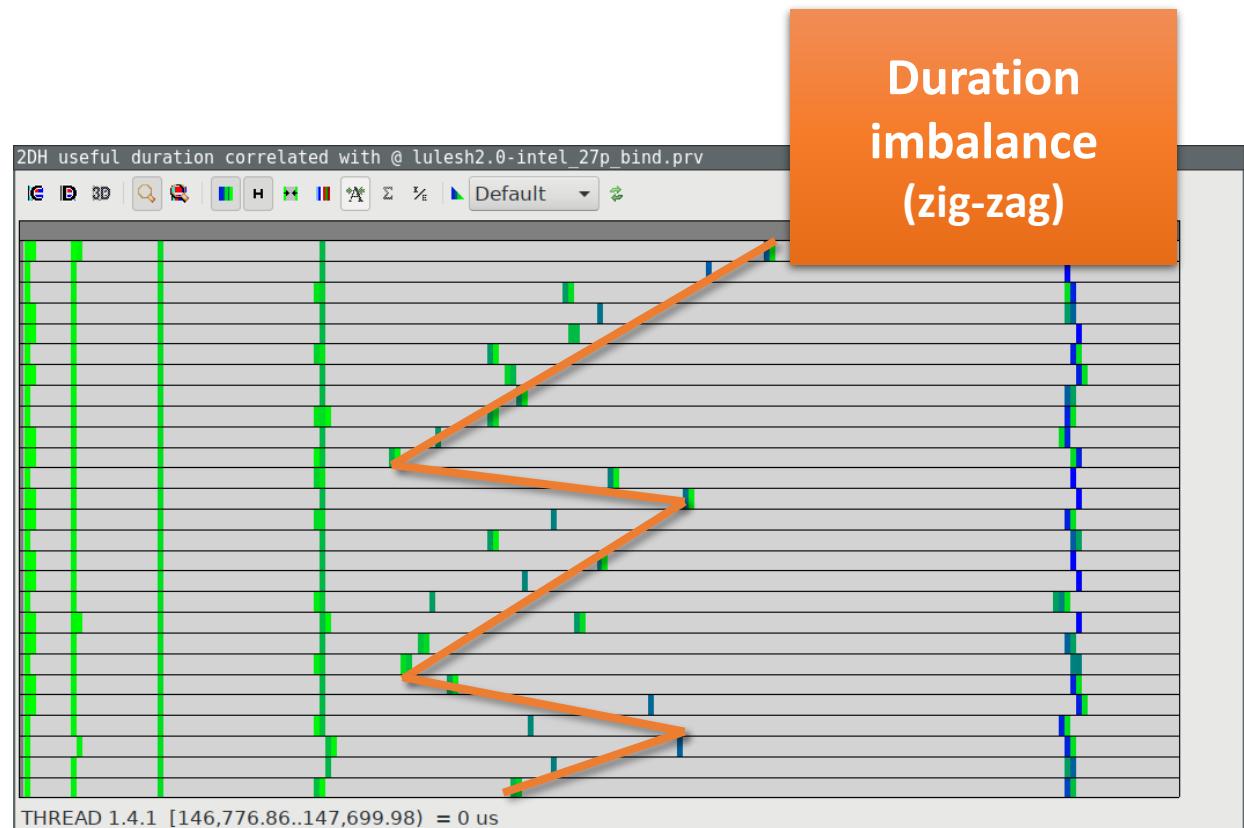
Computation time distribution

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



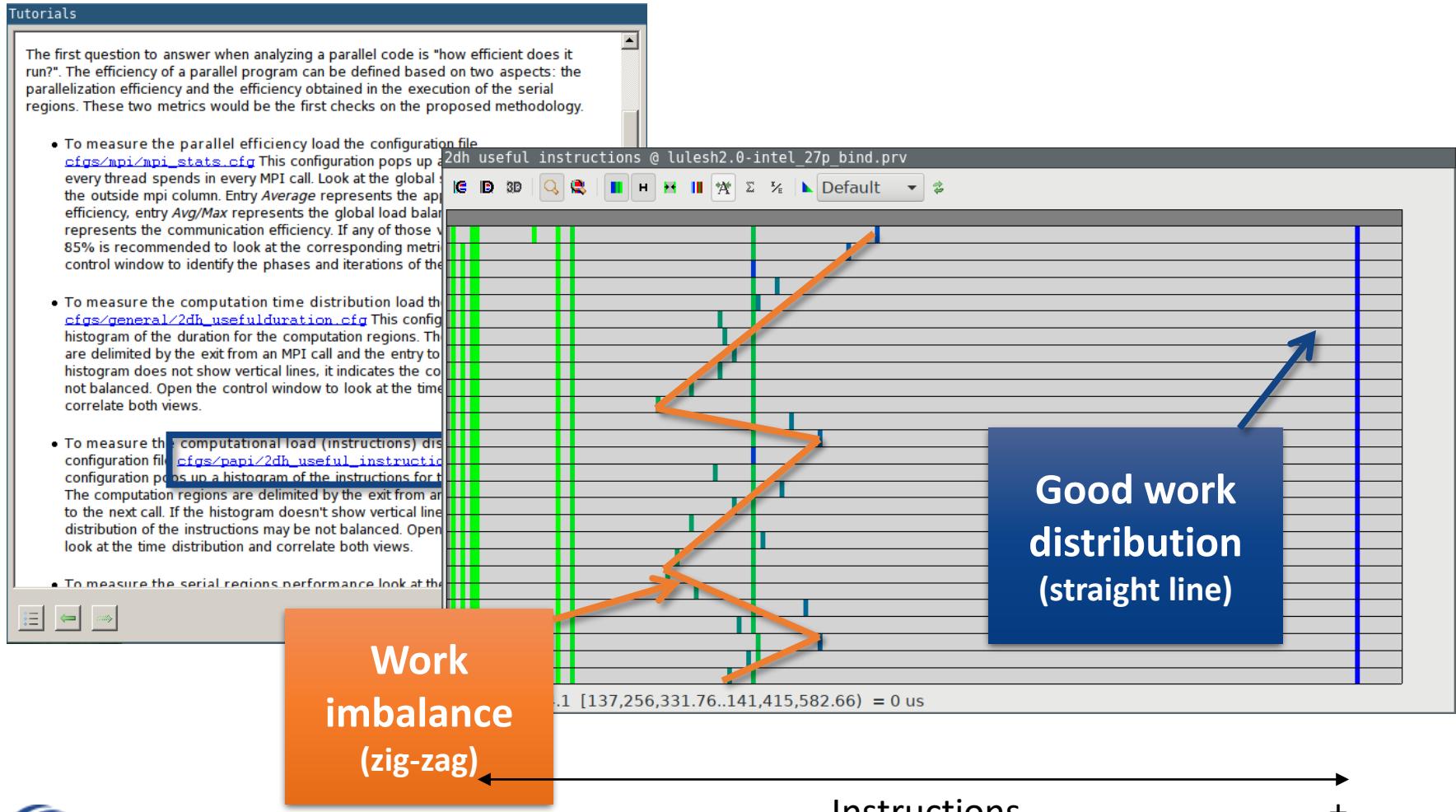
Computation time distribution

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



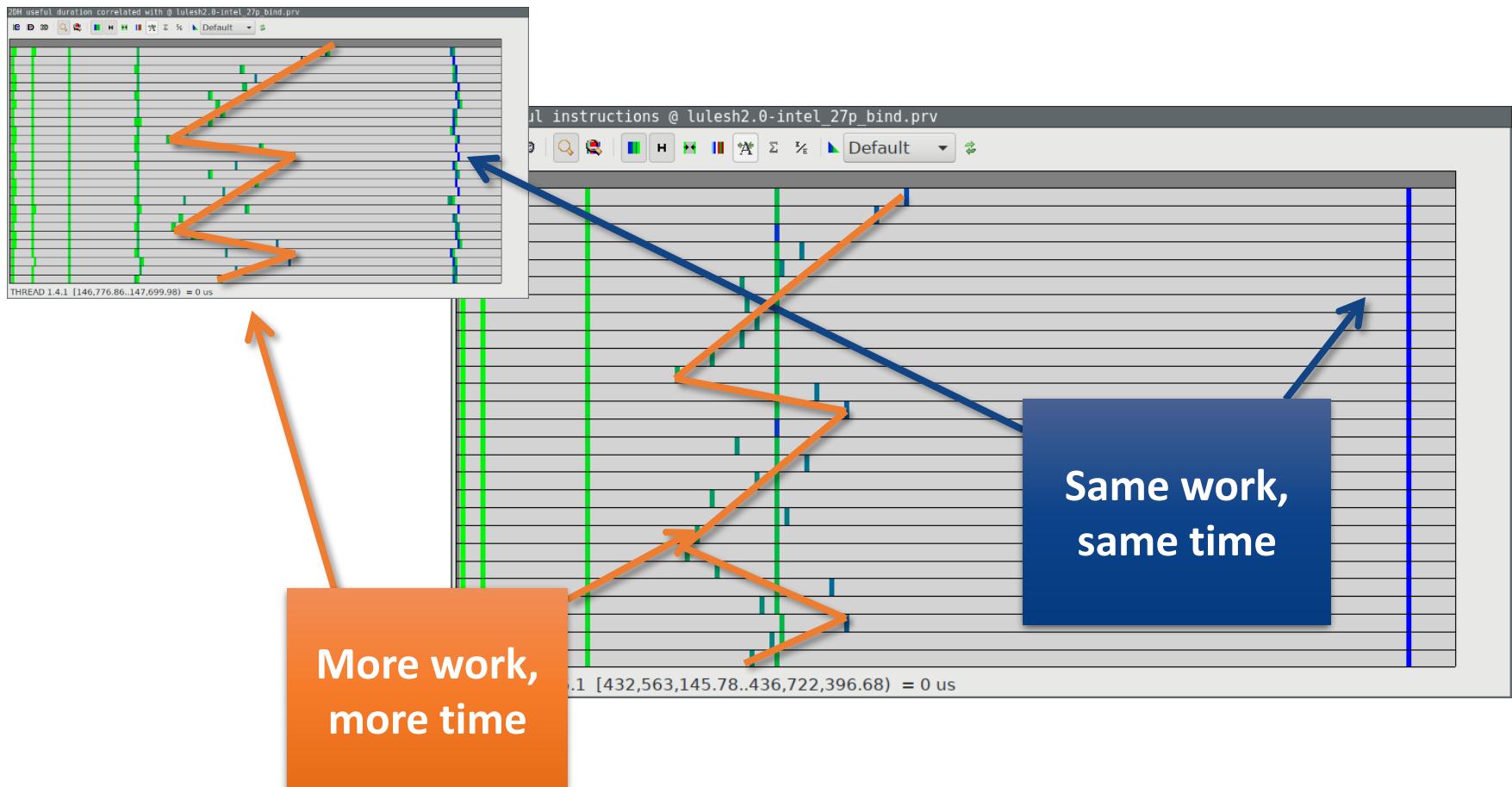
Computation load distribution

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



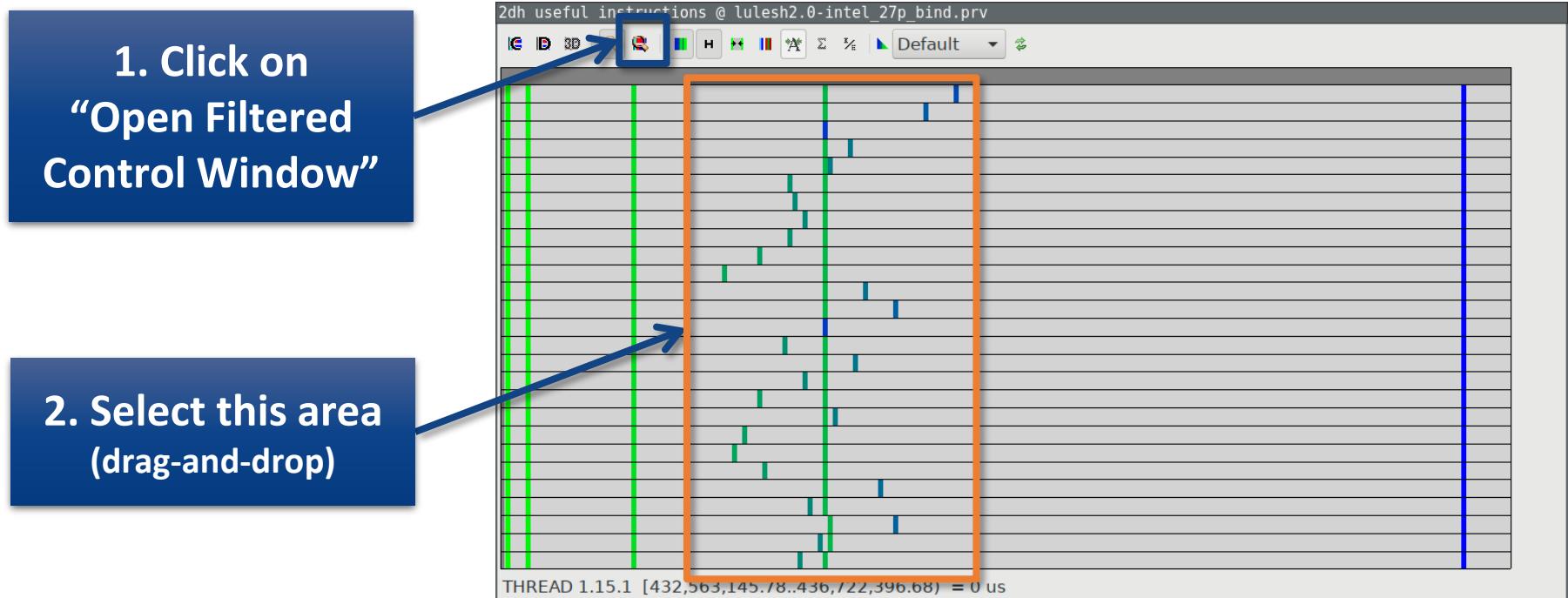
Correlate two histograms

- Clear correlation between the **amount of work** and the **time computing**



Where does this happen?

- Go from the table to the timeline

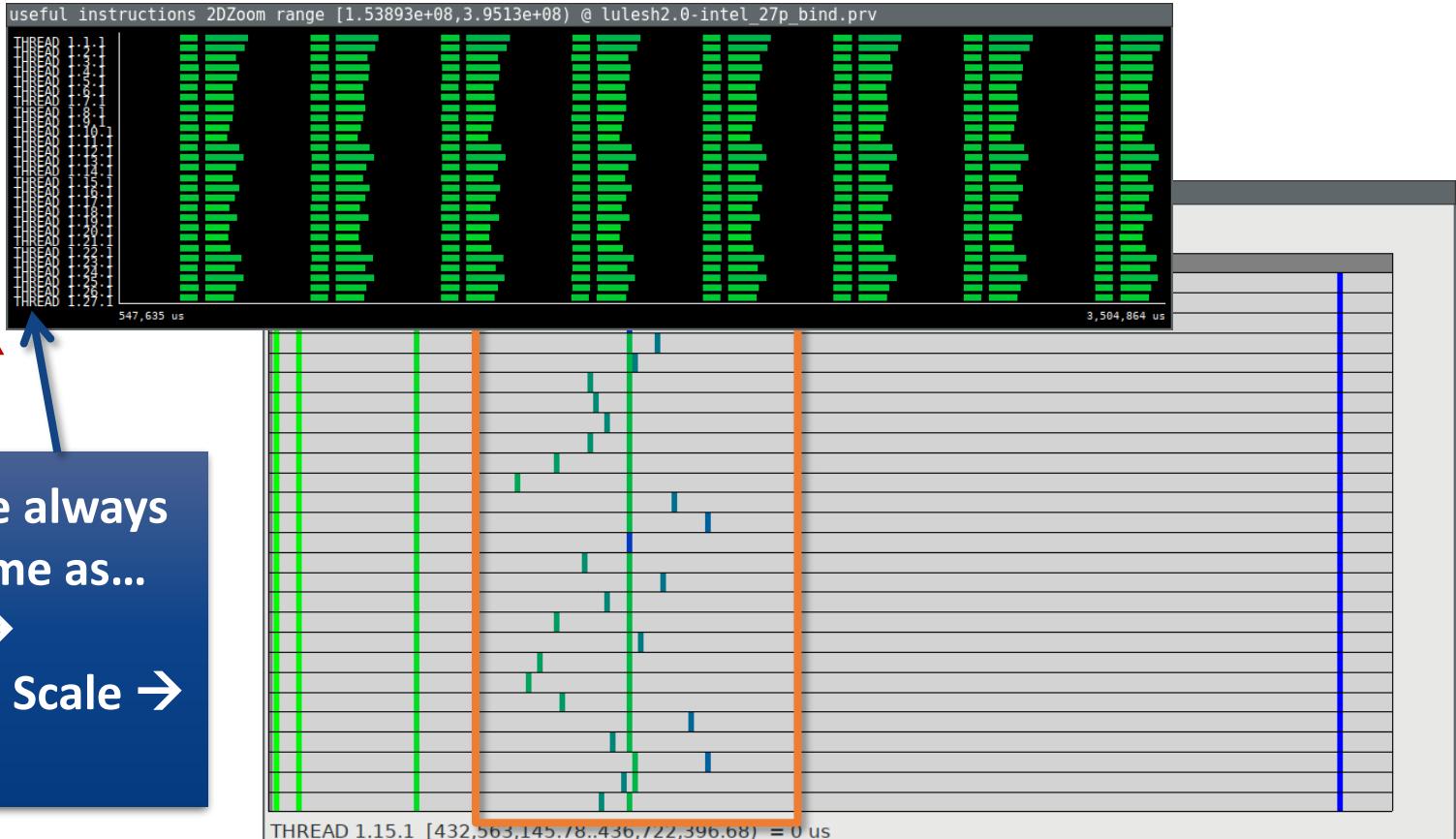


Where does this happen?

- Go from the table to the timeline



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



Where does this happen?

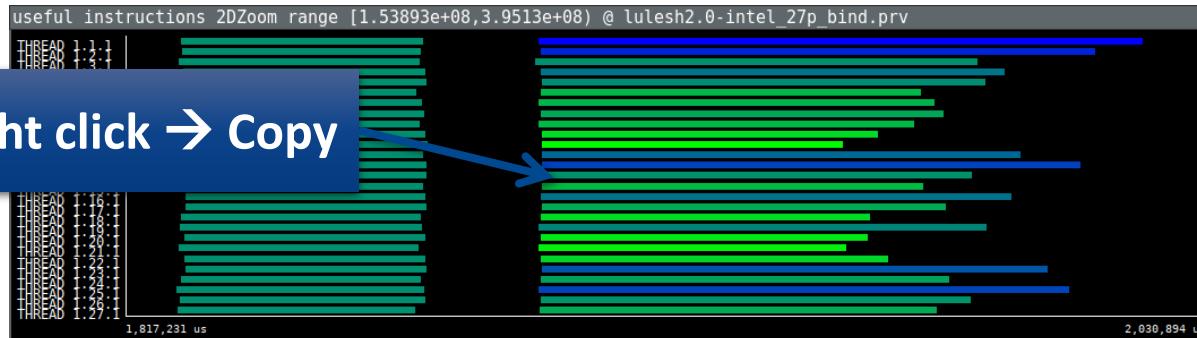
- Slow & Fast at the same time? → Imbalance



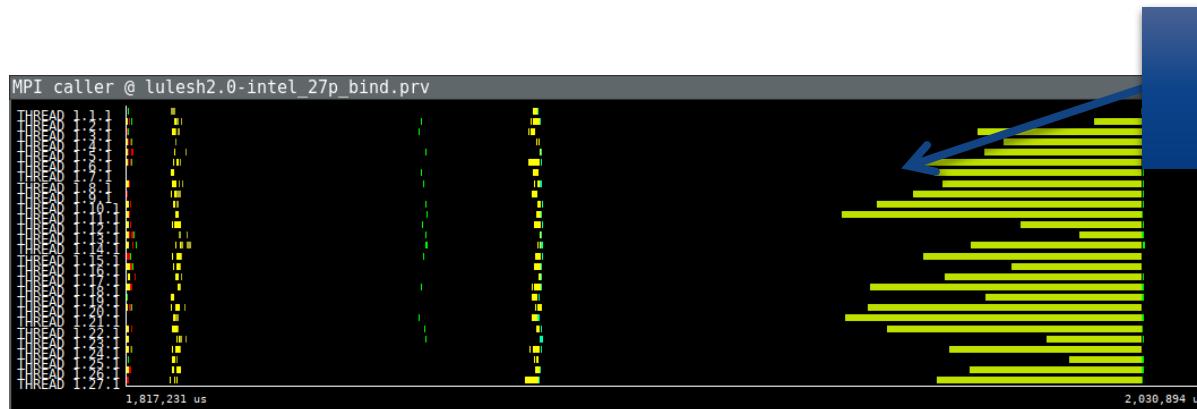
Zoom into
1 of the iterations
(by drag-and-dropping)

Where does this happen?

- Slow & Fast at the same time? → Imbalance



- Hints → Call stack references → Caller function



Where does this happen?

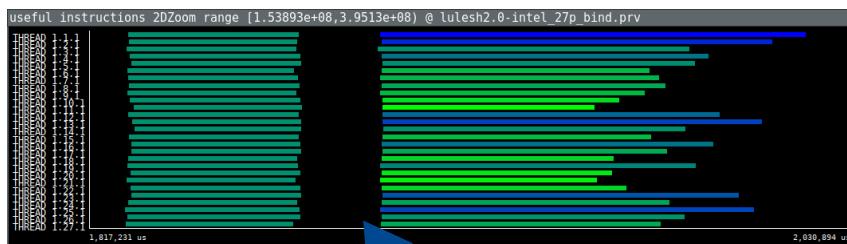
- Slow & Fast at the same time? → Imbalance



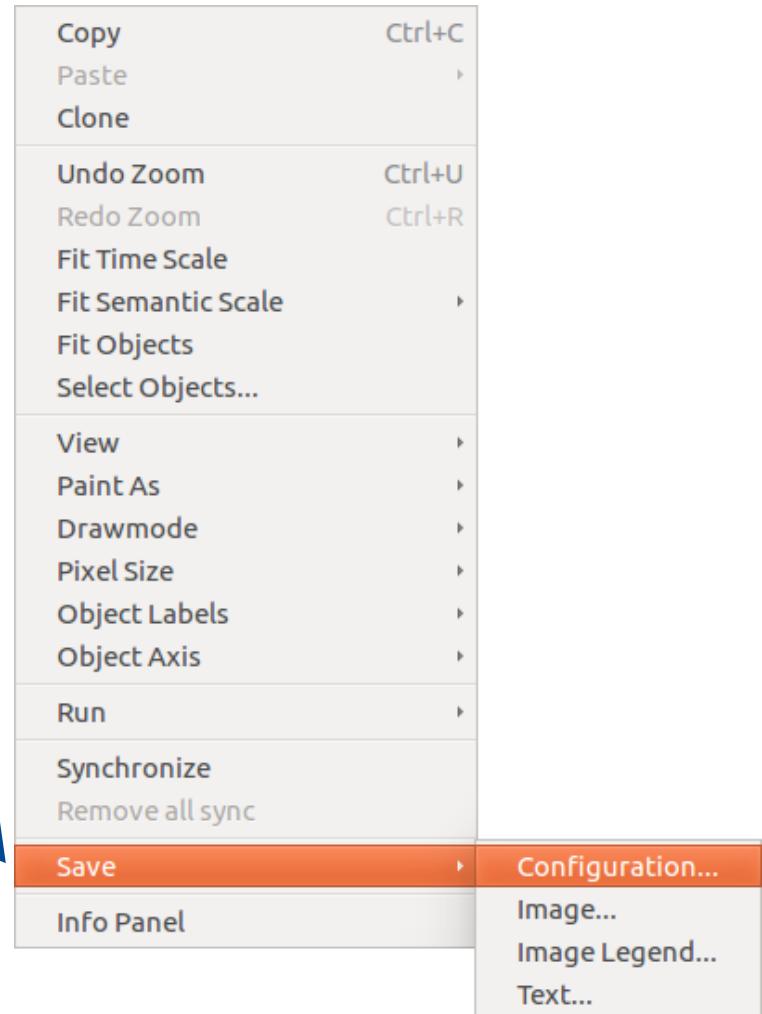
- Hints → Call stack references → Caller function



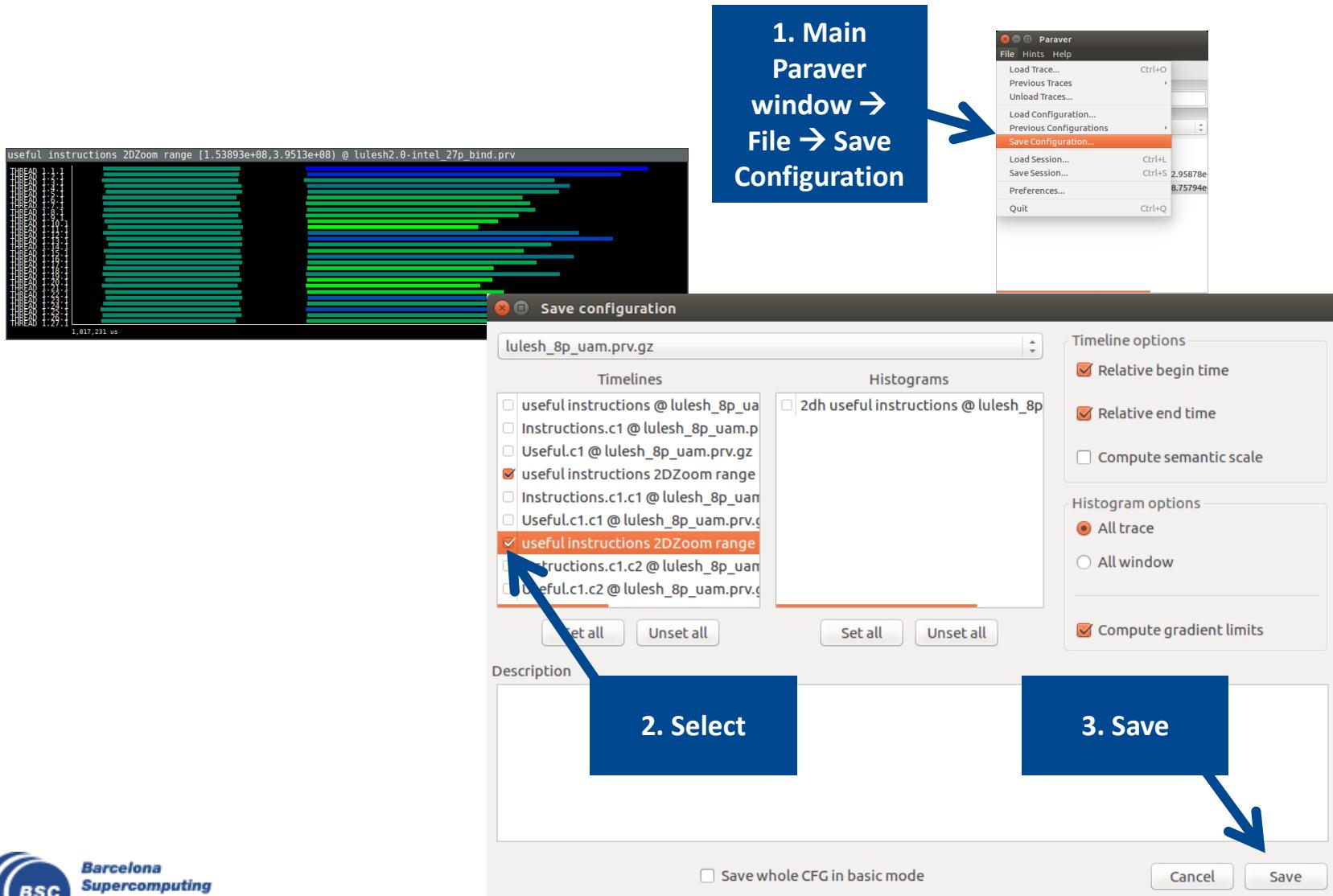
Save CFG's (method 1)



Right click on
timeline

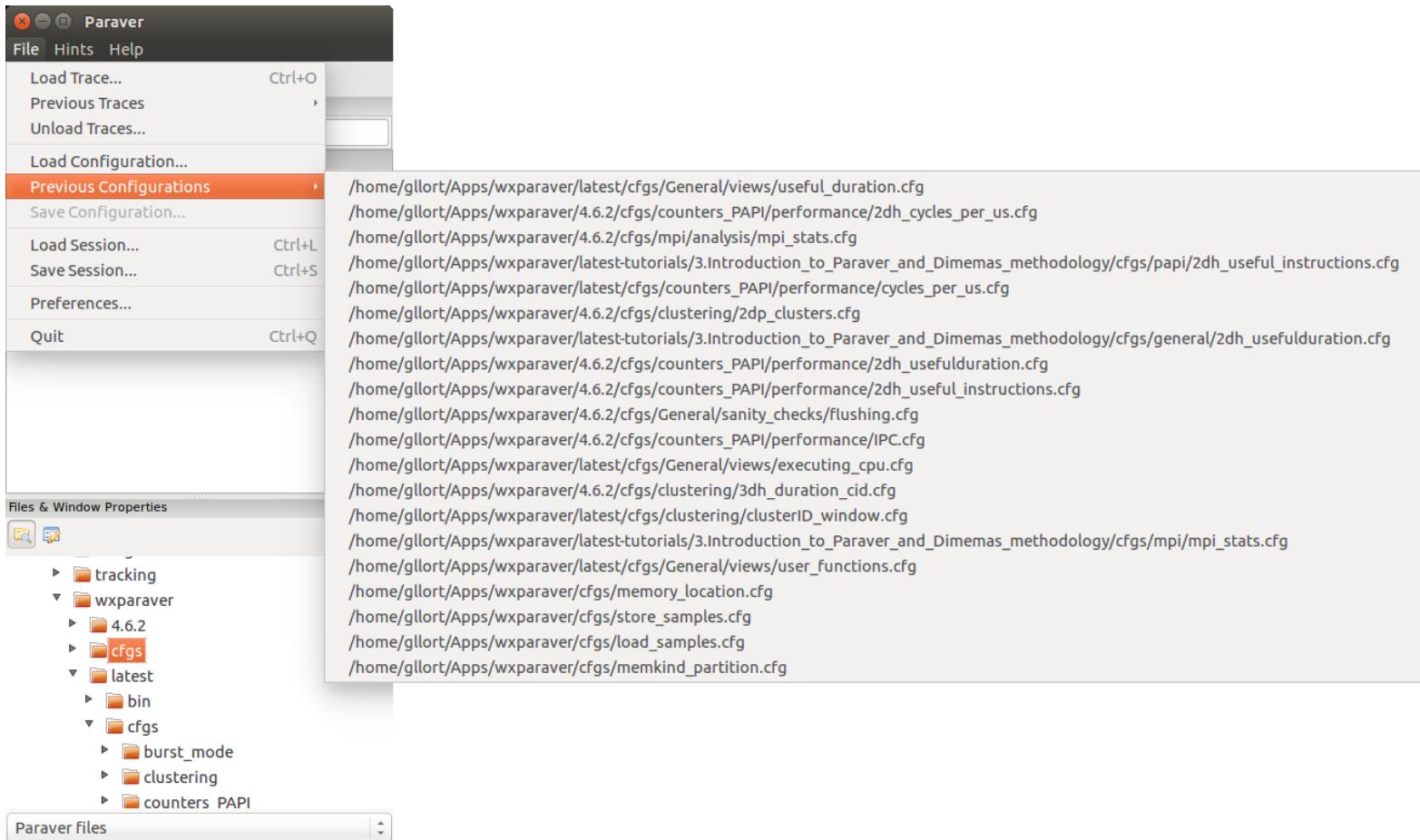


Save CFG's (method 2)



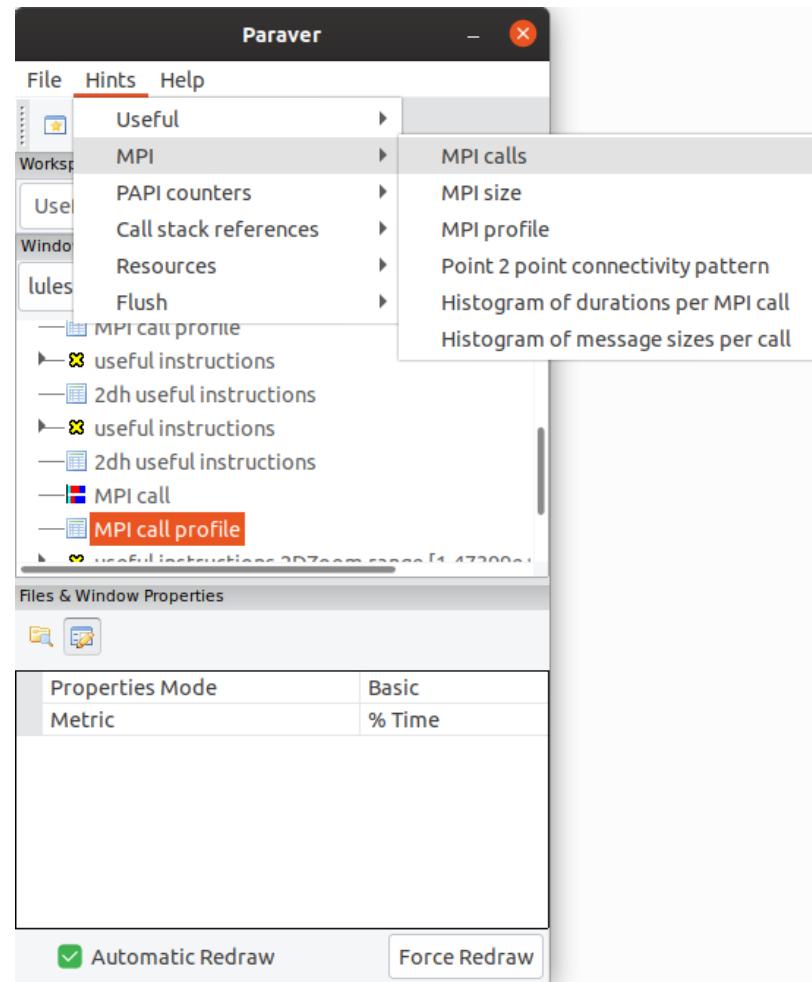
CFG's distribution

- Paraver comes with many more included CFG's



Hints: a good place to start!

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



What is hiding behind CFG's...

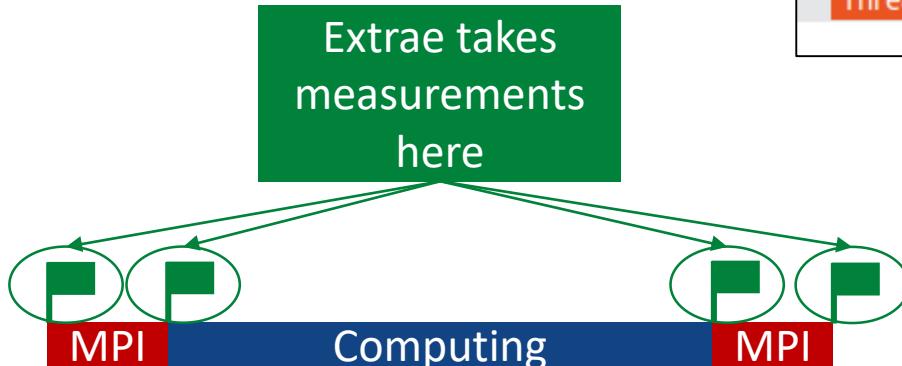
- Paraver is like an equation editor, and CFG's like a formula
- Filter
 - The data I'm interested on seeing (Example: Instructions)

Click & Select from the list

Event Type	Description
40000003 Flushing Traces	Tracing mode:
40000004 Executing CPU	Process IDentifier
40000005 Parent Process IDentifier	
40000038 fork() depth	
40000050 RAW clock() value from system	
40000133 CPU-Event sampling interval	
41000000 Object addresses for task 1.1	
41999999 Active hardware counter set	
42000000 PAPI_L1_DCM [L1D cache misses]	
42000002 PAPI_L2_DCM [L2D cache misses]	
42000008 PAPI_L3_TCM [L3 cache misses]	
42000050 PAPI_TOT_INS [Instr completed]	
42000059 PAPI_TOT_CYC [Total cycles]	

What is hiding behind CFG's...

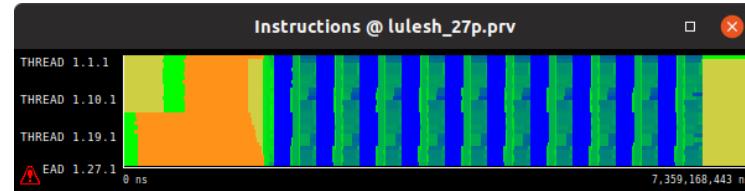
- Semantic
 - How to paint the filtered data
 - Next Evt Val?



Each measurement counts how many Instructions were executed in preceding region

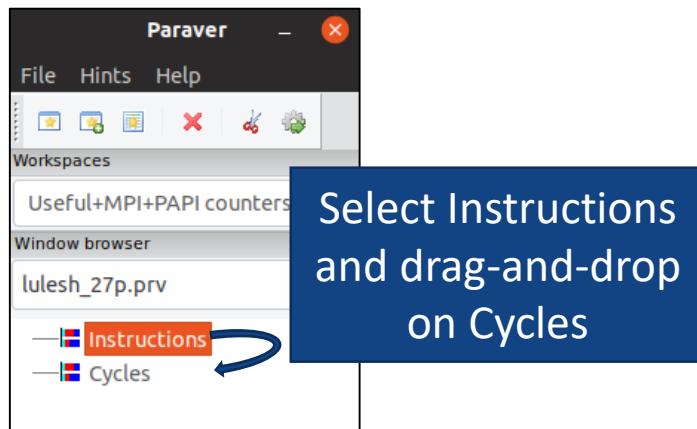
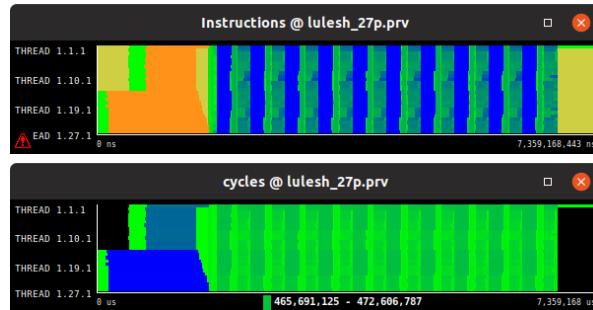
Files & Window Properties	
Semantic	
Top Compose 1	As Is
Top Compose 2	As Is
Compose Thread	As Is
Thread	Next Evt Val

- Instructions count is found at the end of the region
- Next Evt Val = Color is determined by the next event value

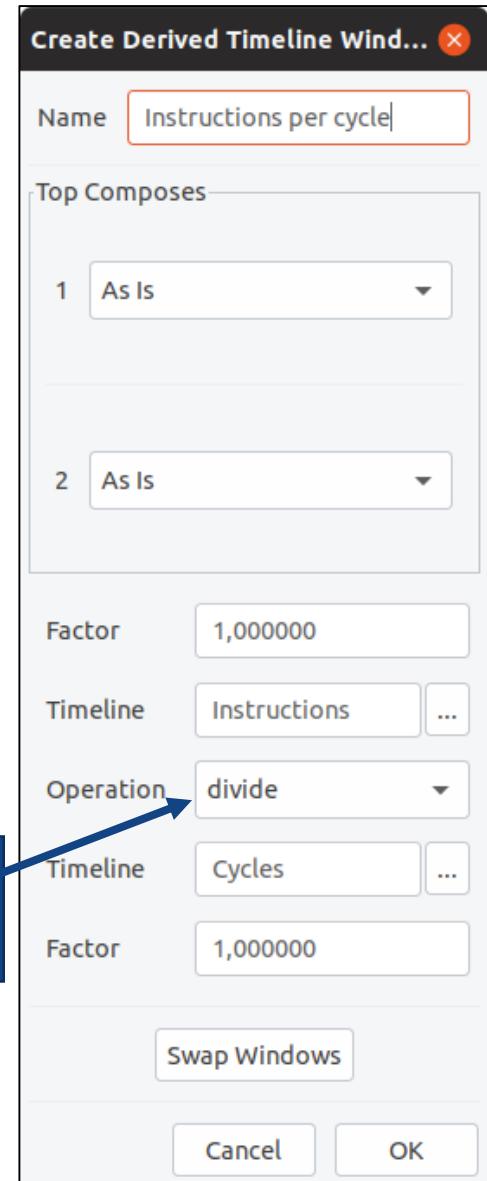


What is hiding behind CFG's...

- Derived
 - Instructions
 - Cycles
 - Instructions / Cycles = IPC

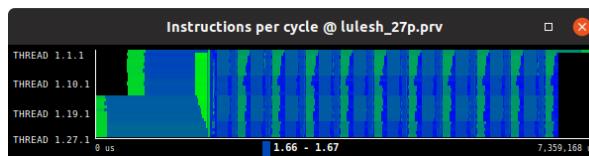
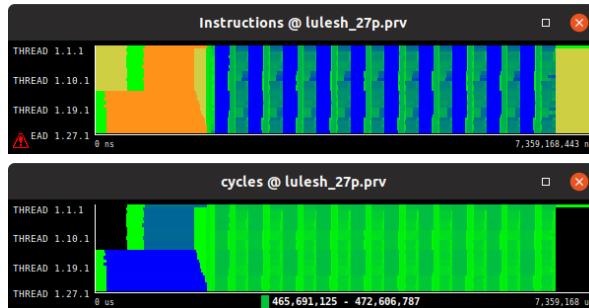


Operation:
divide

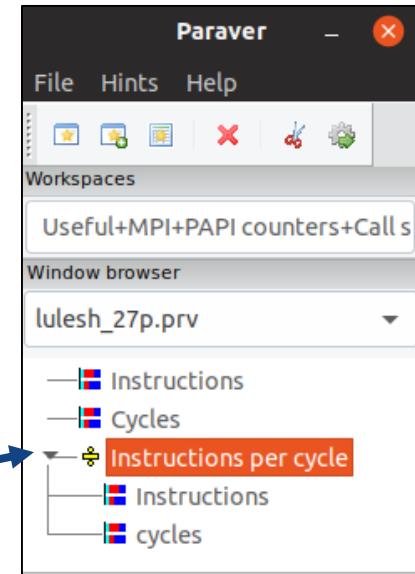


What is hiding behind CFG's...

- Derived
 - Instructions
 - Cycles
 - Instructions / Cycles = IPC



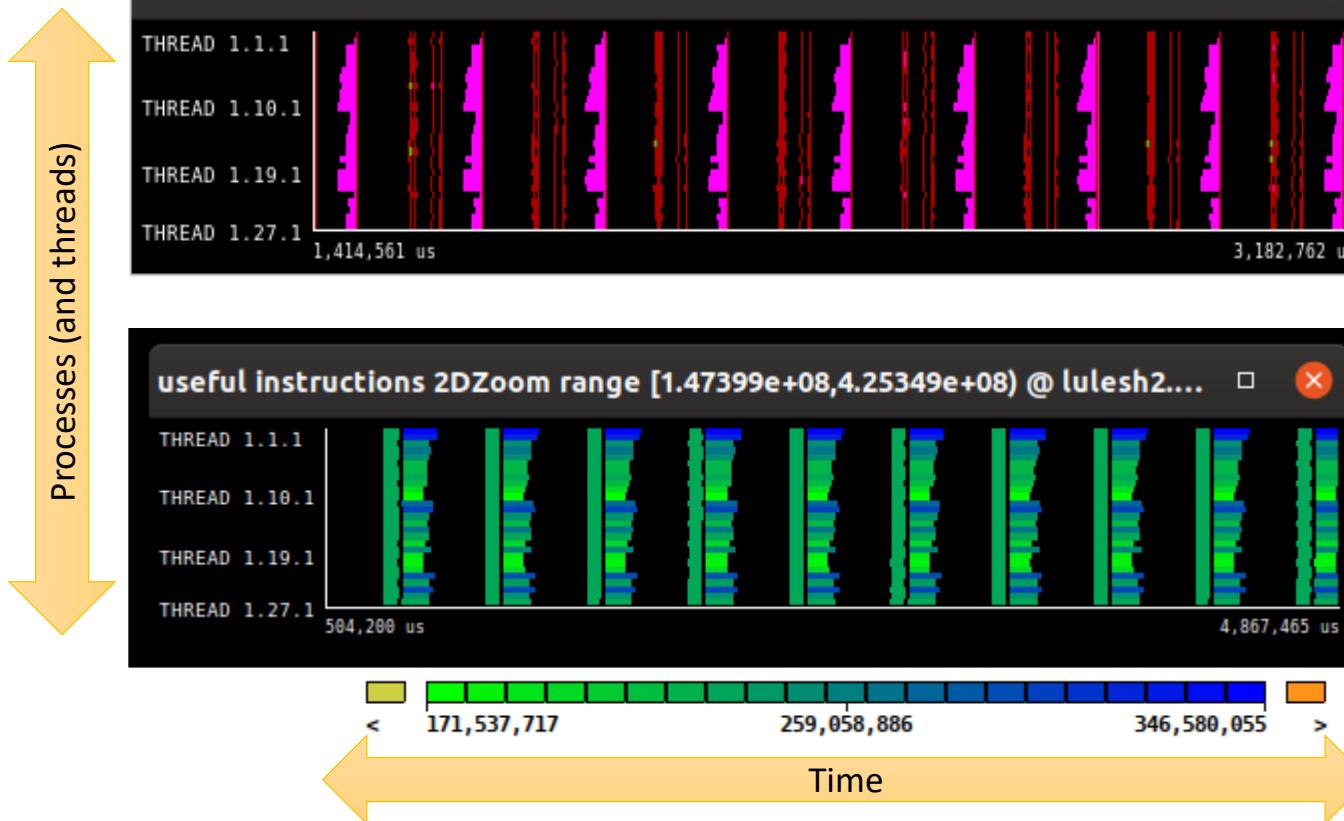
See how a CFG was
constructed



- Build any formula with the data available and save a CFG!

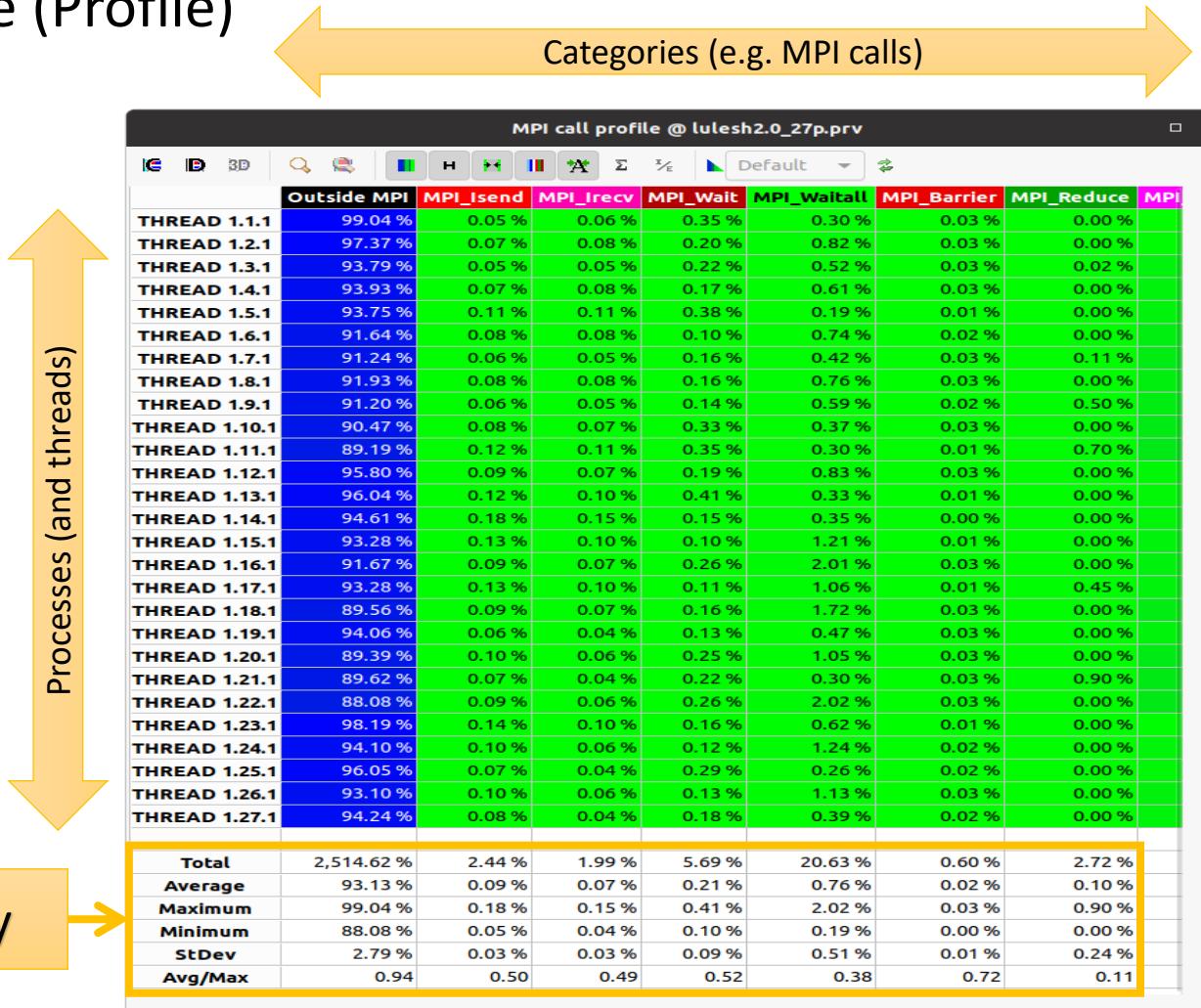
Main views of Paraver (I)

- Timeline



Main views of Paraver (II)

- Table (Profile)



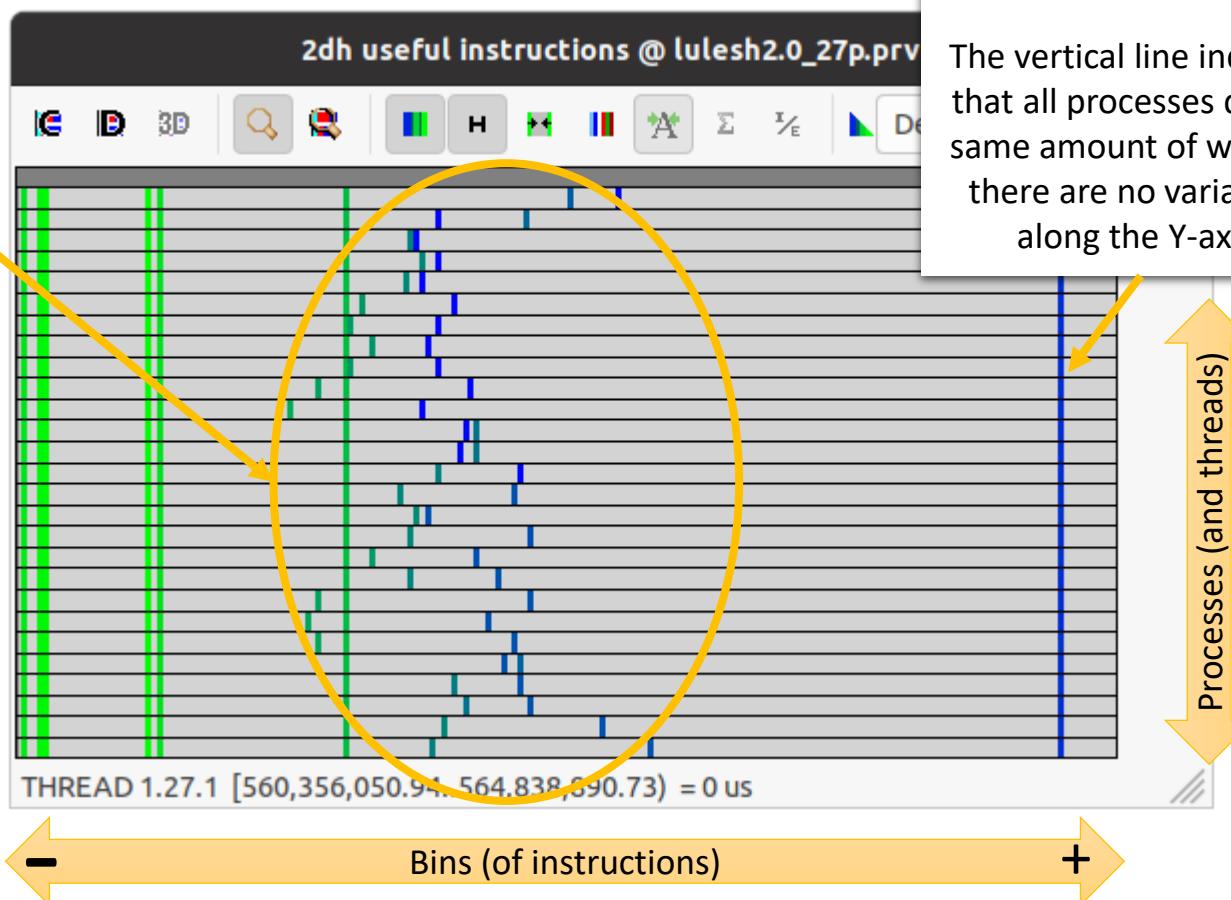
The table can display a variety of statistics (e.g. % of time, # of calls, etc.) with gradient coloring showing from low values to high values

Main views of Paraver (II)

- Histograms

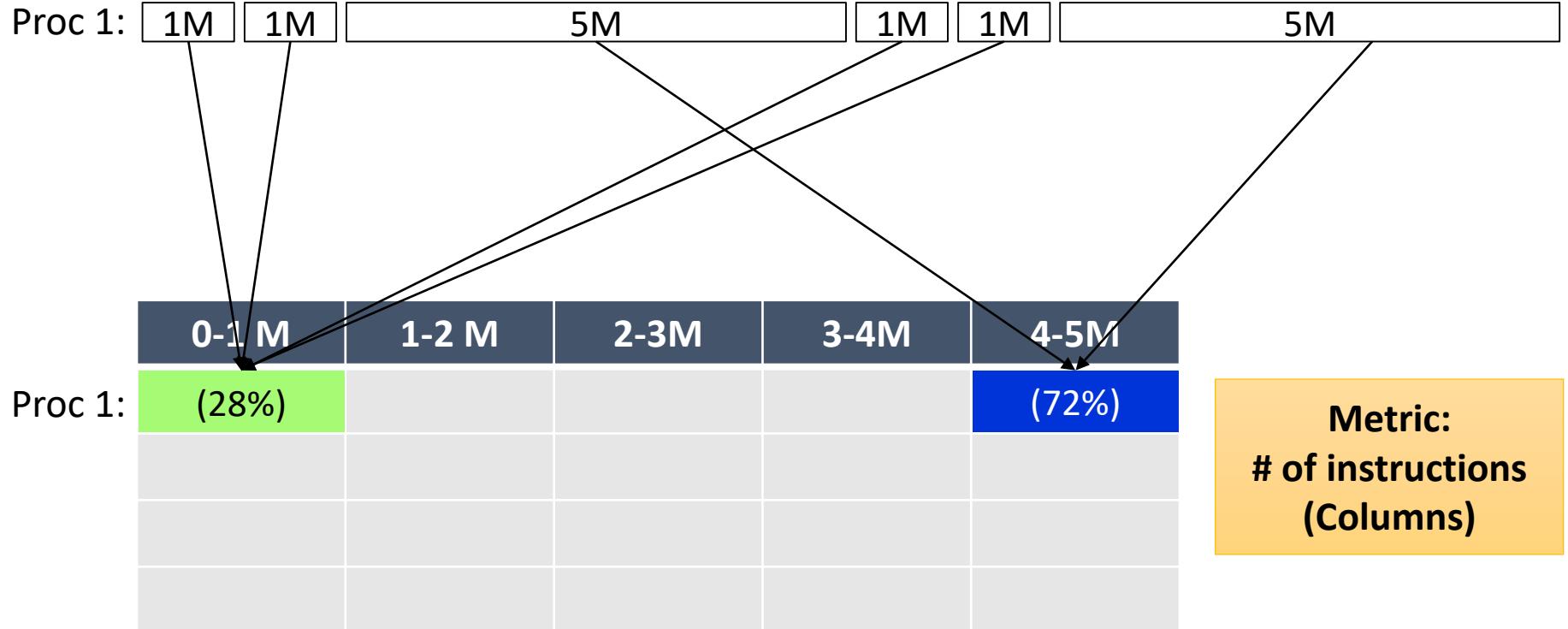
Having data in the middle of the histogram indicates that the program executes a medium-number of instructions.

Being the data scattered (to the left and right), indicates that different processes do different amounts of work (left→less; right→more).

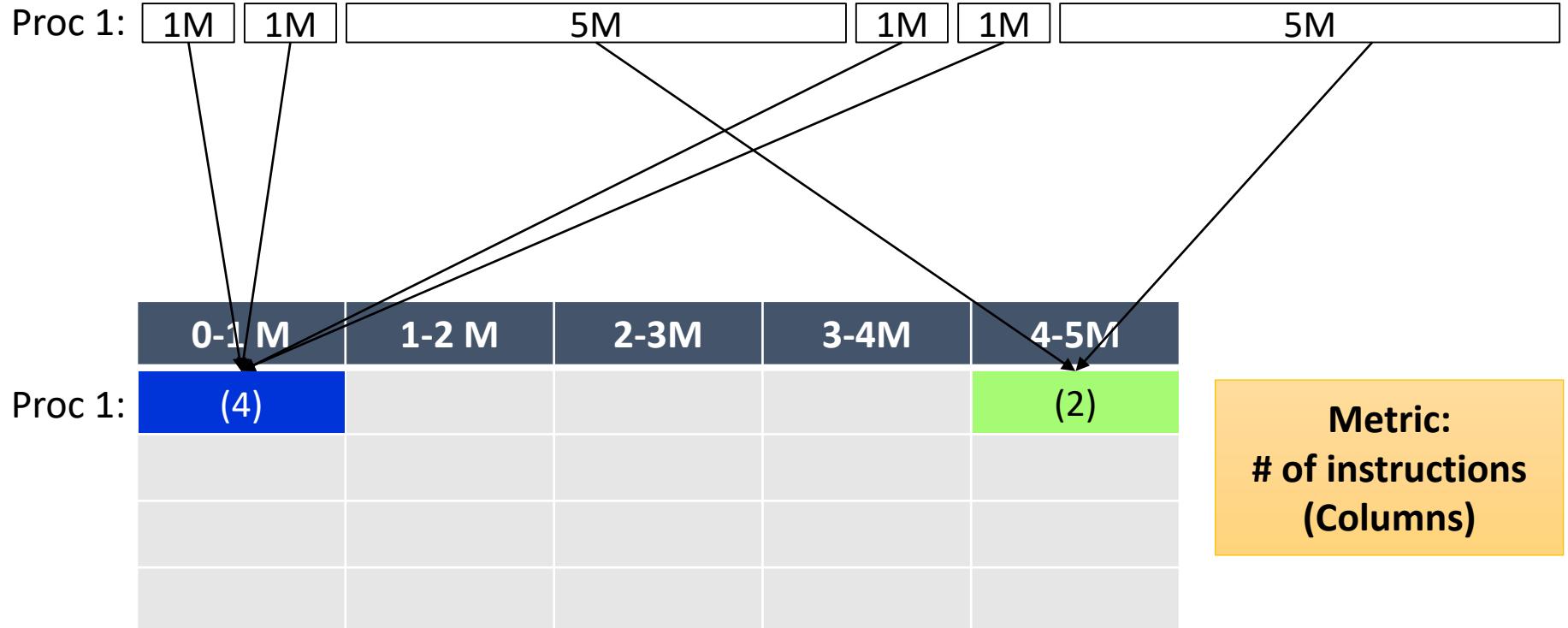


- The more to the left → doing less instructions → greener → low % of exec time
- Data on the right → doing more instructions → bluer → high % of exec time

Histogram: position, metric, statistic... wait, what!?



Histogram: position, metric, statistic... wait, what!?



Histogram: position, metric, statistic... wait, what!?

Proc 1:	1M	1M	5M	1M	1M	5M
Proc 2:	1M	1M	5M	1M	1M	5M
Proc 3:	1M	1M	5M	1M	1M	5M
Proc 4:	1M	1M	5M	1M	1M	5M

	0-1 M	1-2 M	2-3M	3-4M	4-5M
Proc 1:	(4)				(2)
Proc 2:	(4)				(2)
Proc 3:	(4)				(2)
Proc 4:	(4)				(2)

Metric:
**# of instructions
(Columns)**

Statistic: # of bursts (Color)

Histogram: position, metric, statistic... wait, what!?

Proc 1:	1M	1M	5M	1M	1M	5M
Proc 2:	1M	1M	5M	1M	1M	5M
Proc 3:	2M	2M	2M	2M	2M	2M
Proc 4:	1M	1M	5M	1M	1M	5M

	0-1 M	1-2 M	2-3M	3-4M	4-5M	
Proc 1:	(4)				(2)	
Proc 2:	(4)				(2)	
Proc 3:			(7)		(2)	
Proc 4:	(4)				(2)	

Metric:
**# of instructions
(Columns)**

Statistic: # of bursts (Color)



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Getting traces for different scales

Judit Giménez, Lau Mercadal, Germán Llort

✉ tools@bsc.es

2024-09-05

PATC: Performance Analysis and Tools

Let's get more traces!

- What happens if we bind processes to cores?

```
#!/bin/bash

#SBATCH --job-name=lulesh2.0_s65_27p_bound
#SBATCH --output=%x_%j.out
#SBATCH --error=%x_%j.err
#SBATCH --ntasks=27
#SBATCH --time=00:10:00

module load extrae

export TRACE_NAME=lulesh2.0_s65_27p_bound.prv

srun --cpu-bind=cores -m cyclic:cyclic \
    ./trace.sh lulesh2.0 -i 10 -s 65 -p
```

Let's get more traces!

- And if we run with 8 and 64 processes?

```
#!/bin/bash

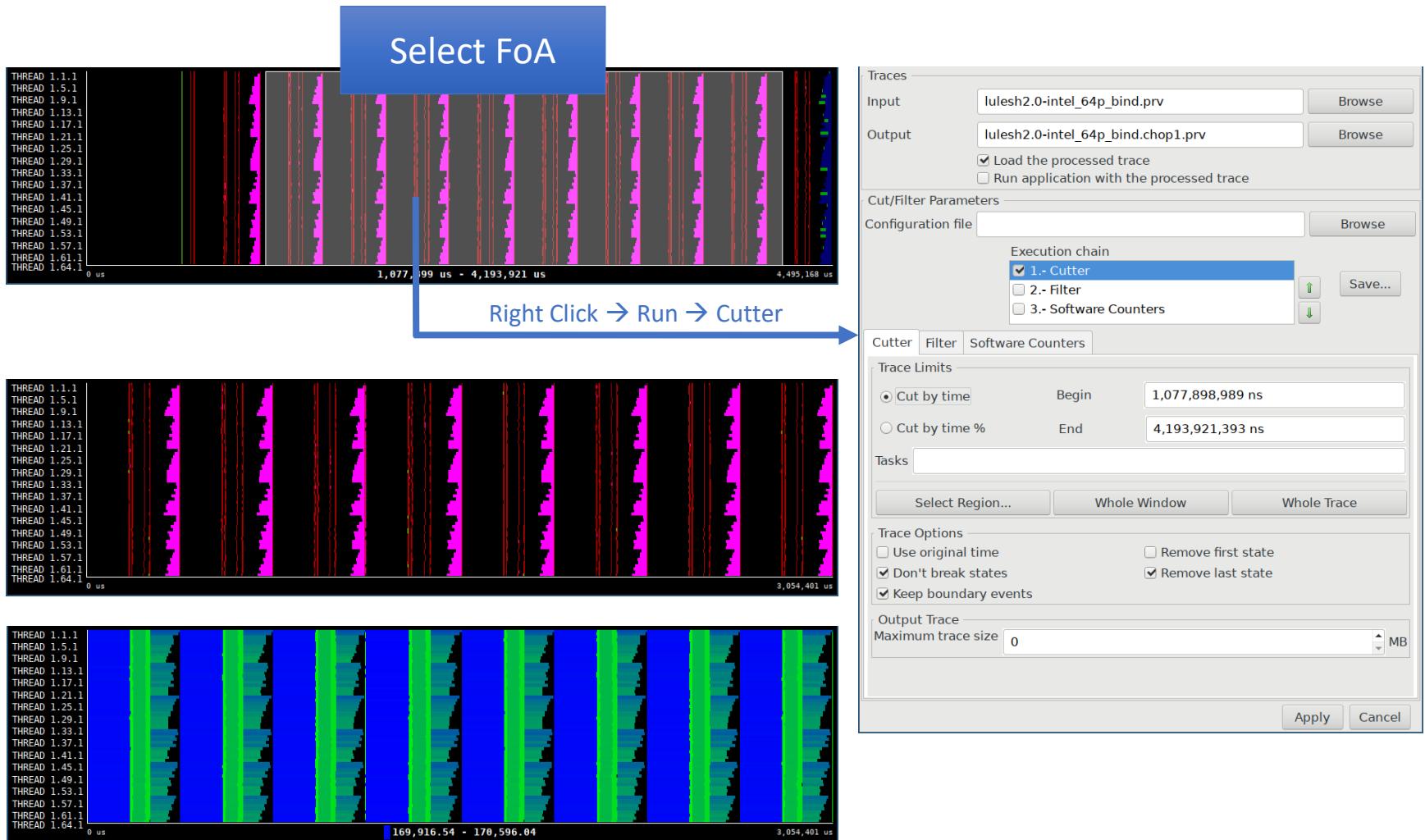
#SBATCH --job-name=lulesh2.0_s65_8p_bound
#SBATCH --output=%x_%j.out
#SBATCH --error=%x_%j.err
#SBATCH --ntasks=8
#SBATCH --time=00:10:00

module load extrae

export TRACE_NAME=lulesh2.0_s65_8p_bound.prv

srun --cpu_bind=cores -m cyclic:cyclic \
./trace.sh ./lulesh2.0 -i 10 -s 65 -p
```

And cut them to select the FoA





**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Getting OpenMP traces

Judit Giménez, Lau Mercadal, Germán Llort

✉ tools@bsc.es

2024-09-05

PATC: Performance Analysis and Tools

Let's get MPI+OpenMP traces!

- What changes do I need in my jobscript?

```
#!/bin/bash

#SBATCH --job-name=lulesh2.0_s65_8p_bound
#SBATCH --output=%x_%j.out
#SBATCH --error=%x_%j.err
#SBATCH --ntasks=8
#SBATCH --cpus-per-task=8
#SBATCH --time=00:10:00

export OMP_NUM_THREADS=8

export TRACE_NAME=lulesh2.0_s65_8p8omp_bound.prv
srun --cpu-bind=cores -m cyclic:cyclic \
    ./trace.sh lulesh2.0-omp -i 10 -s 65 -p
```

Let's get MPI+OpenMP traces!

- What changes do I need in my XML config file?

```
<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>
```

Let's get MPI+OpenMP traces!

- What changes do I need in my `trace.sh` file?

```
#!/bin/bash

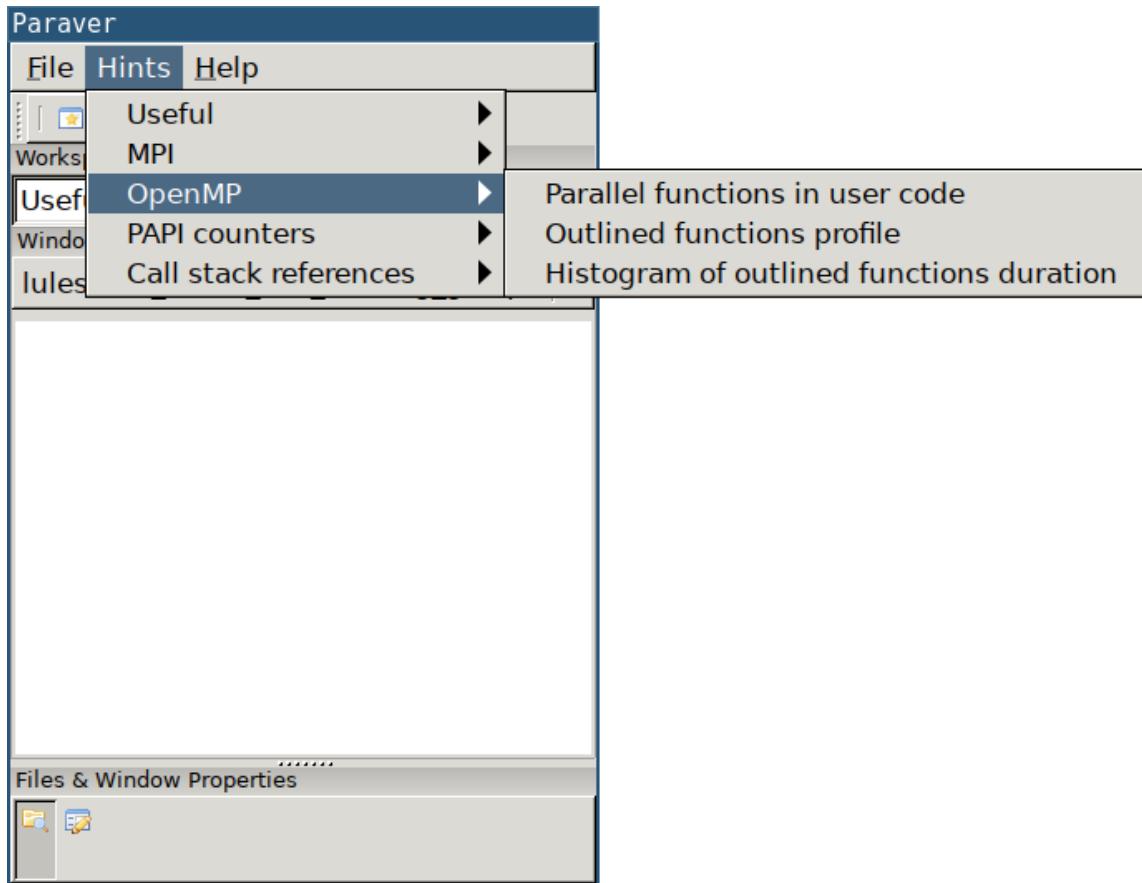
# Configure Extrae
export EXTRAE_HOME=<installation_base_path>
export EXTRAE_CONFIG_FILE=./extrae.xml

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

# Run the program
$*
```

OpenMP Workspace

- Paraver detects OpenMP and suggests views through Hints





**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Thank you!

Judit Giménez, Lau Mercadal, Germán Llort

✉ tools@bsc.es

2024-09-05

PATC: Performance Analysis and Tools