



Selected performance assessments using Extrae and Paraver

Marta Garcia-Gasulla, BSC

HORIZON-EUROHPC-JU-2023-COE



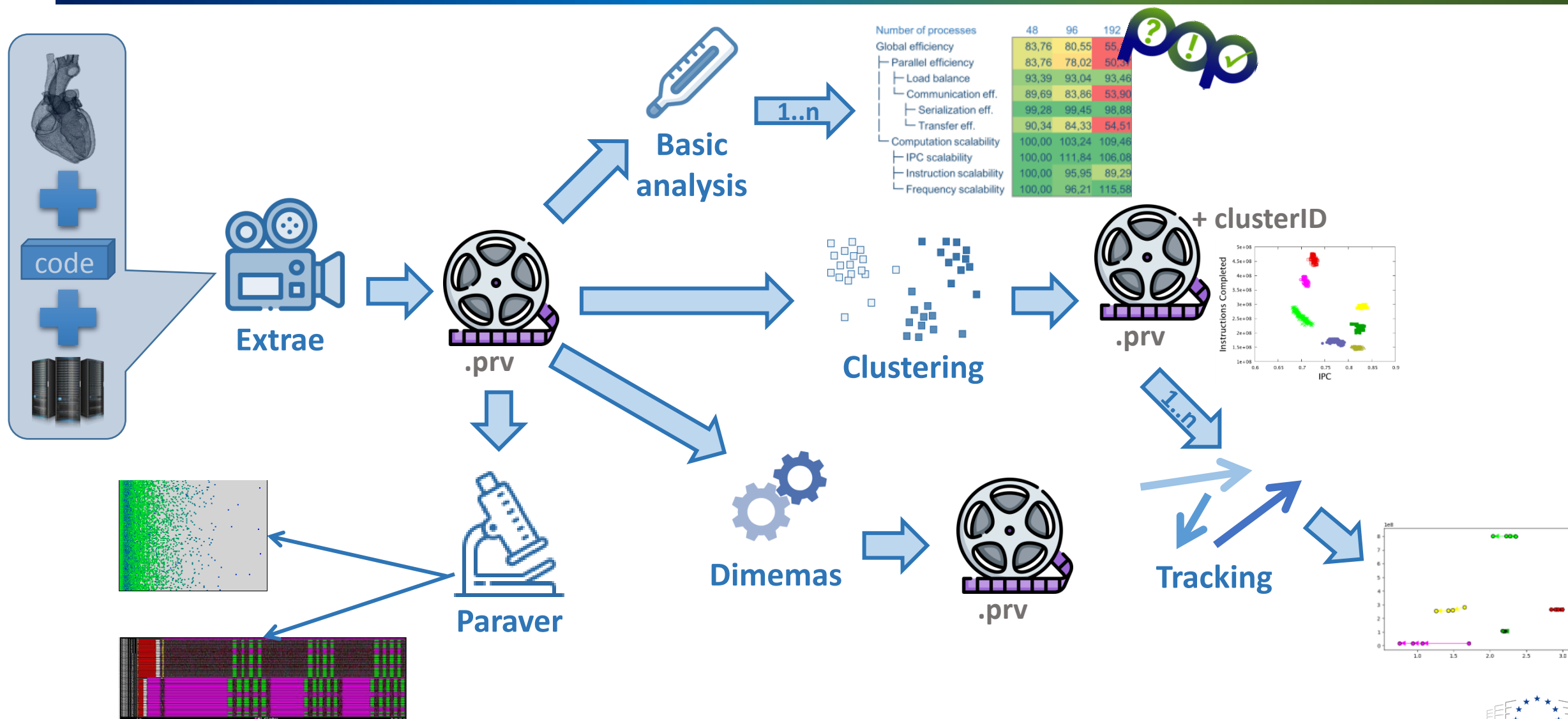
EuroHPC
Joint Undertaking

Grant Agreement No 101143931

1 January 2024– 31 December 2026

- The ecosystem of BSC performance tools
 - Extrae
 - Paraver
 - Dimemas
- Show case of a performance analysis using the BSC performance tools

The ecosystem of BSC performance tools



Determining Parallel Application Execution Efficiency and Scaling using the POP Methodology (ISC24)



Extrae

No need to
recompile
or relink

- Tracing library transparent to the application
- Platform agnostic
- Parallel programming models:
 - MPI, OpenMP, pthreads, OmpSs, CUDA, OpenCL, Java, Python
- Hardware counters
 - Through PAPI
- Link to source:
 - Callstack at MPI routines
 - OpenMP outlined routines
 - Selected user functions (Dyninst)
- Periodic sampling
- User friendly API to annotate your code with custom events

How to use Extrae



- Symbol substitution through LD_PRELOAD

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

- Specific libraries for each runtime and combinations

- MPI
- OpenMP
- OpenMP+MPI
- ...

libmpitrace.so	libompitrace.so
libmpitracecf.so	libptmpitrace.so
libmpitracecf.so	libptmpitracecf.so
libompitrace.so	libptmpitracecf.so
libompitracecf.so	libpttrace.so
libompitracecf.so	libseqtrace.so

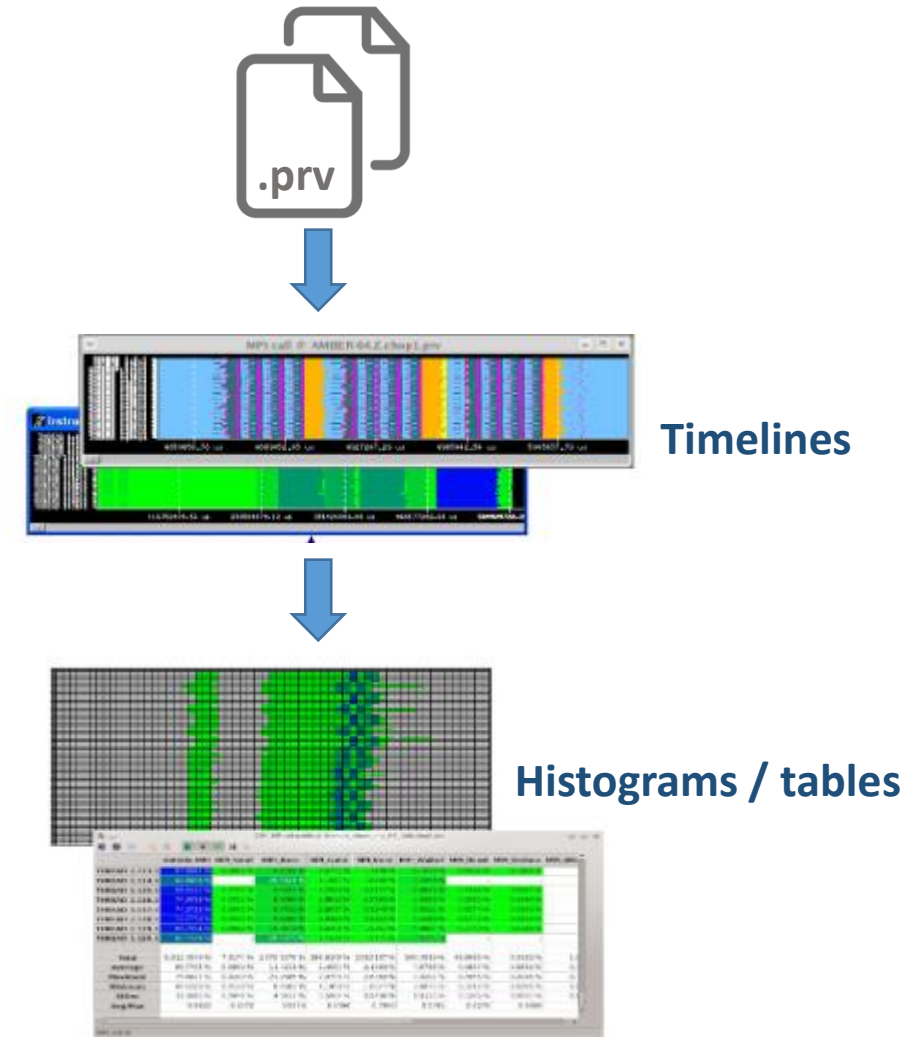
- Detailed configuration in XML file

```
export EXTRAE_CONFIG_FILE=../extrae.xml
```



Paraver

- (Performance) Data Browser
 - Any kind of timestamped data
 - Trace Visualization and analysis
 - Trace manipulation
 - Flexible
 - No pre-assumed semantics
 - Fully programmable
- 2 Kind of views:
 - Timeline
 - Histograms, 2D and 3D tables
- Multiple loaded traces
 - Allow comparative analysis



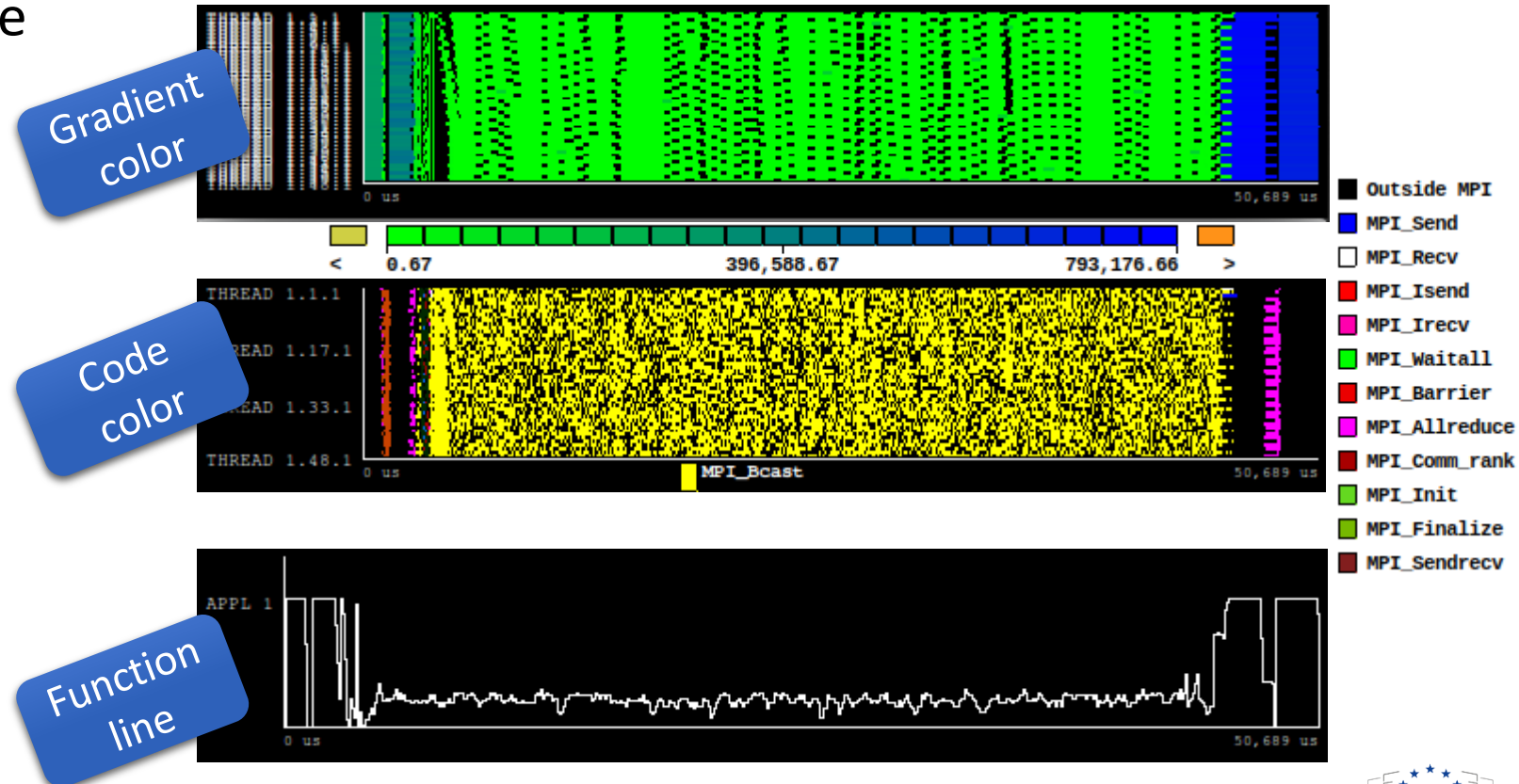
Timelines



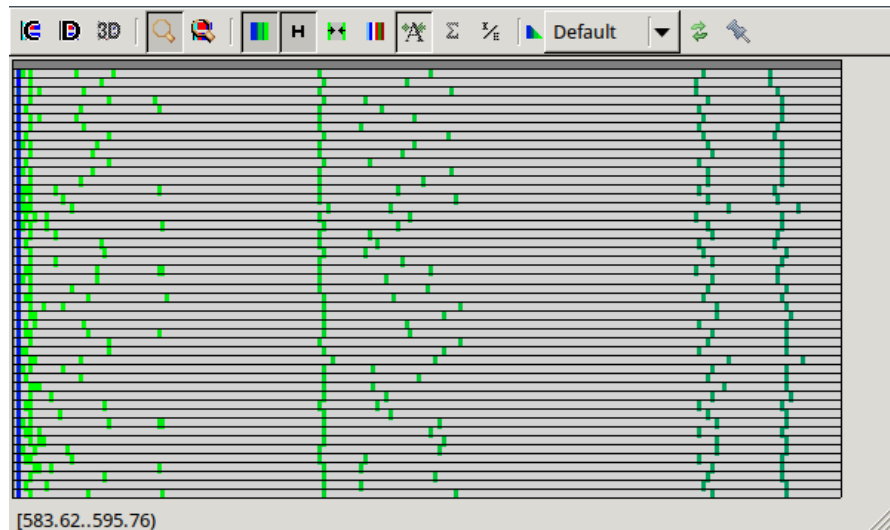
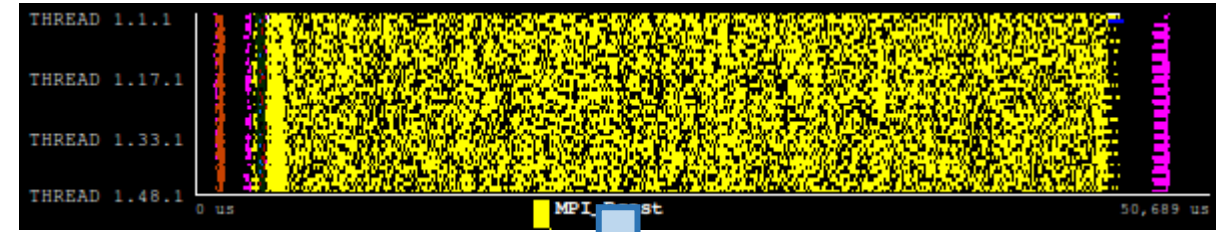
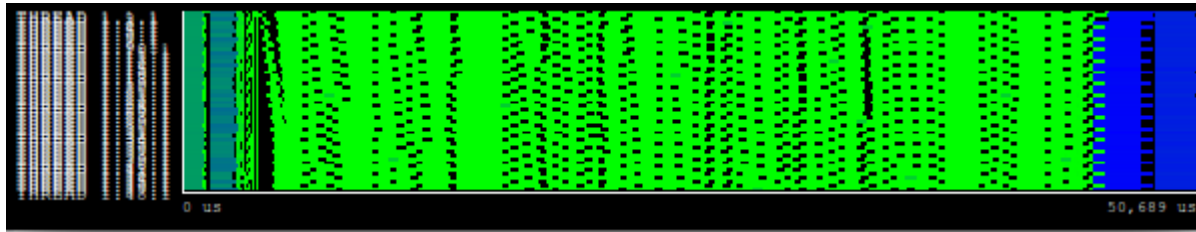
Separate the *What* from the *How*

- *What* -> Semantics
 - Which event to visualize
 - Which value to use

- *How* -> Visual representation



Histograms



	Outside MPI	MPI_Send	MPI_Recv	MPI_Bcast	MPI_Allreduce	MPI_Comm_rank	MPI_Comm_size
THREAD 1.1.1	23.74 %	0.26 %	1.40 %	72.80 %	0.48 %	0.03 %	0.03 %
THREAD 1.2.1	24.55 %	0.19 %	0.23 %	72.50 %	0.77 %	0.03 %	0.03 %
THREAD 1.3.1	22.63 %	1.47 %	0.34 %	73.42 %	0.38 %	0.03 %	0.03 %
THREAD 1.4.1	23.79 %	-	-	72.27 %	2.29 %	0.03 %	0.03 %
THREAD 1.5.1	23.47 %	-	-	73.83 %	1.06 %	0.03 %	0.02 %
THREAD 1.6.1	23.34 %	-	-	72.99 %	2.03 %	0.03 %	0.03 %
THREAD 1.7.1	22.74 %	-	-	73.19 %	2.36 %	0.04 %	0.03 %
THREAD 1.8.1	23.58 %	-	-	74.02 %	0.69 %	0.03 %	0.03 %
THREAD 1.9.1	22.66 %	-	-	73.59 %	2.05 %	0.04 %	0.03 %
THREAD 1.10.1	22.93 %	-	-	73.34 %	2.10 %	0.04 %	0.02 %
THREAD 1.11.1	23.02 %	-	-	74.22 %	1.11 %	0.03 %	0.03 %
THREAD 1.12.1	23.94 %	-	-	72.64 %	1.78 %	0.03 %	0.03 %
THREAD 1.13.1	23.04 %	-	-	73.28 %	1.99 %	0.04 %	0.03 %
THREAD 1.14.1	24.67 %	-	-	72.64 %	0.99 %	0.03 %	0.03 %
THREAD 1.15.1	22.88 %	-	-	73.61 %	1.80 %	0.03 %	0.03 %
THREAD 1.16.1	22.96 %	-	-	73.46 %	1.98 %	0.04 %	0.03 %
THREAD 1.17.1	22.96 %	-	-	74.58 %	0.88 %	0.03 %	0.03 %
THREAD 1.18.1	23.72 %	-	-	72.61 %	2.09 %	0.03 %	0.03 %
THREAD 1.19.1	23.33 %	-	-	72.71 %	2.28 %	0.04 %	0.03 %
THREAD 1.20.1	22.98 %	-	-	74.26 %	1.09 %	0.03 %	0.03 %
THREAD 1.21.1	23.03 %	-	-	73.04 %	2.27 %	0.03 %	0.03 %
THREAD 1.22.1	22.96 %	-	-	73.35 %	2.06 %	0.04 %	0.02 %

Determining Parallel Application Execution Efficiency and Scaling using the POP Methodology (ISC24)

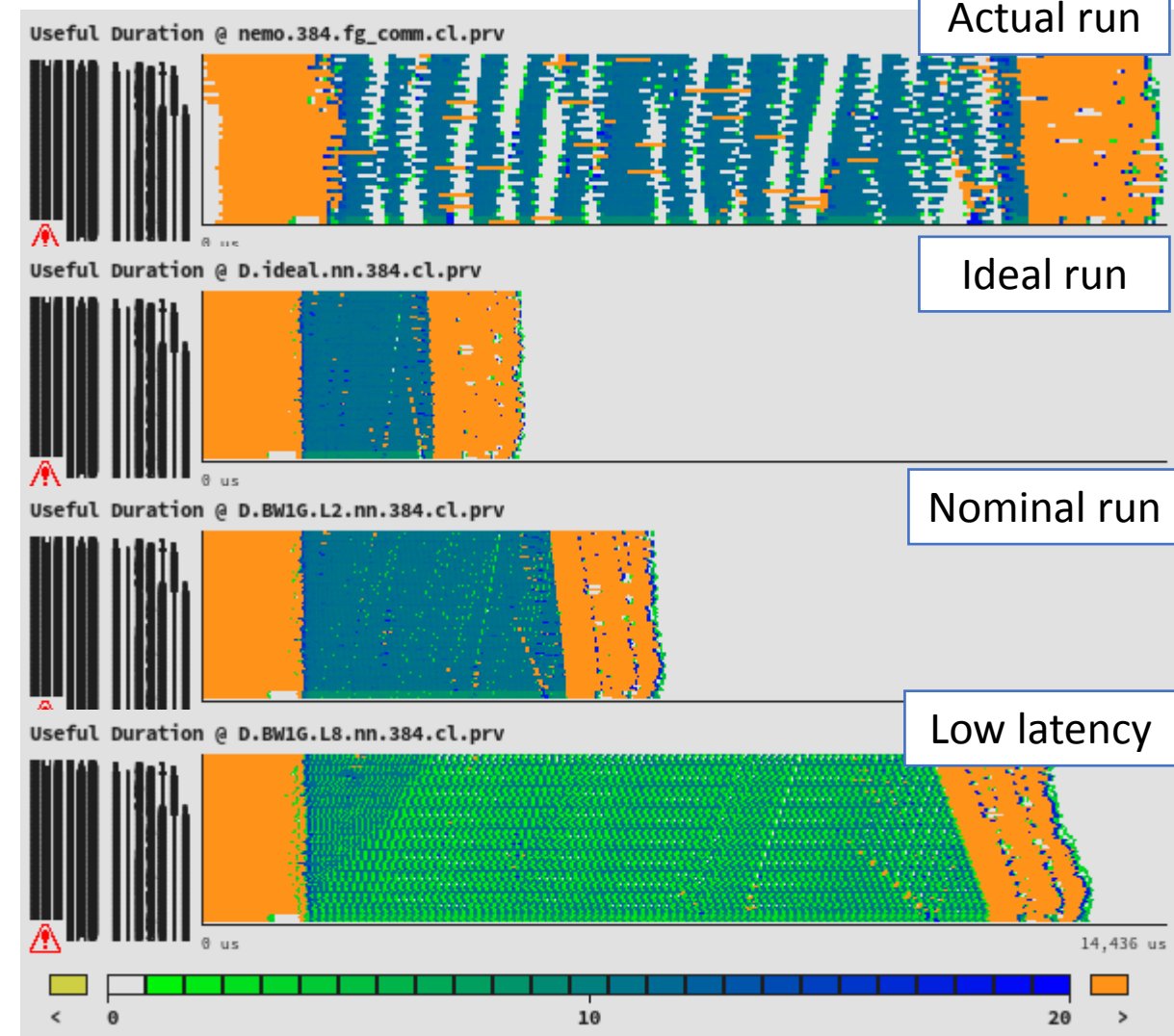
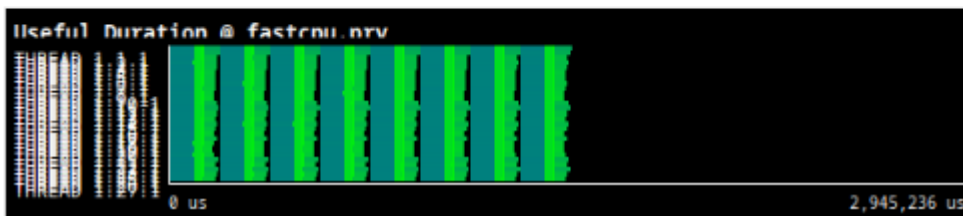
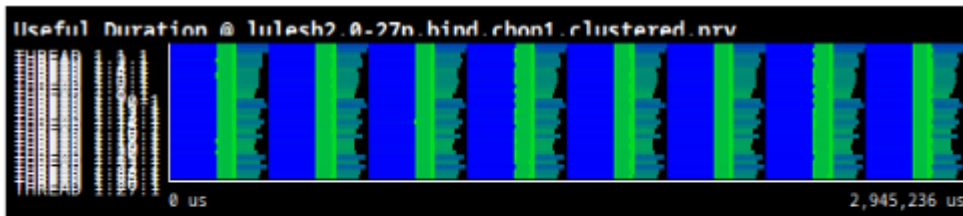
Trace Manipulation



- Traces can be manipulated
 - Chop: cut in time or processes
 - Filter: Subset of records based on...
 - ... duration, type, value...
- The output is still a Paraver trace



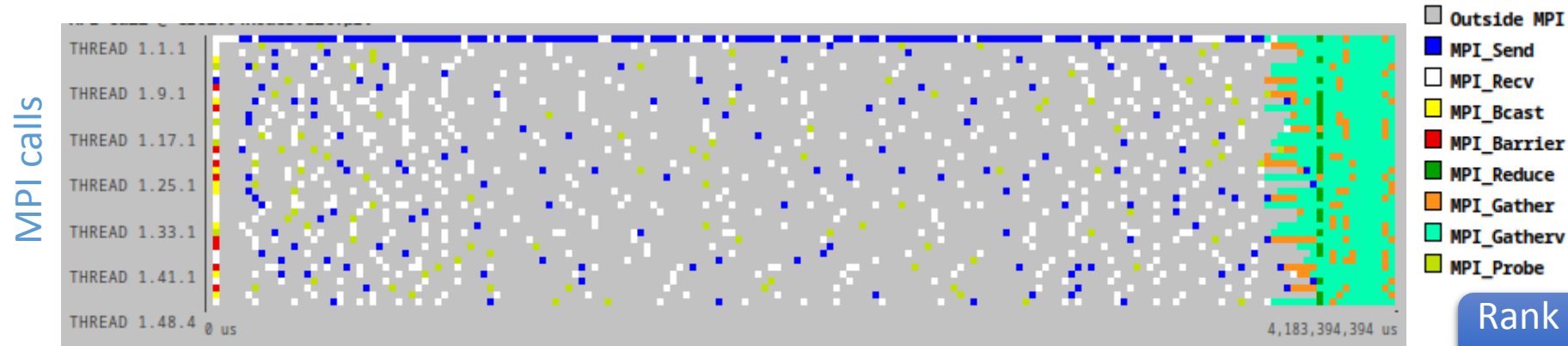
- Change parameters from the machine:
 - Network: Latency, bandwidth
 - CPU frequency



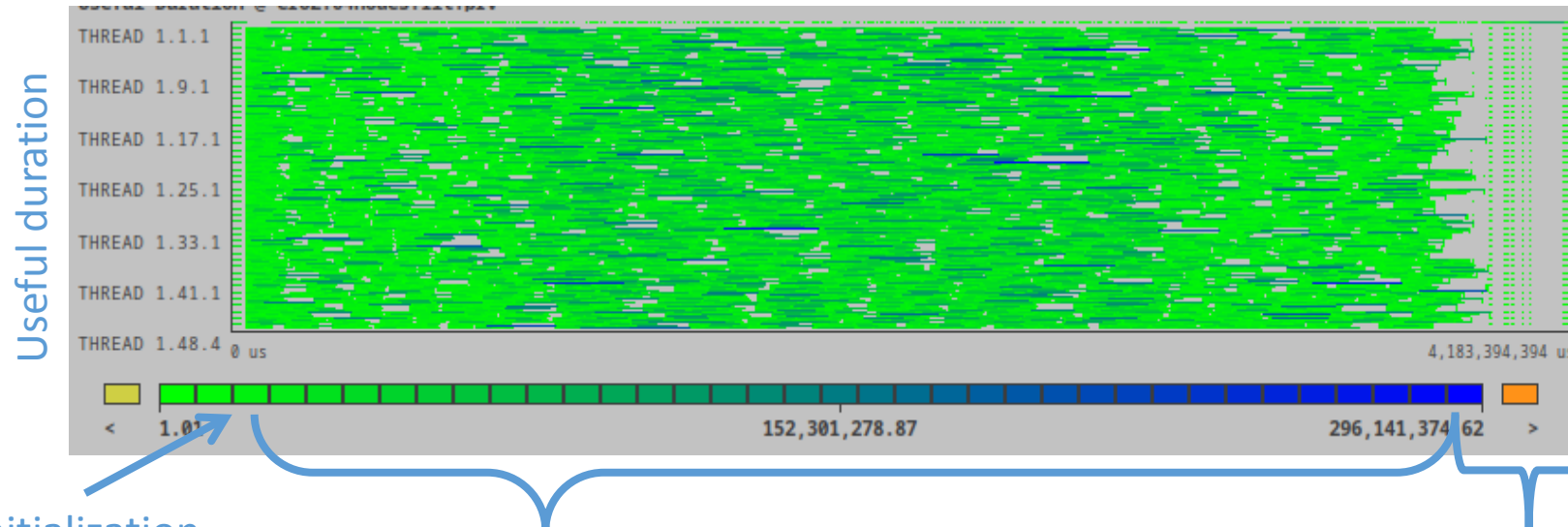


Performance analysis of a hybrid code with BSC tools

Structure and Focus of Analysis (FoA)



Rank 1 seems to be always in communication



Randomly distributed communication for other ranks.

Initialization

Computation

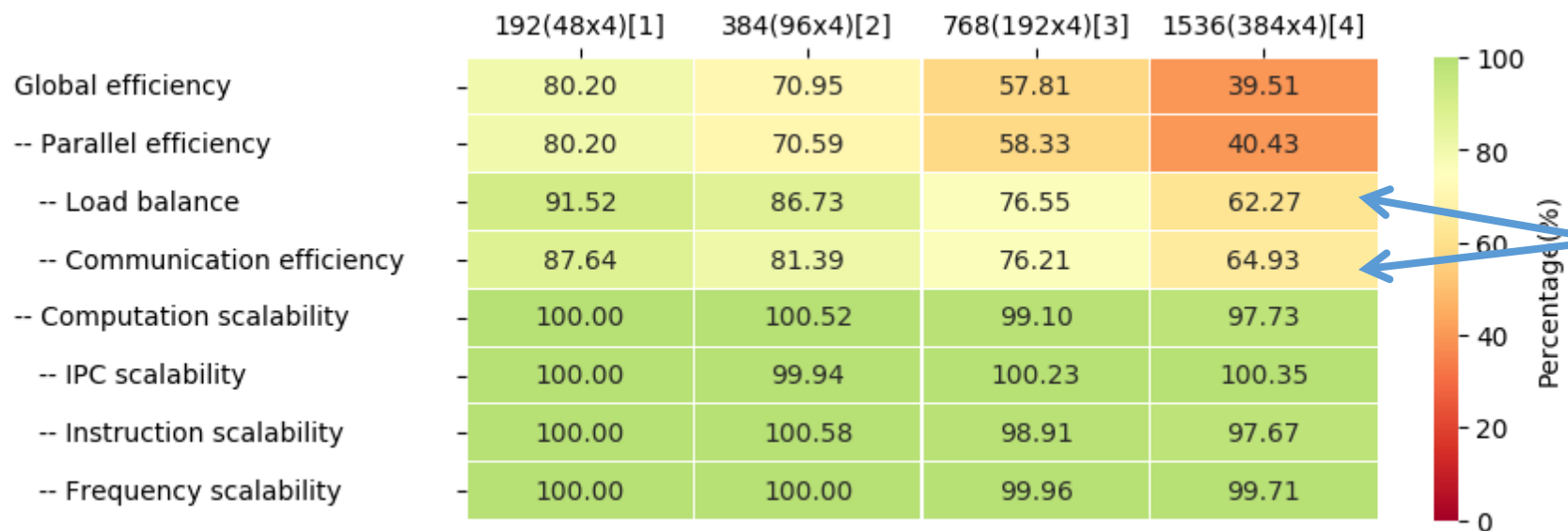
Gather/Reduction

Determining Parallel Application Execution Efficiency and Scaling using the POP Methodology (ISC24)

Efficiency metrics

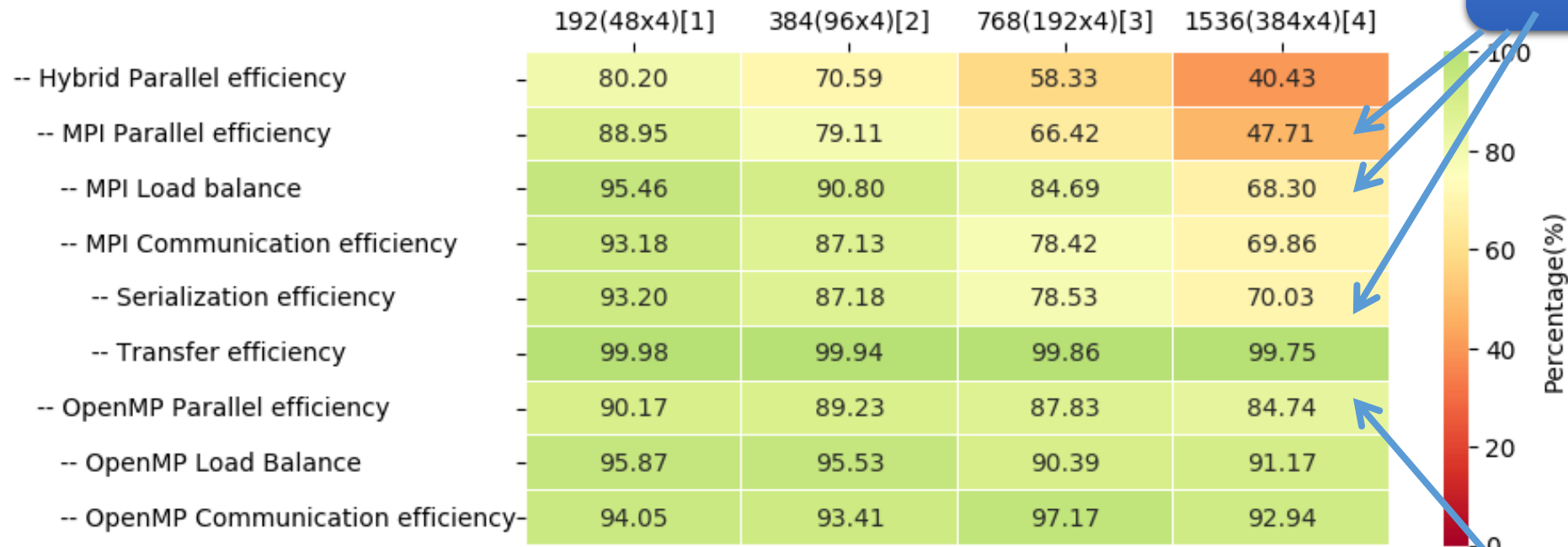


Scaling number of processes (48 → 384)
Keeping constant number of threads (4)



Main factors affecting scalability:
- Load Balance
- Communication efficiency

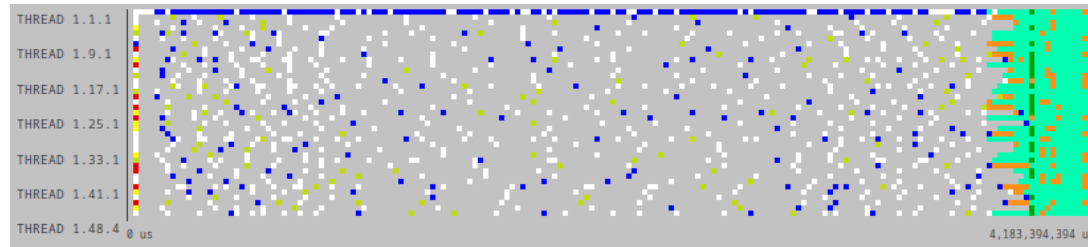
Who to blame MPI or OpenMP?



Very low MPI parallel efficiency due to Load Balance and Serialization

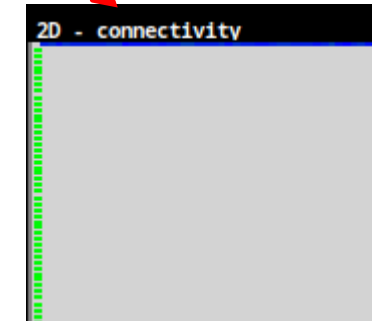
OpenMP better parallel efficiency, still a trend to decrease

Communication pattern

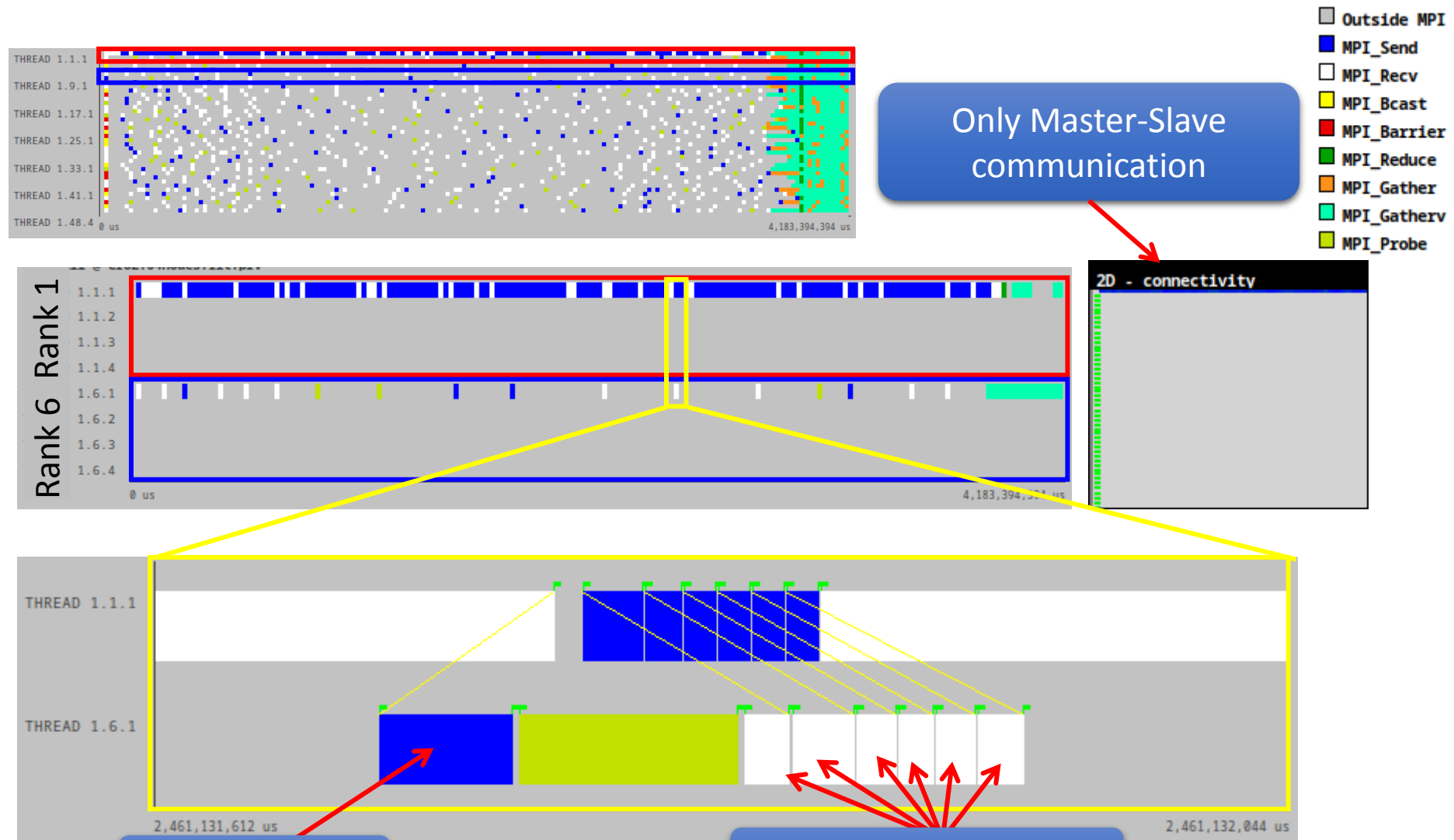


Only Master-Slave
communication

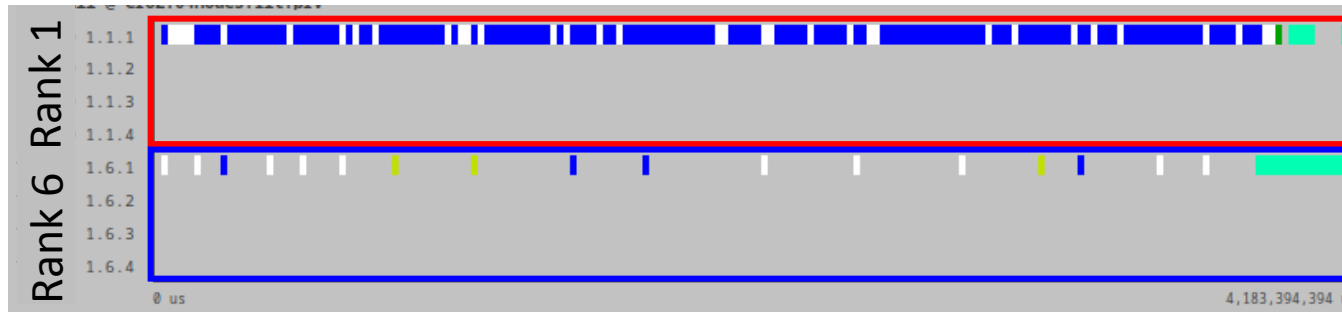
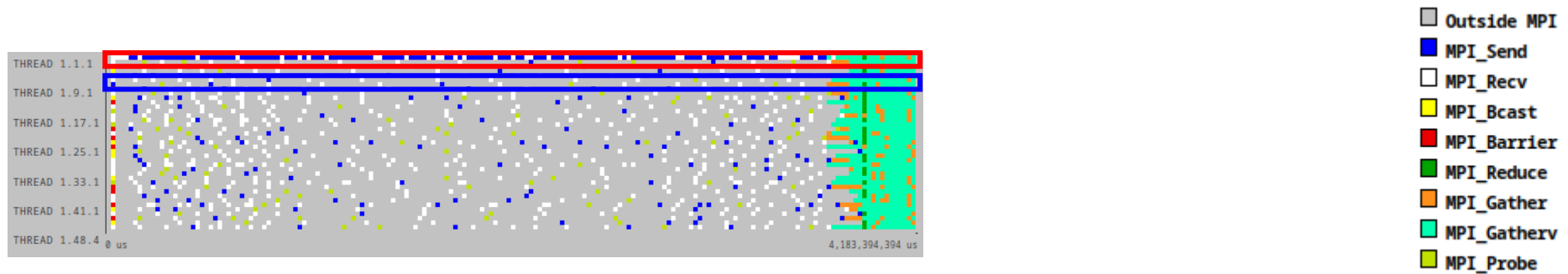
- Outside MPI
- MPI_Send
- MPI_Recv
- MPI_Bcast
- MPI_Barrier
- MPI_Reduce
- MPI_Gather
- MPI_Gatherv
- MPI_Probe



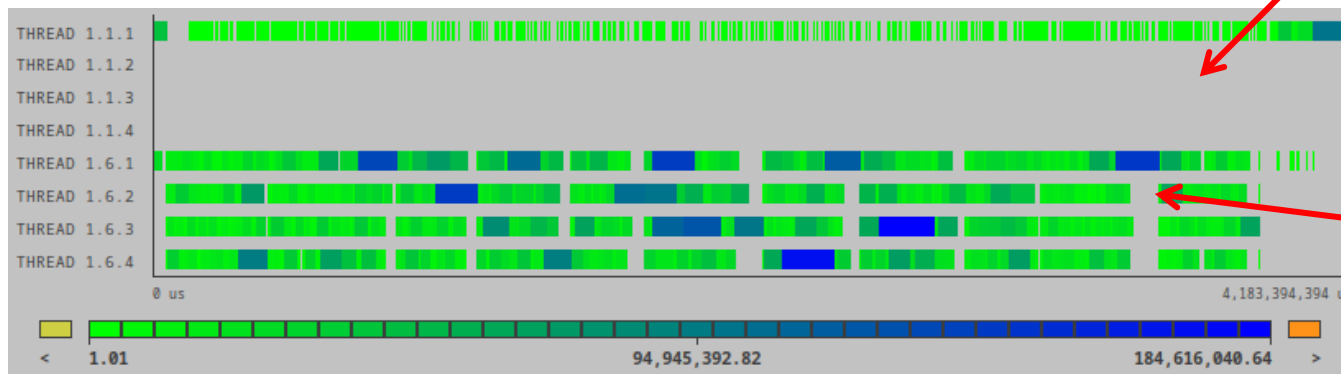
Communication pattern



Communication pattern



The threads of the master process don't perform useful computation

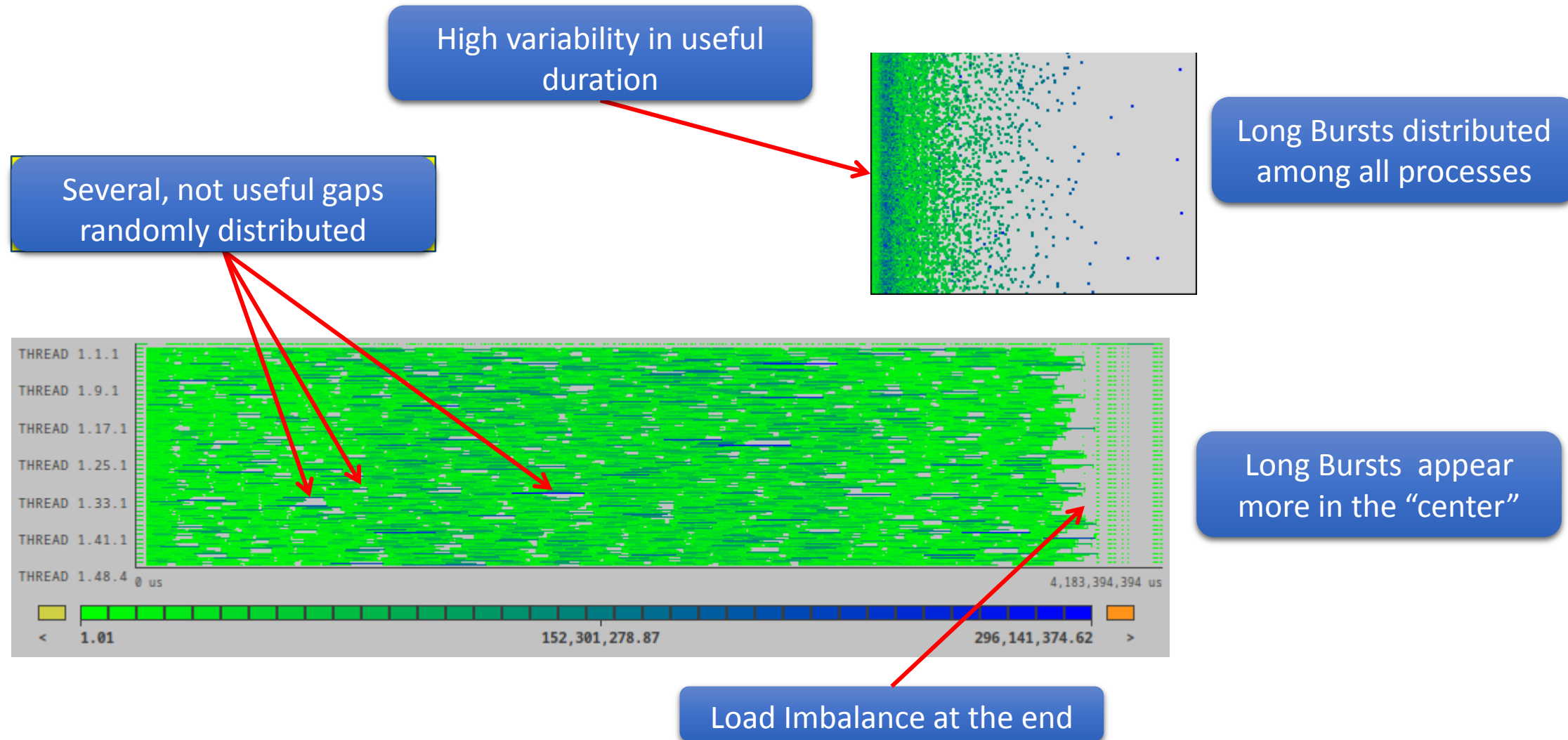


Worker rank do useful computation at all threads

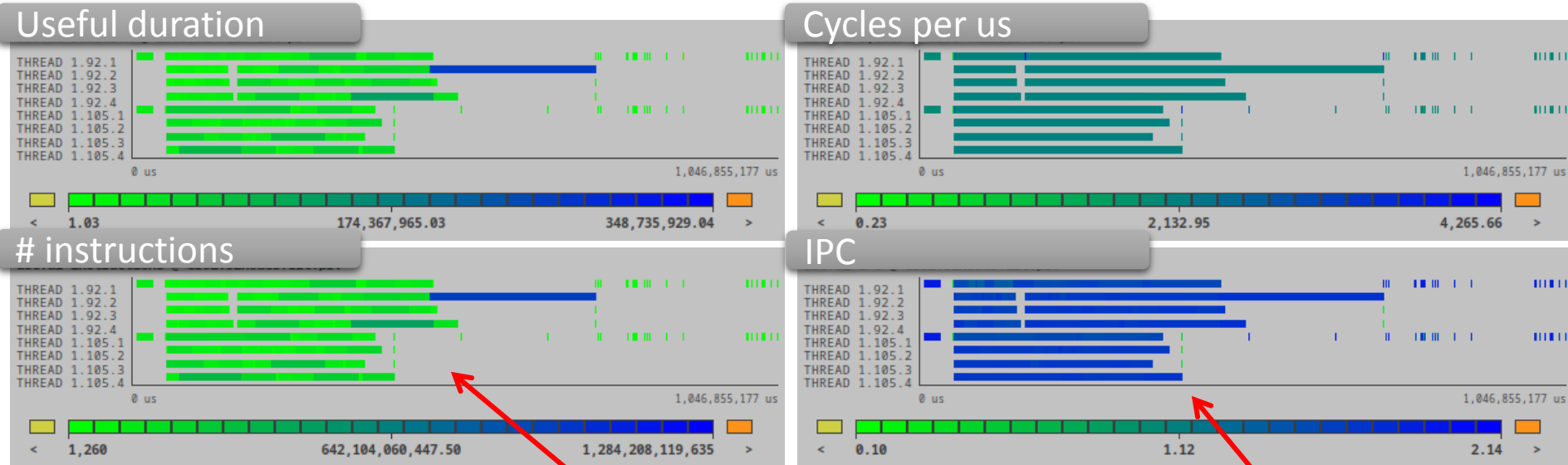
Gaps happen when one thread has longer useful computation.

Determining Parallel Application Execution Efficiency and Scaling using the POP Methodology (ISC24)

Load Imbalance



The source of the Load imbalance

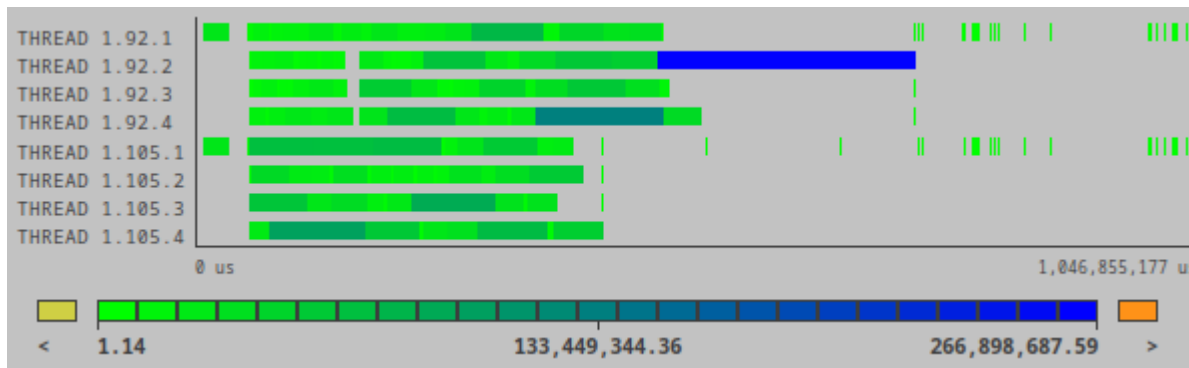
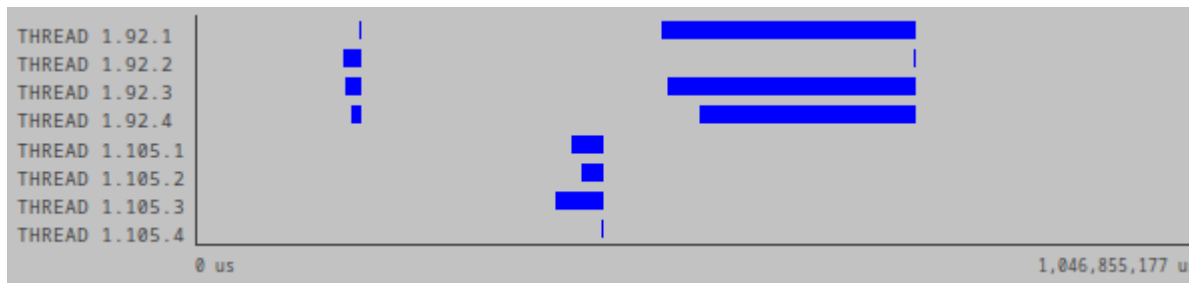
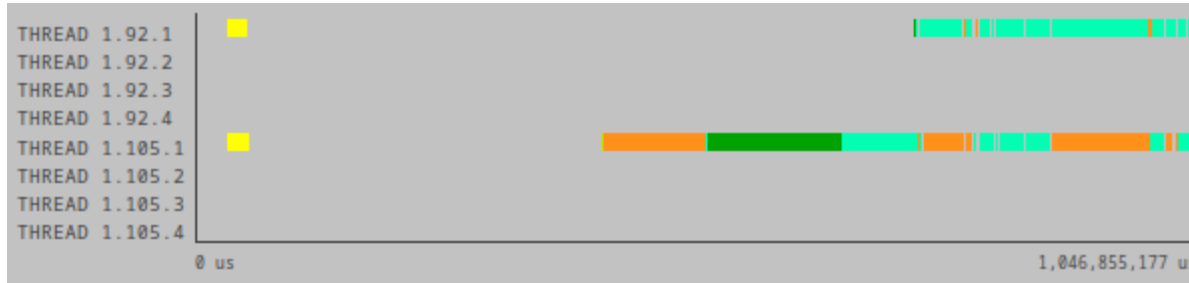


The amount of instructions correlates with the burst duration.

Load balance seems to come due to more instructions. Hence more code execution.

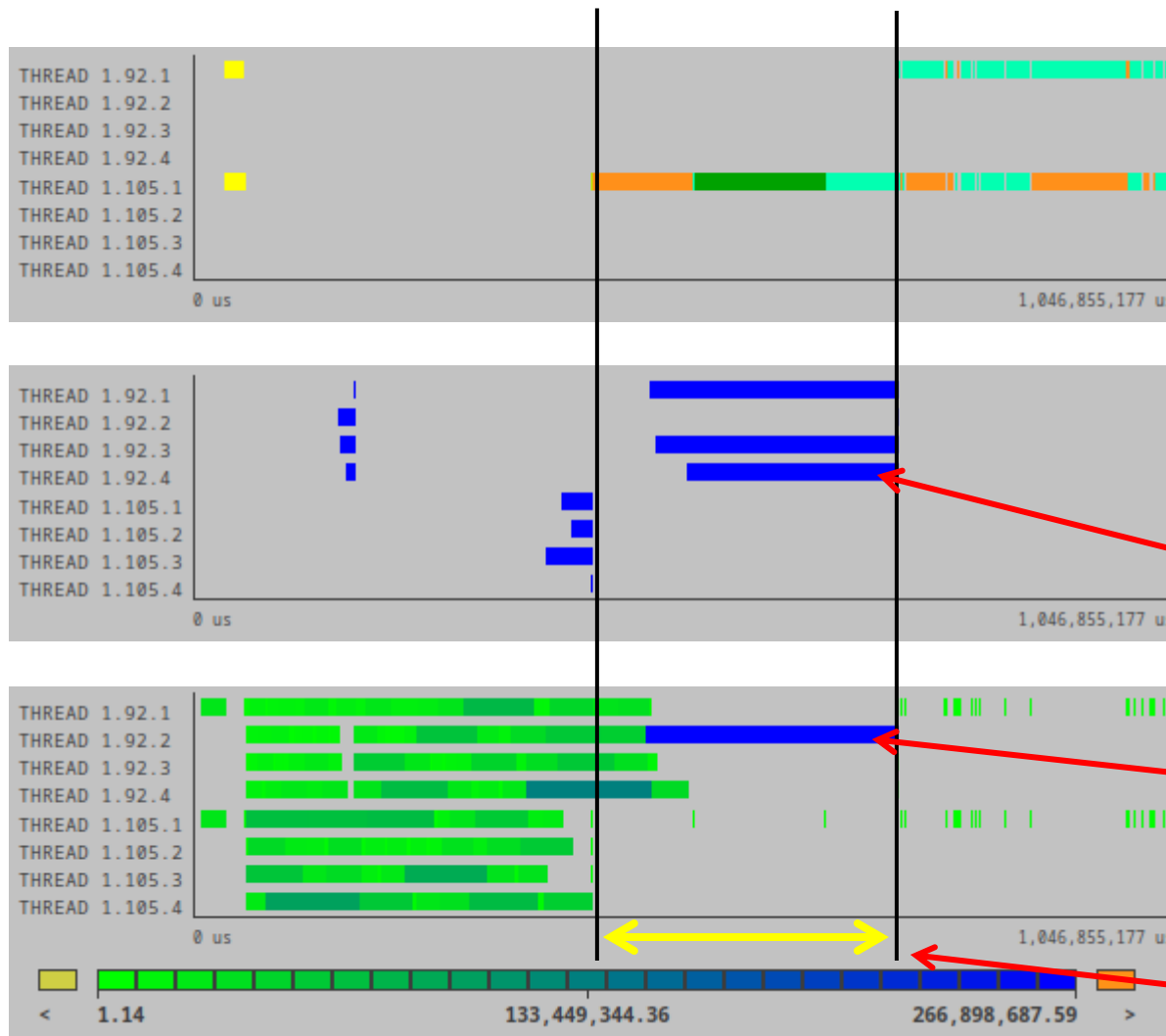
On the other hand, the cycles and IPC are stable for most bursts.

Load imbalance in detail



Determining Parallel Application Execution Efficiency and Scaling using the POP Methodology (ISC24)

Load imbalance in detail



OpenMP Load Imbalance

One burst of compute takes a lot much longer than the others

MPI Load Imbalance

Load imbalance



Application parameter to set size of chunk sent.
Reduce it to improve MPI
Load imbalance

maxMpiChunkSize:	50	10
	1536(384x4)[1]	1536(384x4)[2]
Global efficiency	40.43	37.57
-- Parallel efficiency	40.43	33.39
-- Load balance	62.27	62.37
-- Communication efficiency	64.93	53.53
-- Computation scalability	100.00	112.51
-- IPC scalability	100.00	100.31
-- Instruction scalability	100.00	112.47
-- Frequency scalability	100.00	99.73
	1536(384x4)[1]	1536(384x4)[2]
-- Hybrid Parallel efficiency	40.43	33.39
-- MPI Parallel efficiency	47.71	53.22
-- MPI Load balance	68.30	74.46
-- MPI Communication efficiency	69.86	71.48
-- Serialization efficiency	70.03	71.65
-- Transfer efficiency	99.75	99.77
-- OpenMP Parallel efficiency	84.74	62.73
-- OpenMP Load Balance	91.17	83.77
-- OpenMP Communication efficiency	92.94	74.89

MPI grain reduction

Overall Worst Load Balance

We tradeoff the OpenMP
Load Balance for the MPI

- Main issue caused by unpredictable long chunks of computation
 - Affects both MPI and OpenMP
- Suggestions:
 - OpenMP: Avoid synchronization between threads when asking for more work to the master
 - MPI and OpenMP: Implement a “guided like” distribution of work by the master process. Bigger chunks at the beginning of the execution smaller at the end
 - MPI: “Kill” long computations at the end of the execution if there is no more work to do.
- After the implementation of the optimizations the execution in 16 nodes was 3x times faster



Performance Optimisation and Productivity 3

A Centre of Excellence in HPC

Contact:

 <https://www.pop-coe.eu>

 pop@bsc.es

 [@POP_HPC](#)

 [youtube.com/POPHPC](https://www.youtube.com/POPHPC)



This project has received funding from the European High Performance Computing Joint Undertaking (JU) under grant agreement No 101143931. The JU receives support from the European Union's Horizon Europe research and innovation programme and Spain, Germany, France, Portugal and the Czech Republic.

