



# POP Studies of Earth Sciences Codes

Jesús Labarta, BSC

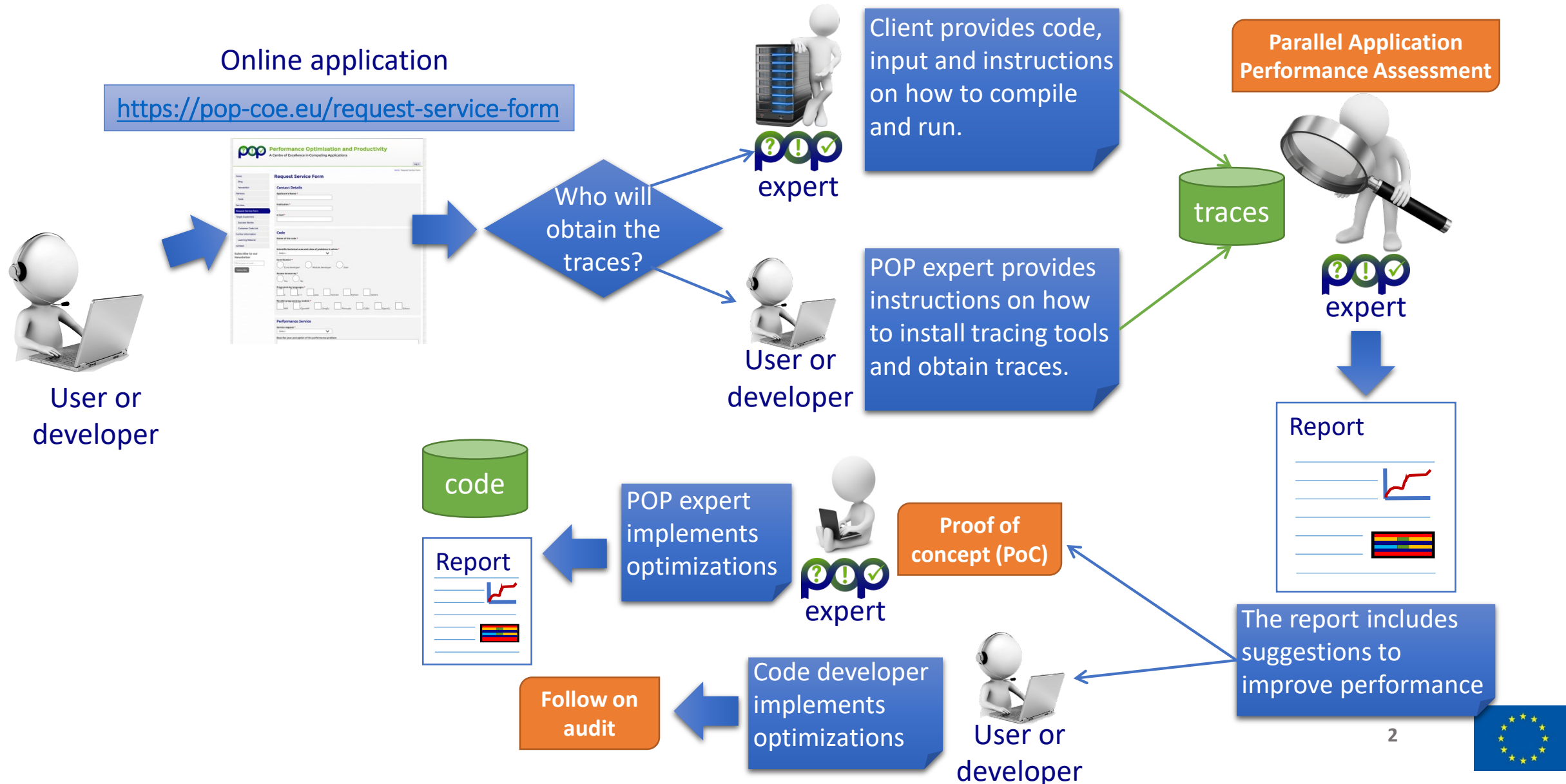
EU H2020 Centre of Excellence (CoE)



Grant Agreement No 824080

1 December 2018 – 30 November 2021

# POP services

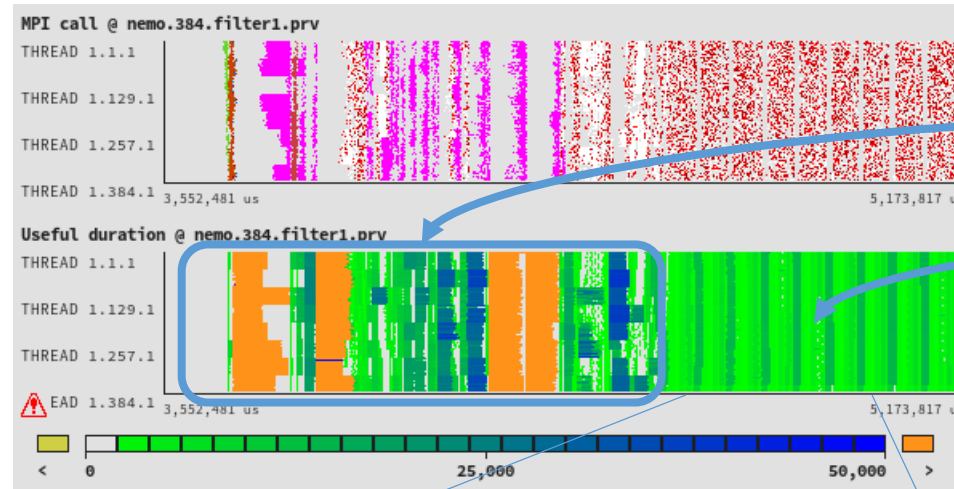


- Analyzed several codes ...
  - IFS, FVM
  - NEMO
  - MONARCH
  - ICON
- Towards Best Practices in
  - Performance Analysis Methodology and Tools
  - Programming Practices

# Structure

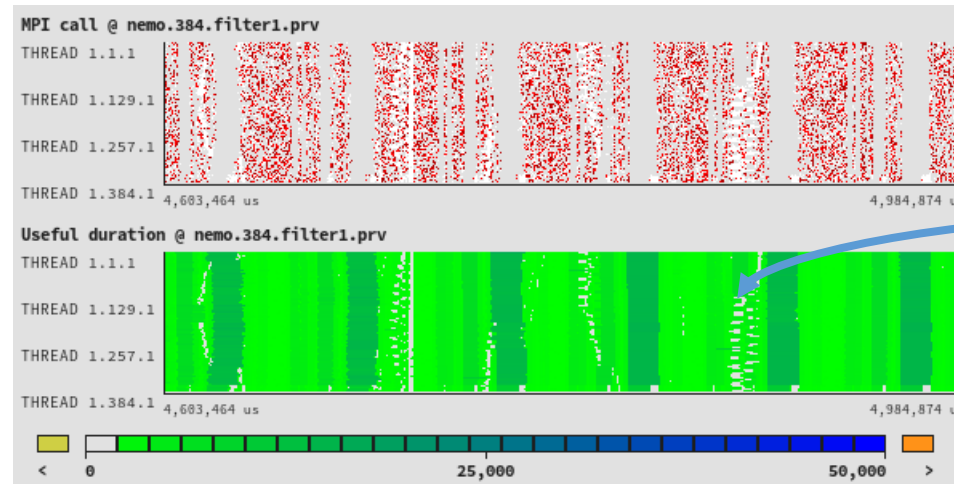


384 cores



Initialization

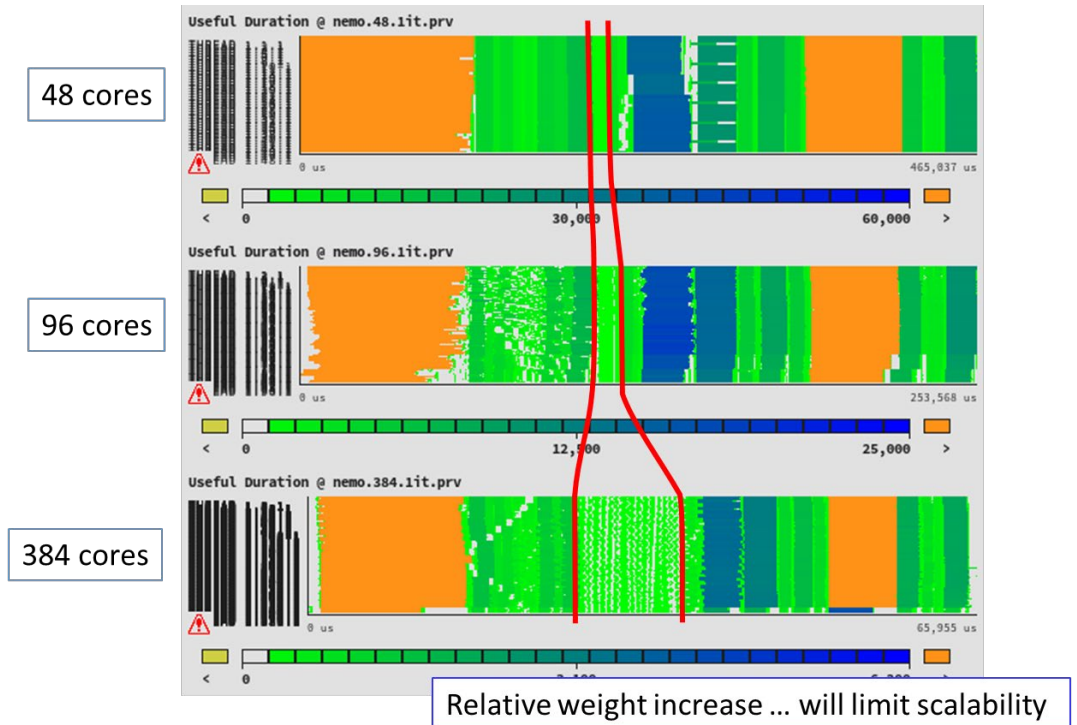
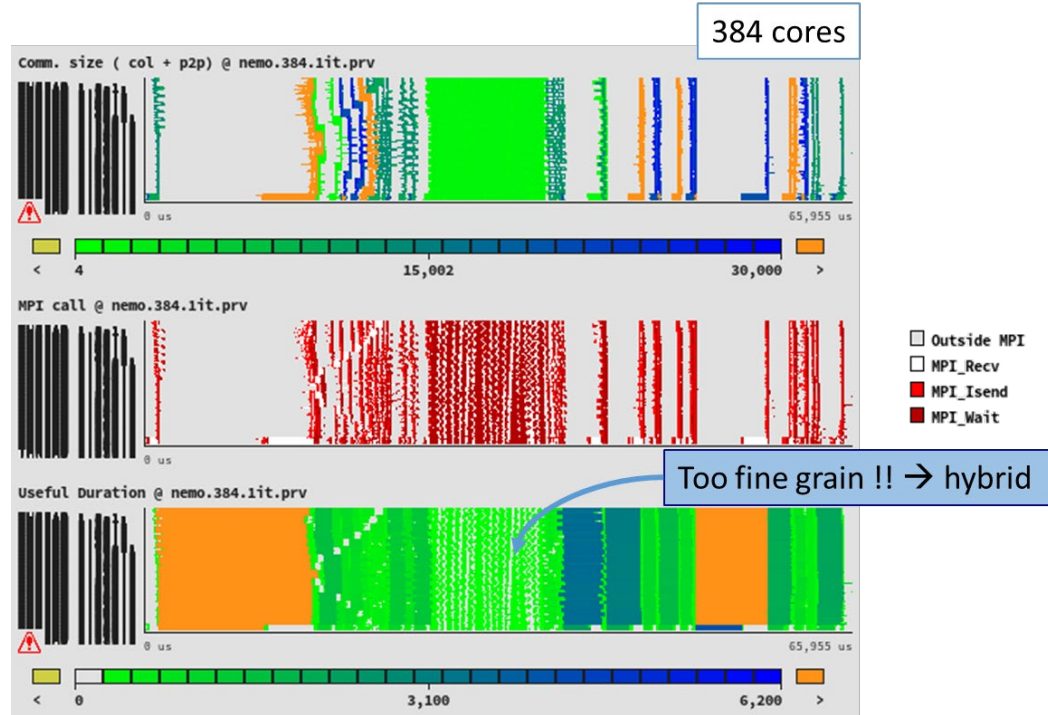
Iterative structure



perturbation



# Structure

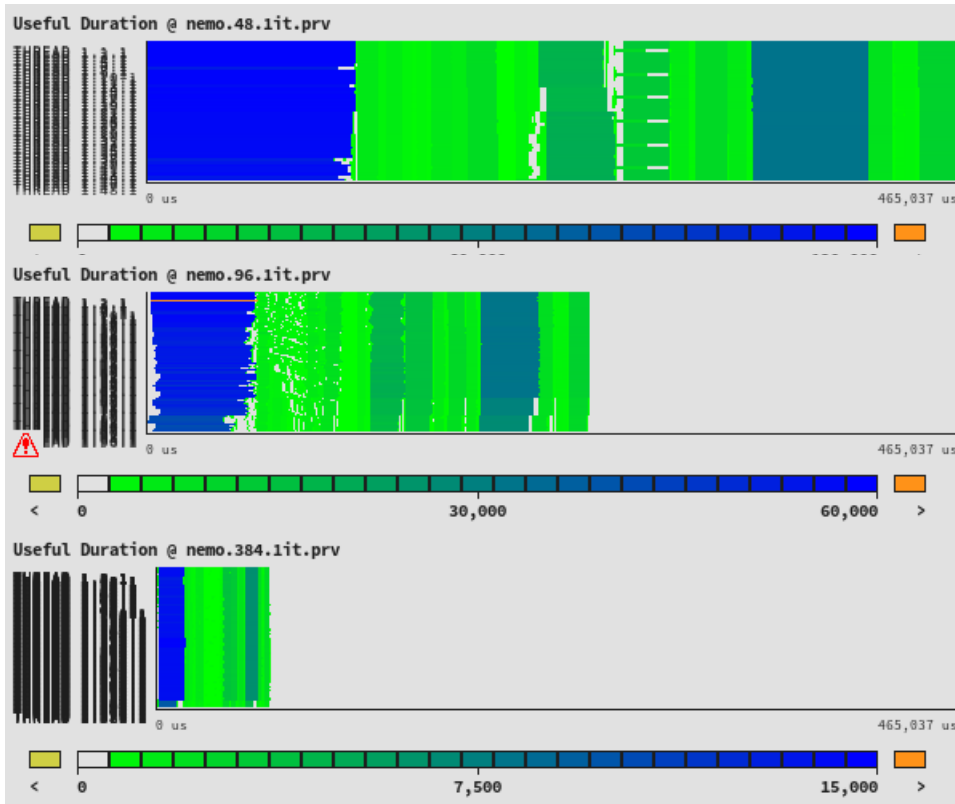


# Scaling



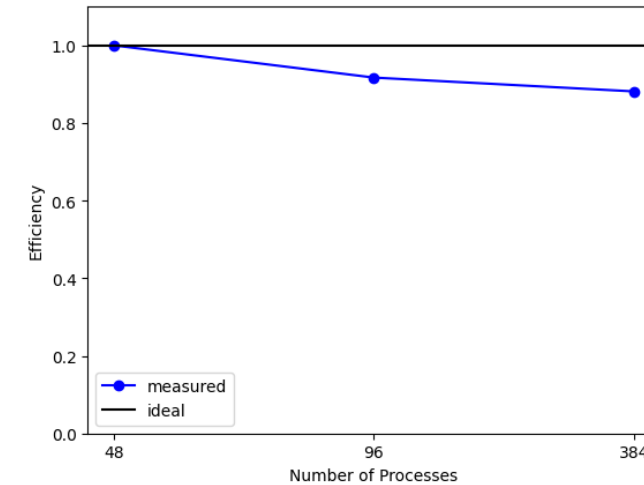
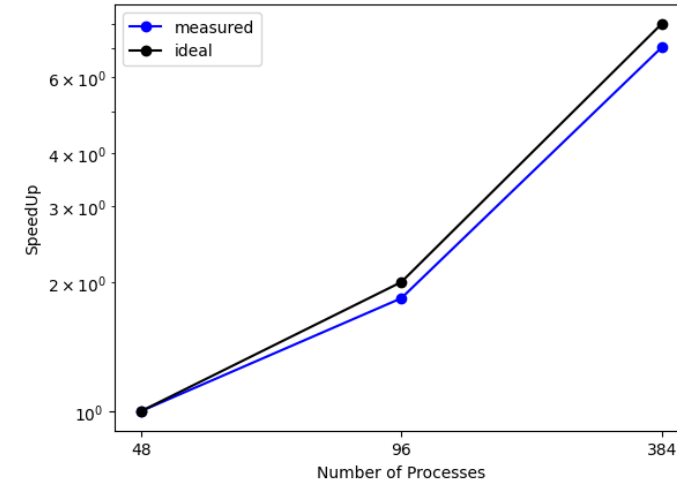
$T(48) = 465 \text{ ms}$

48

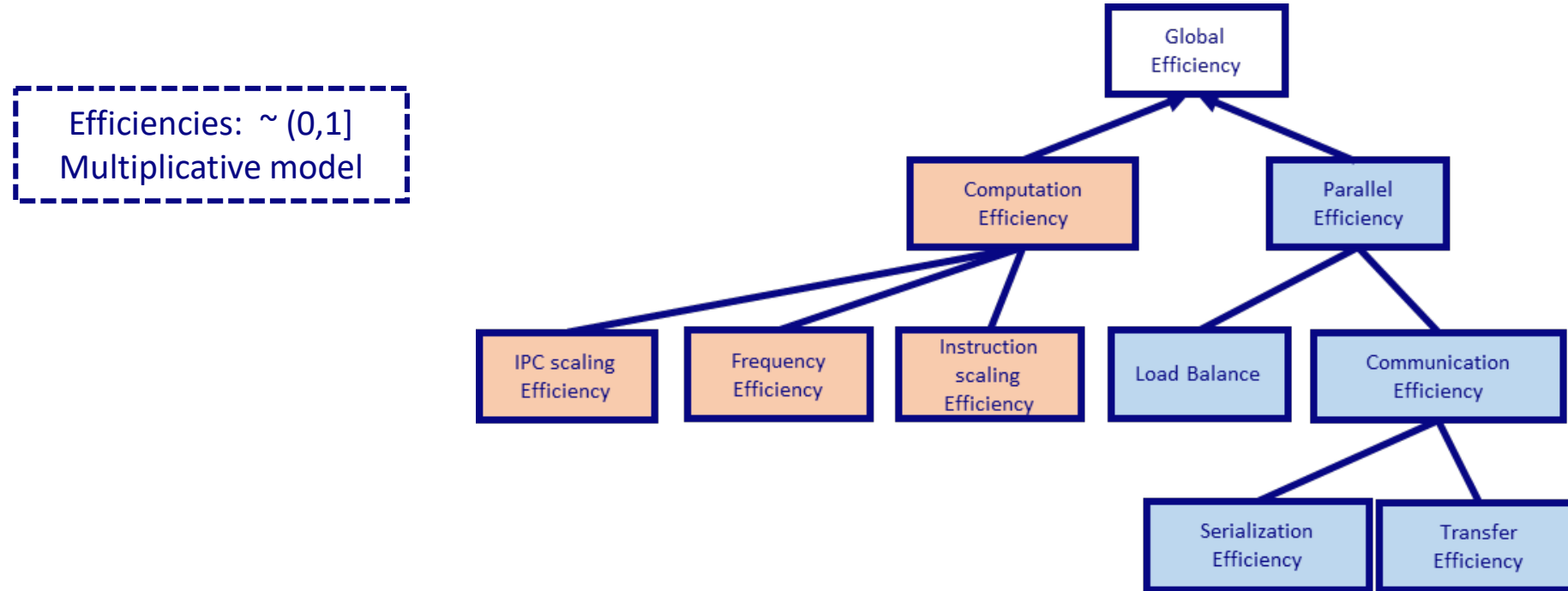


96

384



# Hierarchical Performance Model



$$CompEff = Ieff * IPCeff * Feff$$

$$\eta_{\parallel} = LB * \overset{CommEff}{\text{Ser} * Trf}$$



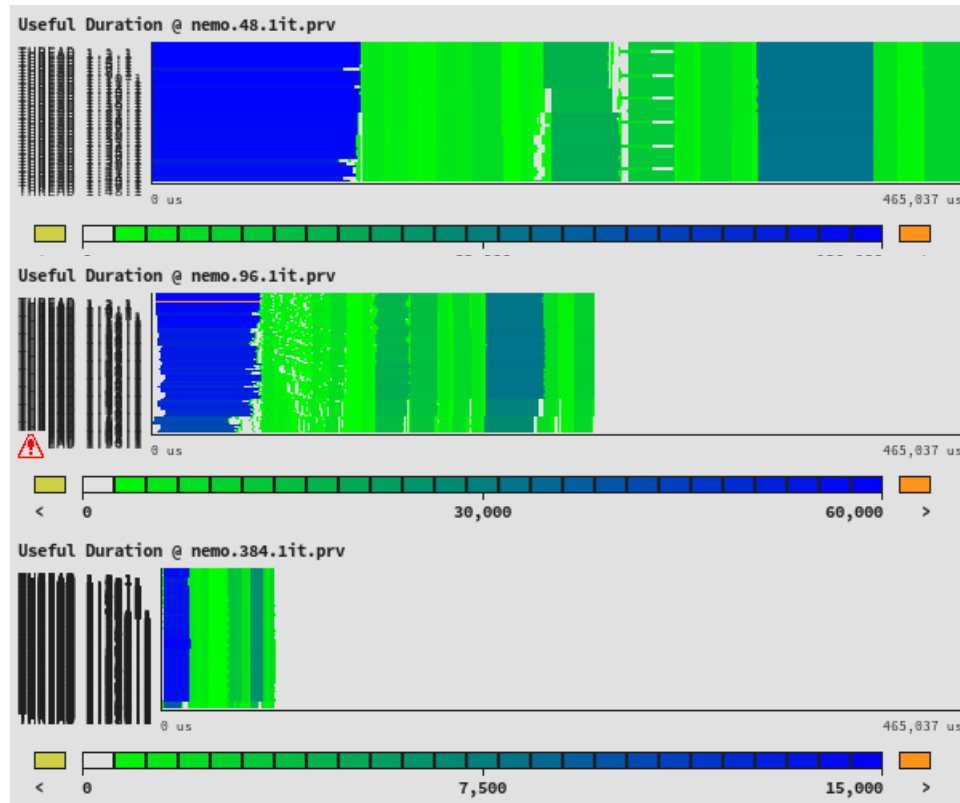
# Efficiency model



48

96

384



	48	96	384	
Global efficiency	94.16	86.35	82.99	
-- Parallel efficiency	94.16	84.11	74.80	
-- Load balance	98.67	95.34	95.77	
-- Communication efficiency	95.43	88.22	78.10	←
-- Serialization efficiency	97.86	93.05	89.61	
-- Transfer efficiency	97.52	94.81	87.15	
-- Computation scalability	100.00	102.66	110.95	
-- IPC scalability	100.00	115.98	182.39	←
-- Instruction scalability	100.00	91.37	62.62	←
-- Frequency scalability	100.00	96.87	97.13	

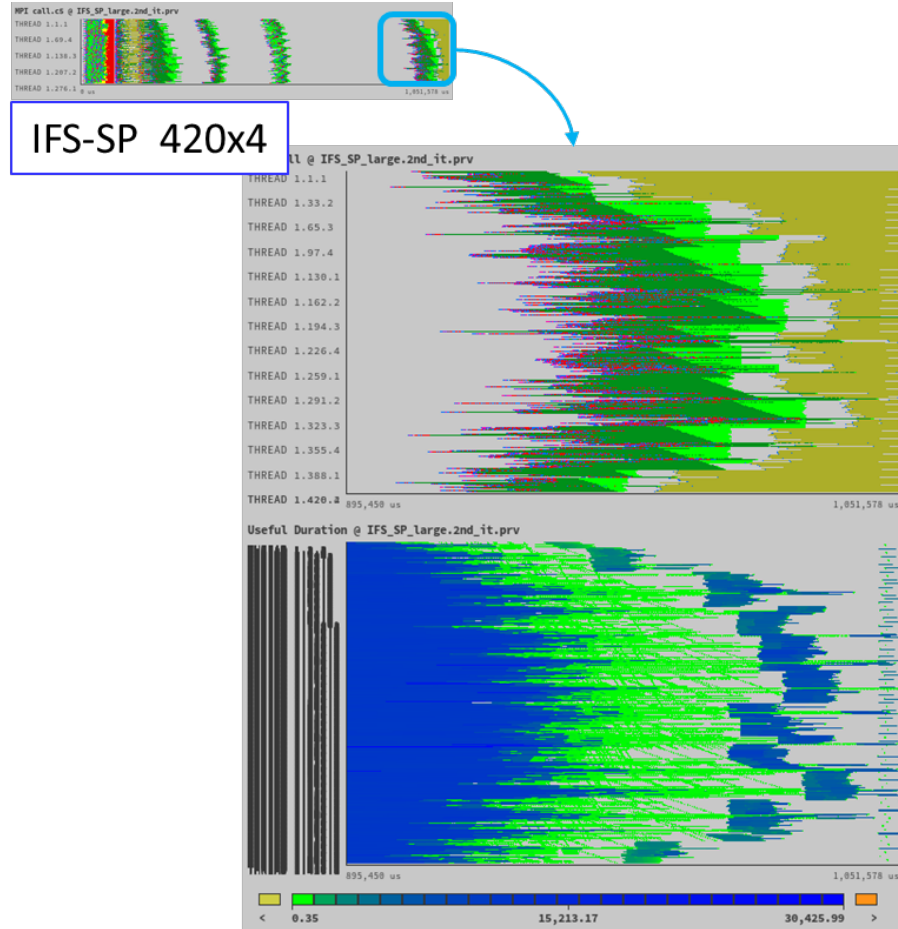
Avg Useful IPC(48) = 0.67

Avg Useful Frequency(48) = 2.061 GHz

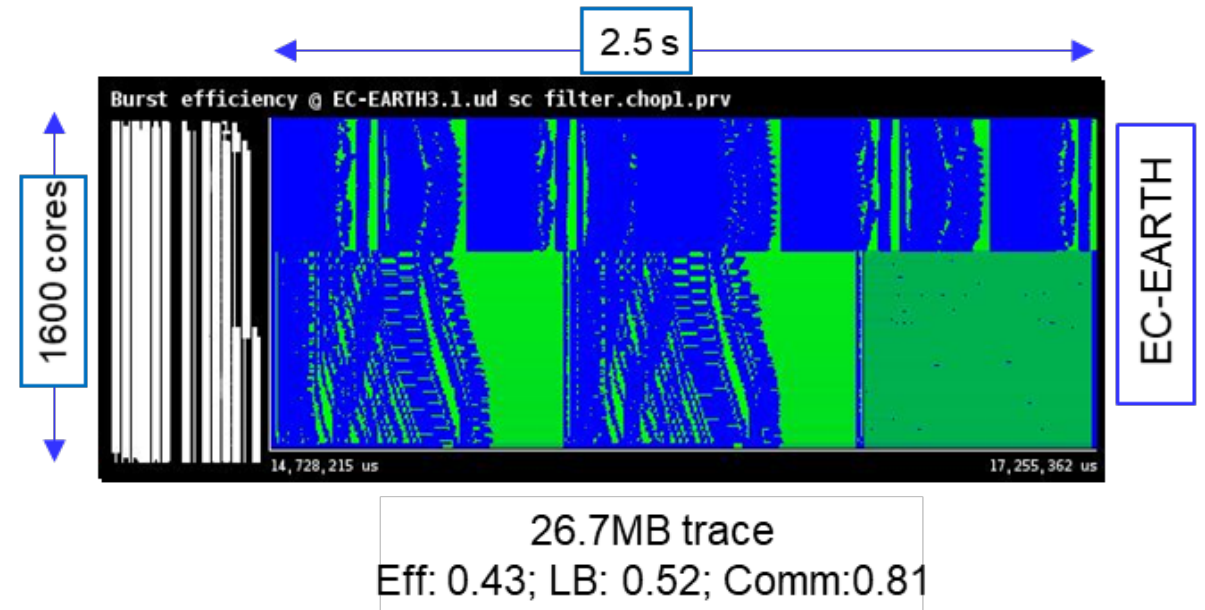




# Load Balance



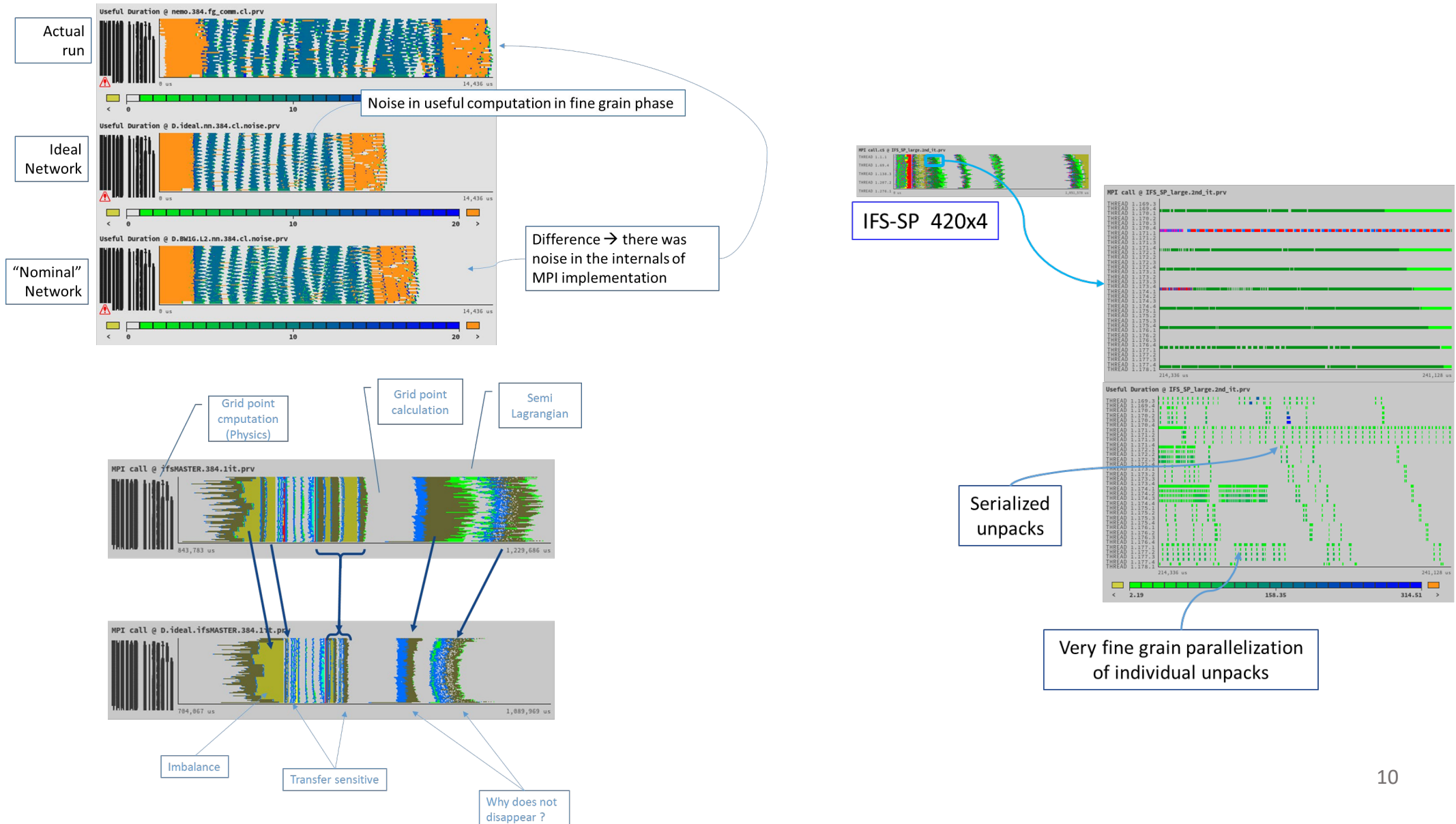
Within a model ...



... and coupled runs



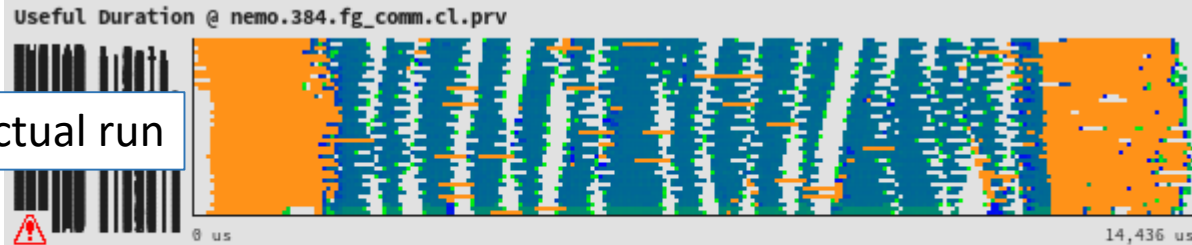
# Communication analysis



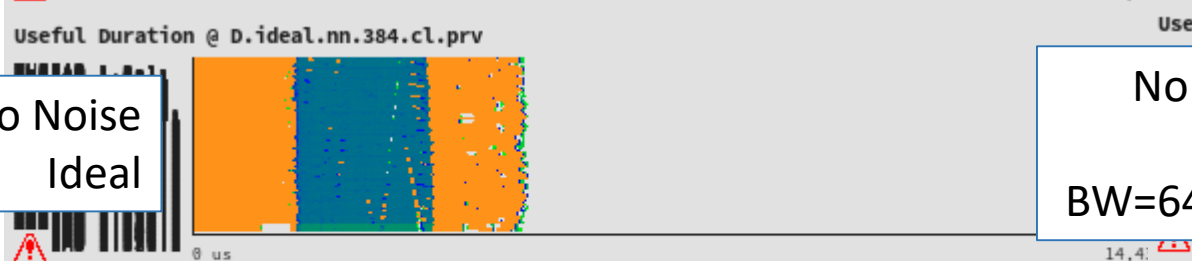
# What ifs



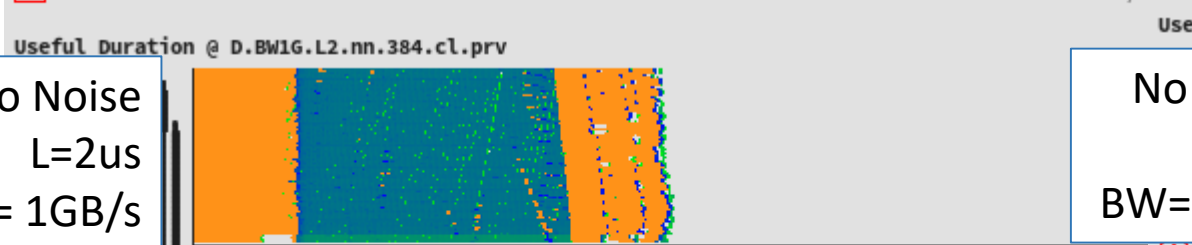
Actual run



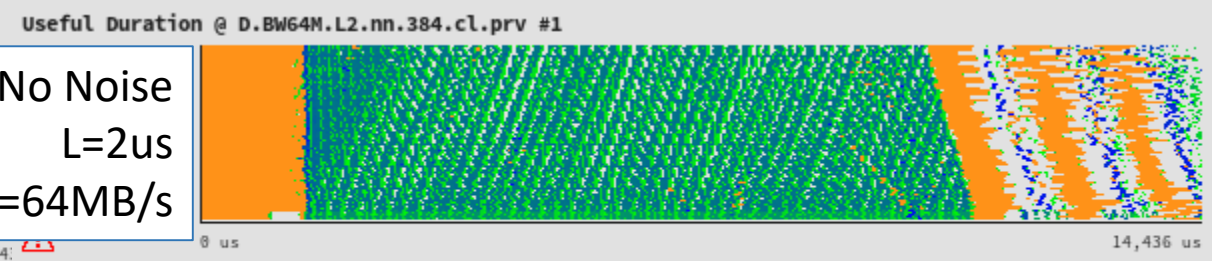
No Noise  
Ideal



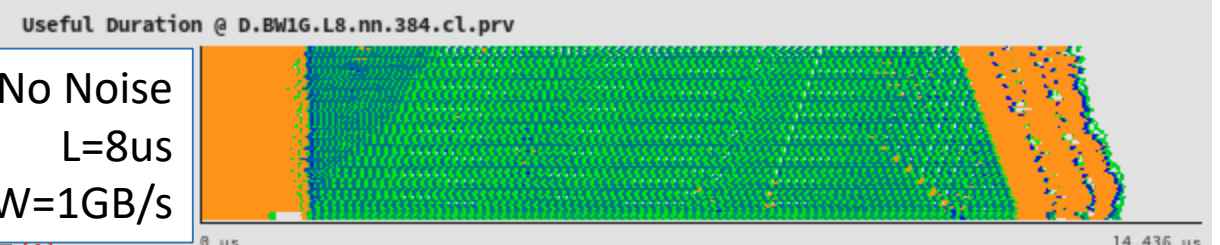
No Noise  
L=2us  
BW= 1GB/s



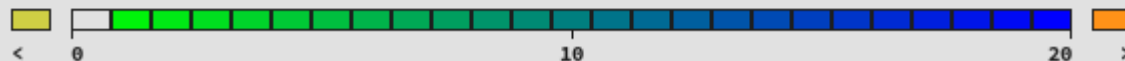
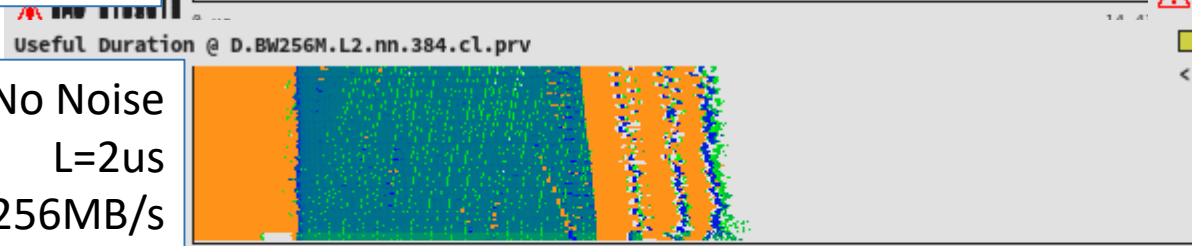
No Noise  
L=2us  
BW=64MB/s



No Noise  
L=8us  
BW=1GB/s



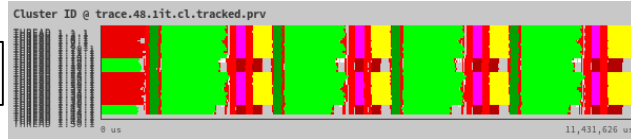
No Noise  
L=2us  
BW=256MB/s



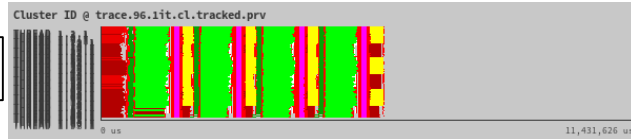
# MPI strong scaling



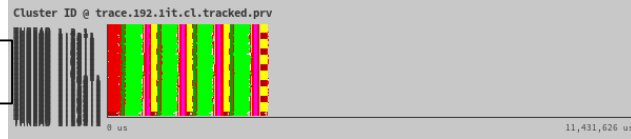
48



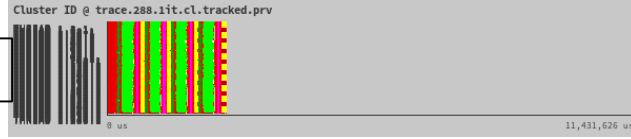
96



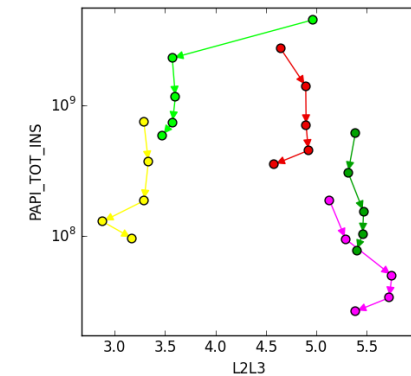
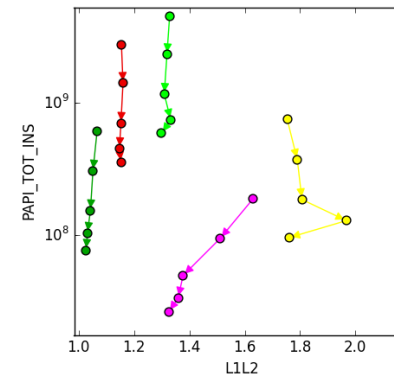
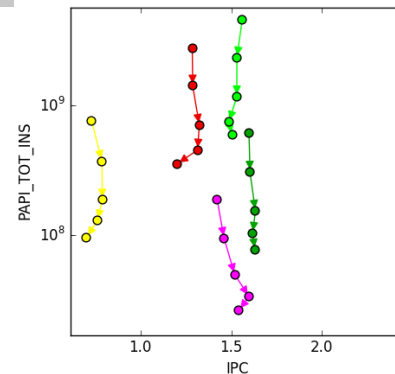
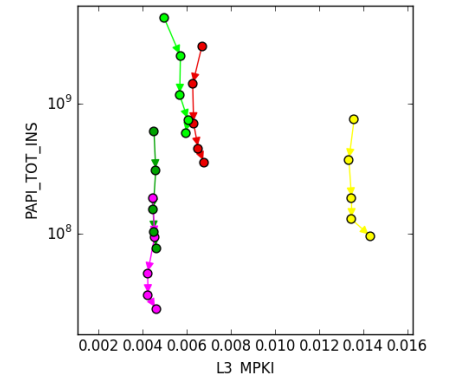
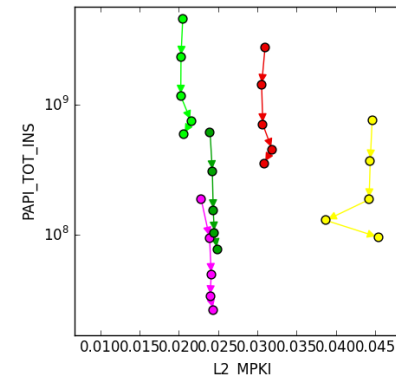
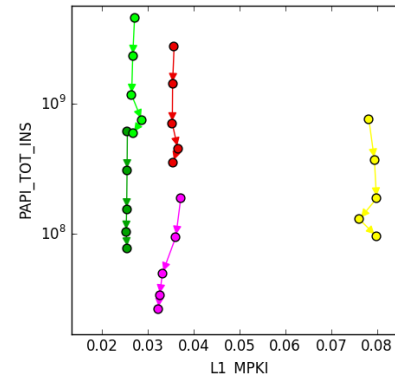
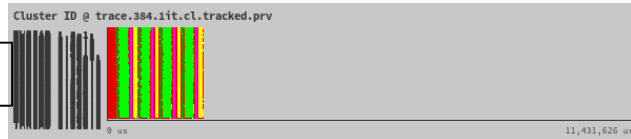
192



288



384



# MPI+OMP strong scaling



48x1

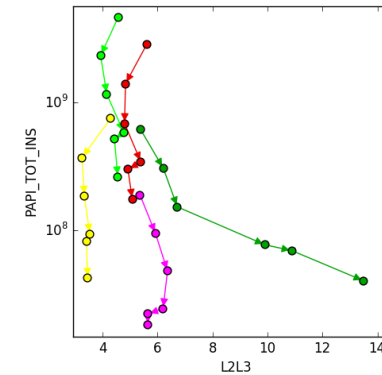
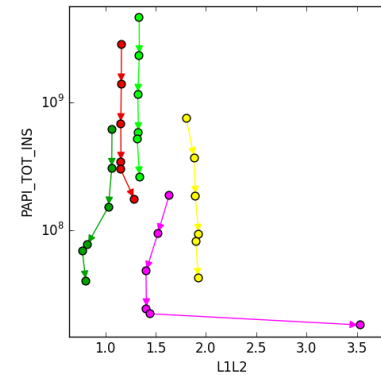
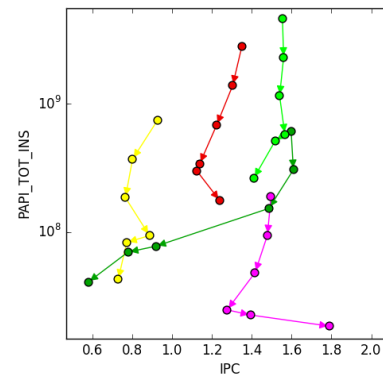
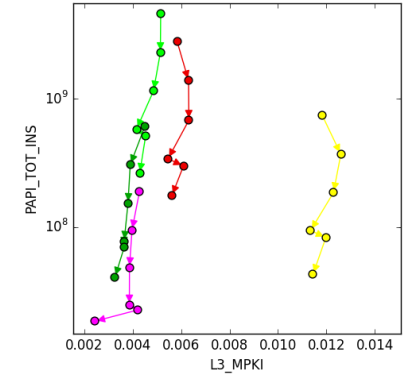
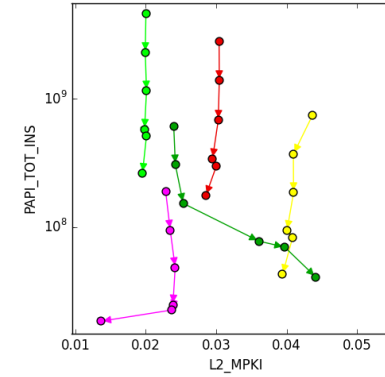
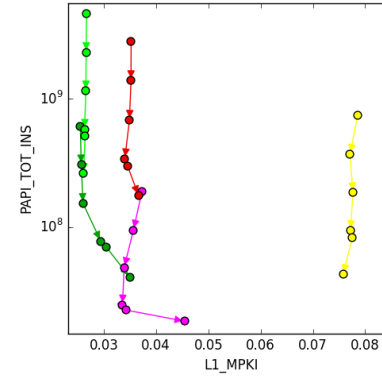
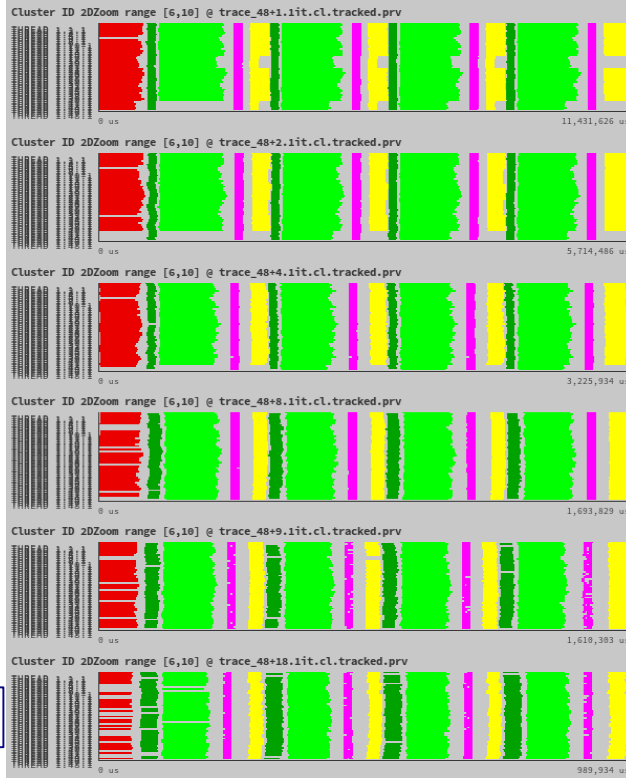
48x2

48x4

48x8

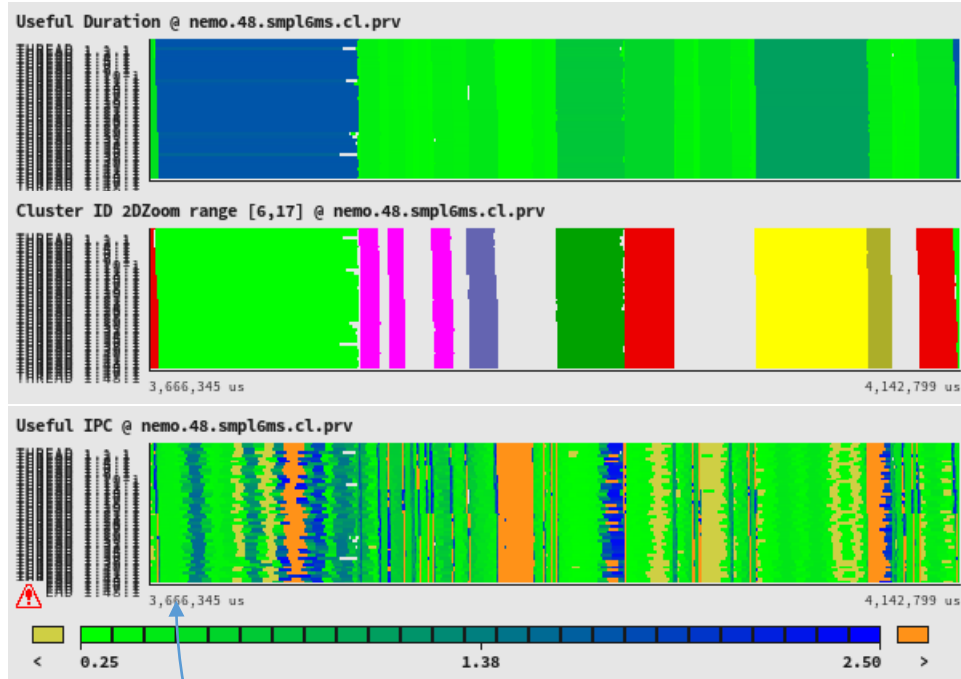
48x9

48x18

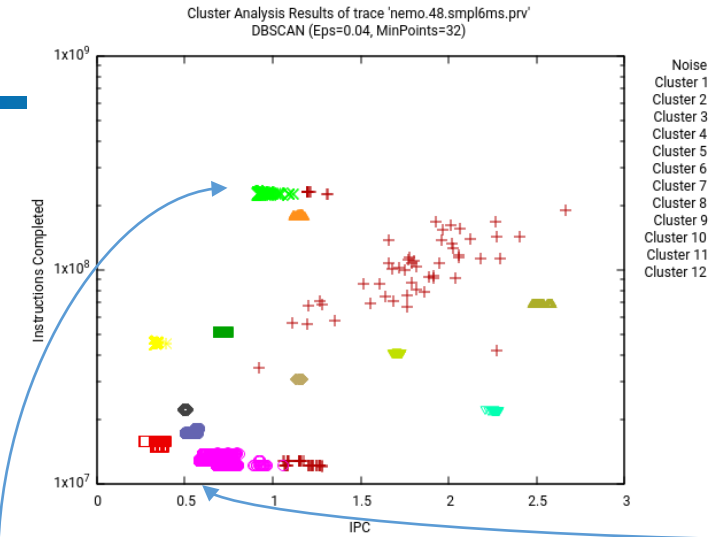




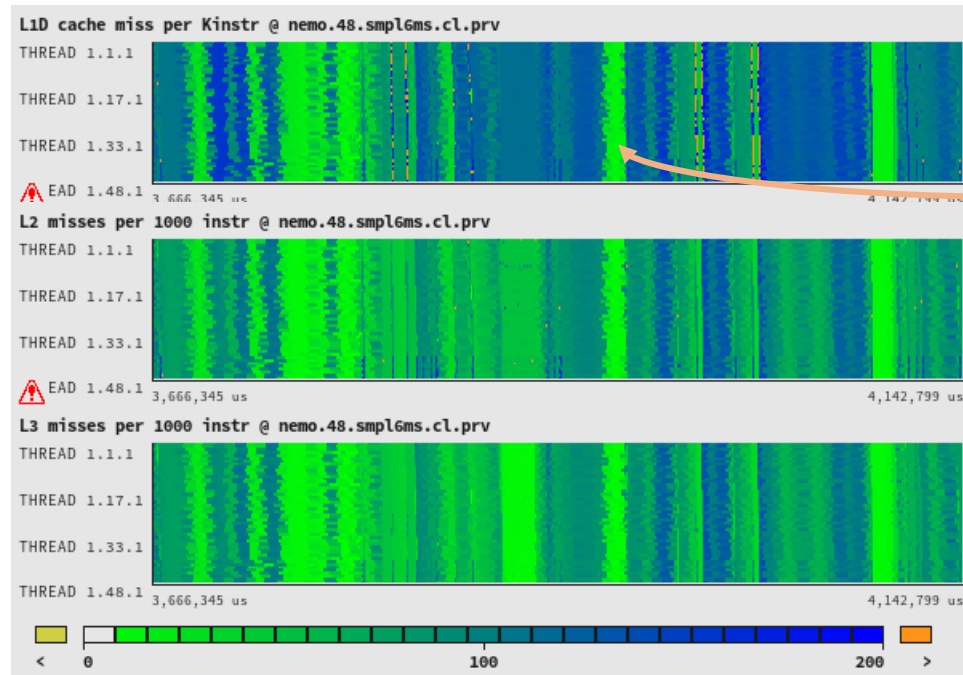
# Sampled traces



Very poor IPC sub regions within region of moderate average IPC



Regions with poor IPC

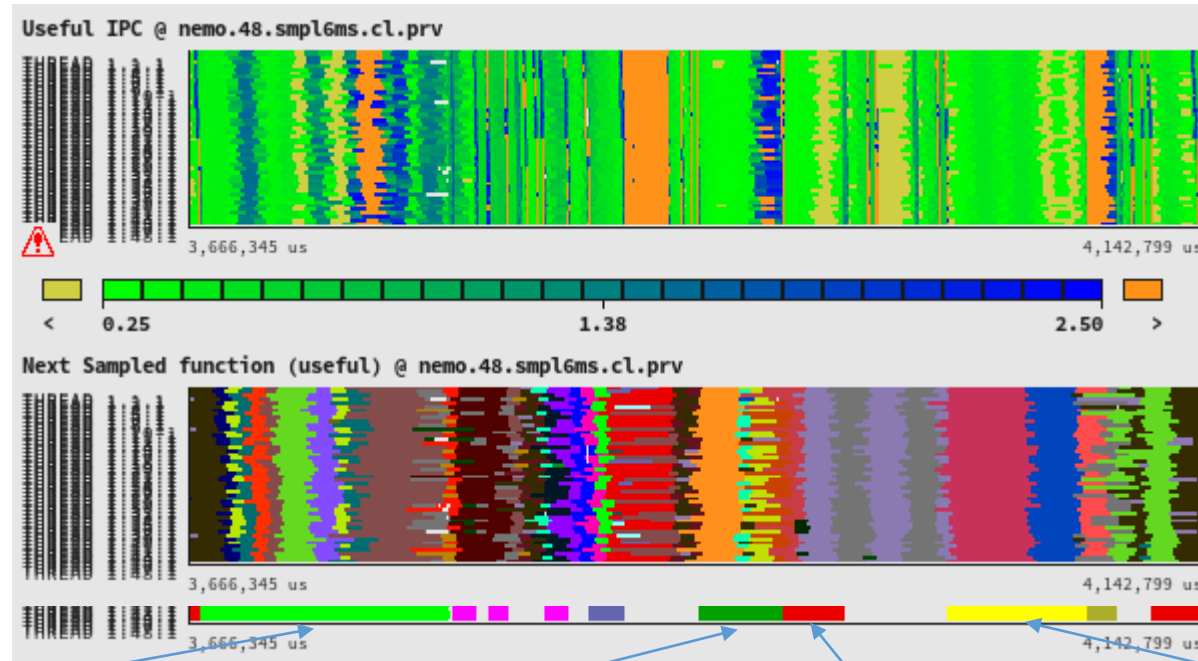


Cache miss ratios  
"explaining" IPC

Limited  
benefit of L3



# Sampled traces



End

- domvvl\_m..nterpol\_ [domvvl\_mp\_dom\_vvl\_interpol\_]
- eosbn2\_mp\_bn2\_
- zdfcke\_m..tke\_tke\_ [zdfcke\_mp\_tke\_tke\_]
- zdfcke\_m..tke\_avn\_ [zdfcke\_mp\_tke\_avn\_]
- zdfiwm\_m..zdf\_iwm\_ [zdfiwm\_mp\_zdf\_iwm\_]

End

- dynzdf\_m..dyn\_zdf\_ [dynzdf\_mp\_dyn\_zdf\_]
- sshwzv\_mp\_wzv\_
- traqsr\_m..tra\_qsr\_ [traqsr\_mp\_tra\_qsr\_]
- traadv\_m..tra\_adv\_ [traadv\_mp\_tra\_adv\_]

End

- lbcInk\_m..\_3d\_ptr\_ [lbcInk\_mp\_mpp\_lnk\_3d\_ptr\_]
- ldftra\_m..eiv\_trp\_ [ldftra\_mp\_ldf\_eiv\_trp\_]
- traadv\_f..adv\_fct\_ [traadv\_fct\_mp\_tra\_adv\_fct\_]
- traadv\_f..nonosc\_ [traadv\_fct\_mp\_nonosc\_]

End

- lbcInk\_m..\_3d\_ptr\_ [lbcInk\_mp\_mpp\_lnk\_3d\_ptr\_]
- traadv\_f..adv\_fct\_ [traadv\_fct\_mp\_tra\_adv\_fct\_]
- traldf\_i..ldf\_iso\_ [traldf\_iso\_mp\_tra\_ldf\_iso\_]
- trazdf\_m..zdf\_imp\_ [trazdf\_mp\_tra\_zdf\_imp\_]
- tra\_nxt\_vvl



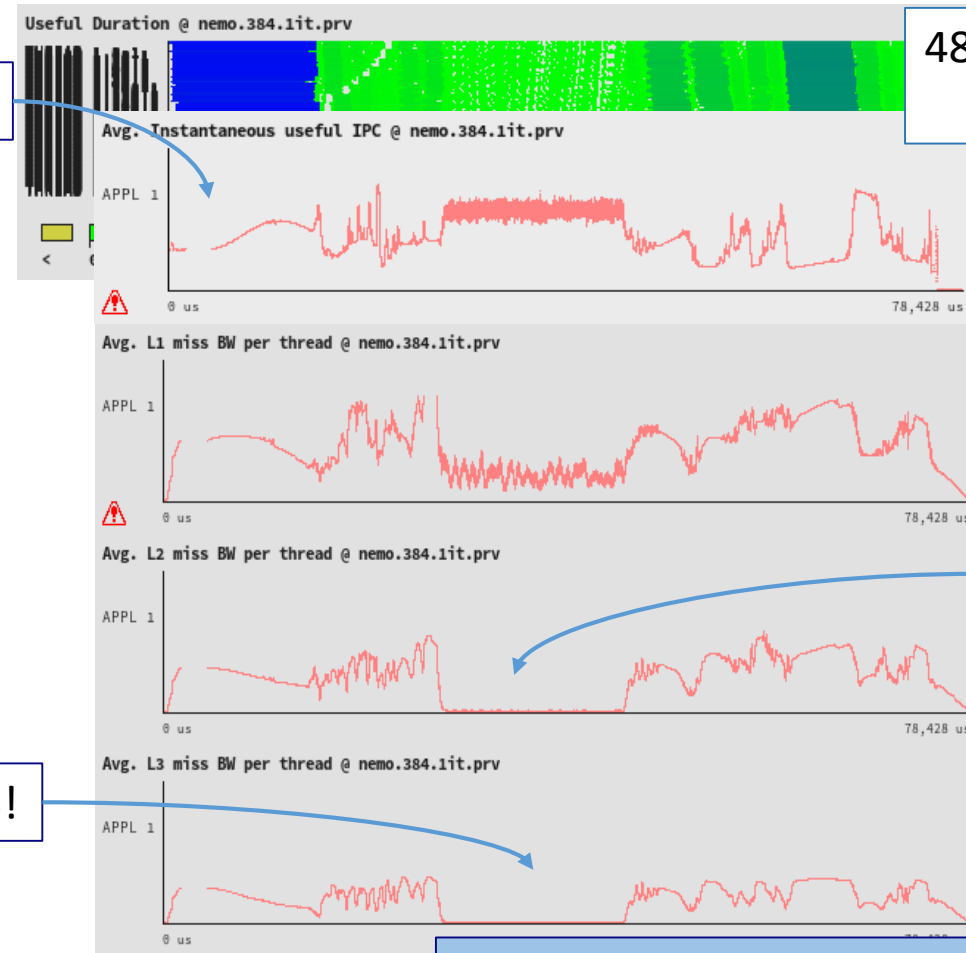
# Aggregated time vaying behavior



Loss of substructure?

Data fits in L3!

Node BW limit



48 vs 384 processes  
IPC - BW

Data fits in L2!

Can refactor code to improve L3

Co-design: no need of L3 ?

use? Blocking?





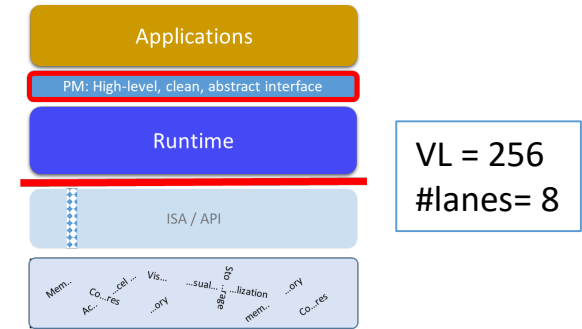
- Observations ...
  - Granularities
  - Instruction scaling
  - IPCs and Memory bandwidth
  - L3 use
  - Pack-unpack
  - False sharing
- Recommendations: Asynchrony and overlap
  - Tasks
- ... and co-design
  - RISC-V vector
  - OpenMP
    - Features: Free agents, precompiled task graphs
    - Libraries: DLB, TAMPI, TALP

# RISC-V & Long vectors



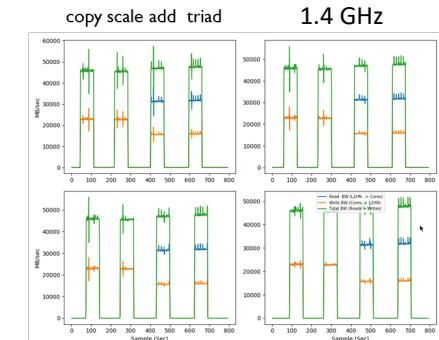
- Raise ISA semantic level
  - Vector instructions == tasks
  - “less words, more work”
  - The importance of ISA
- Parallelism
  - Decouple Front end – back end
    - Less pressure, throughput orientation
  - OoO execution

- Osmotic membrane
  - Convey access pattern semantics to the architecture.
  - Potential to optimize memory throughput.



```
happy@epac$ axpy 1024
Running AXPY Scalar with 1024 array elements
init time: 45060 cycles
axpy scalar reference time: 23555 cycles
done
Result ok !!!
happy@epac$ vaxpy 1024
Running AXPY Vector with 1024 array elements
init time: 45043 cycles
axpy vector time: 932 cycles
done
Result ok !!!
happy@epac$
```

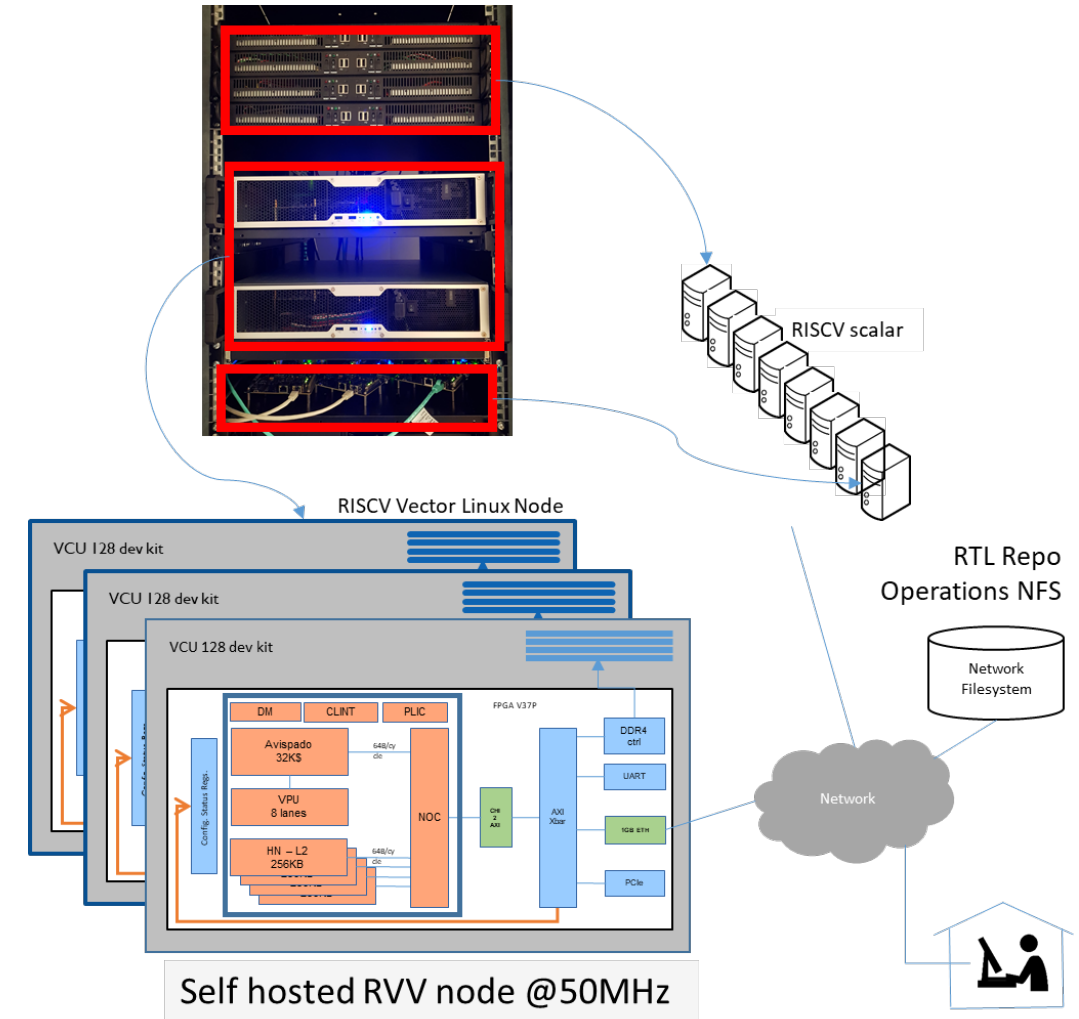
~25x  
while only 8x FPUs  
→ Long vectors !!  
→ Memory Bandwidth



# EPI SDV ecosystem

- RISC-V cluster
  - Commercially available RISC-V platforms
  - Porting and configuring HPC software stack and increase productivity (e.g., SLURM, MPI, OpenMP, BSC tools, SDV1.2)
- SDV: RVV @ FPGA nodes
  - CI Infrastructure: Validation at “scale”
  - Software development and co-design steering
    - Test real “complex” codes @ real RTL
    - EPAC1.5 RTL improvement
    - Give to EPI partners and interested users easy access to the latest EPAC technology
    - Two step procedure

Contact : [filippo.mantovani@bsc.es](mailto:filippo.mantovani@bsc.es)

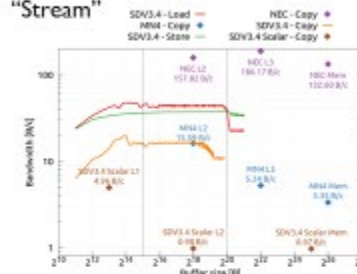


# EPI SDV ecosystem

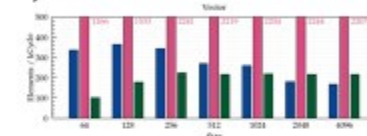
## SDV – VECTOR PERFORMANCE

- EPAC ...
- ... vs. state of the art  
(Xeon, NEC, A64FX, FU740, ...)

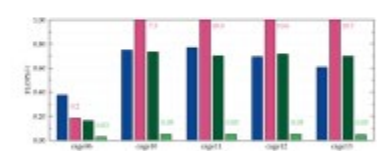
### "Stream"



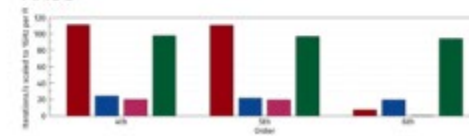
### Jacobi 2D



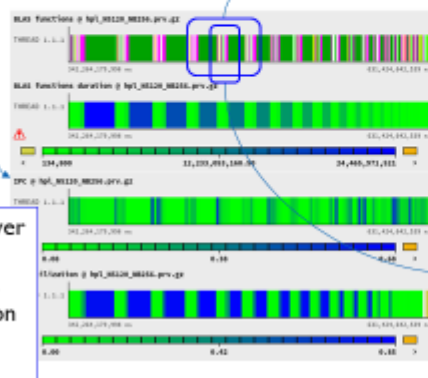
### SpMV



### HACC



## LOD IN BEHAVIOR ANALYSIS

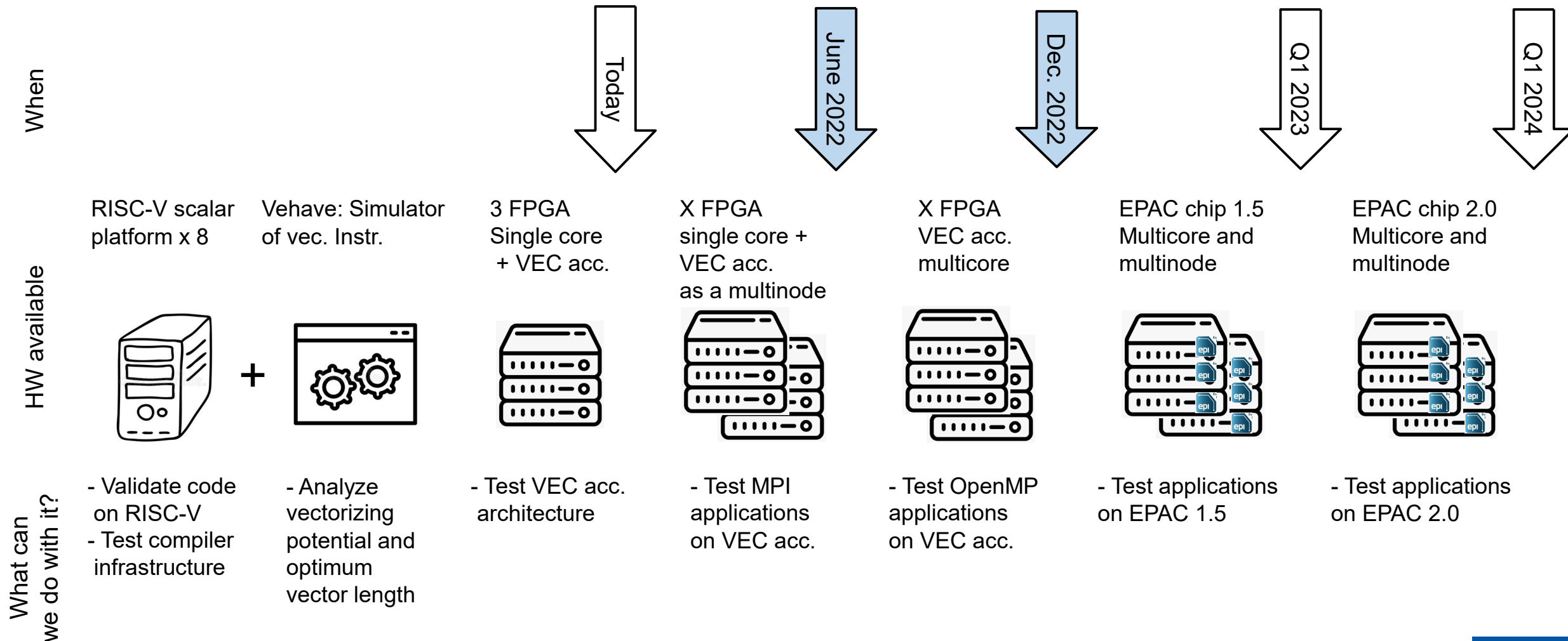


Extrae + Paraver  
Standard HPC  
instrumentation  
analytics and  
visualization

Vehave + MUSA  
Detailed analysis ...  
... of flexible what-ifs  
at ISA/μArch level

SDV@FPGA  
Detailed analysis ...  
... of actual RTL ...  
... at "large" scale

# EPI SDV roadmap





# Performance Optimisation and Productivity

A Centre of Excellence in HPC

Contact:

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

<mailto:pop@bsc.es>

 @POP\_HPC

