



Nemo analyses POP_AR_078

Jesus Labarta, Marta Garcia, Joan Vinyals (BSC)



EU H2020 Center of Excellence (CoE)

November 4th 2020

The code



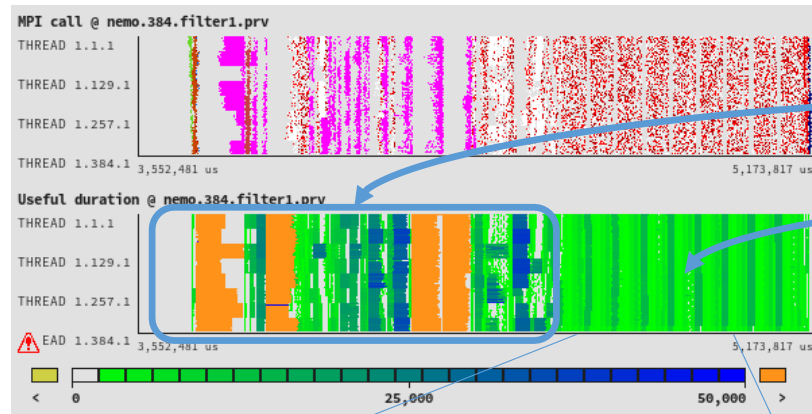
- Code
 - NEMO version 4.0.2
 - MPI
- Problem:
 - Bench ORCA 1 like
- Traces:
 - Run @ MareNostrum, filling node (48 processes/core)
 - 48, 96 and 384 processes



Structure

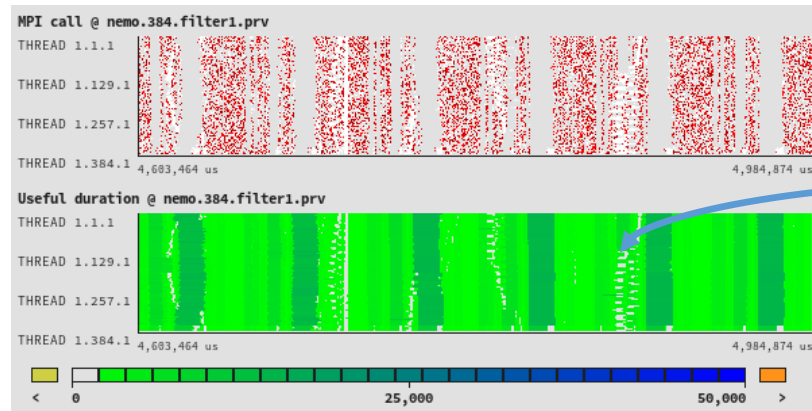


384 cores



Initialization

Iterative structure



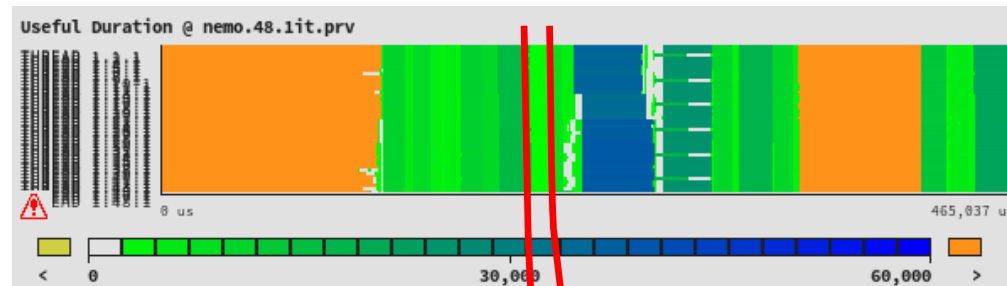
perturbation



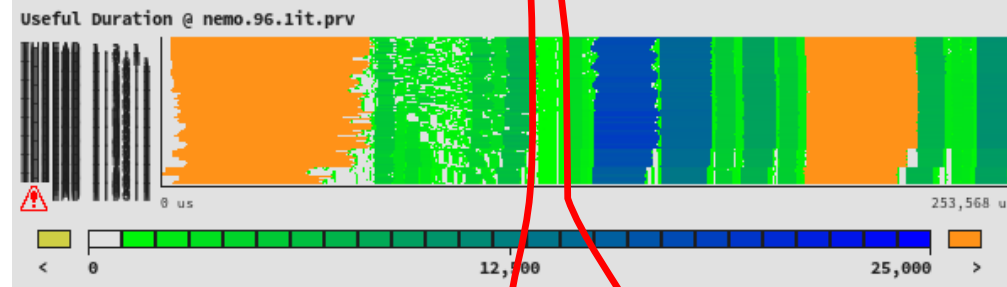
Structure



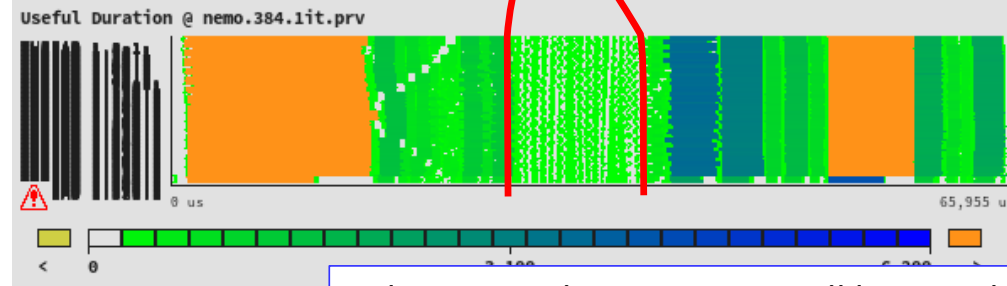
48 cores



96 cores



384 cores



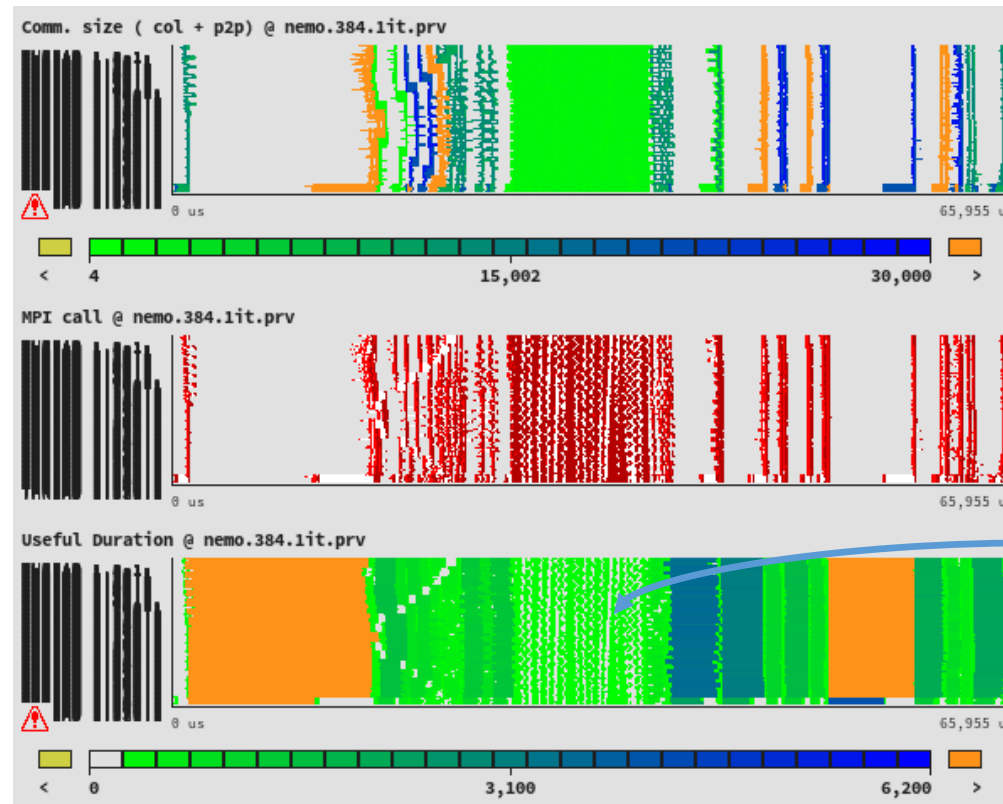
Relative weight increase ... will limit scalability



Structure



384 cores



Too fine grain !! → hybrid

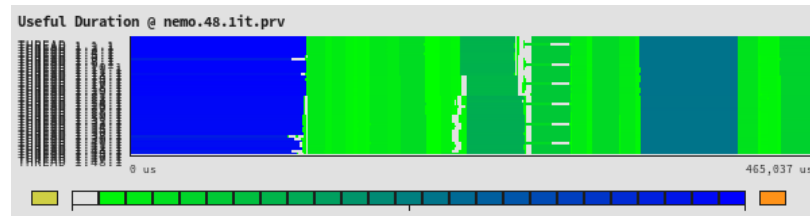


Scaling

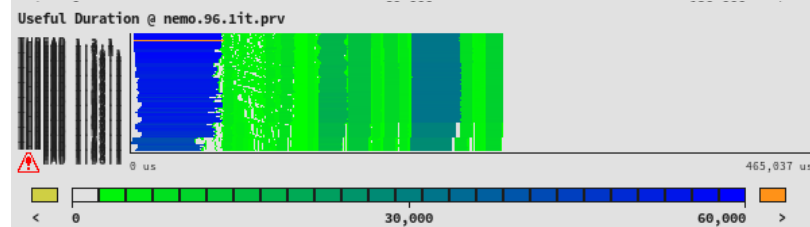


$T(48) = 465 \text{ ms}$

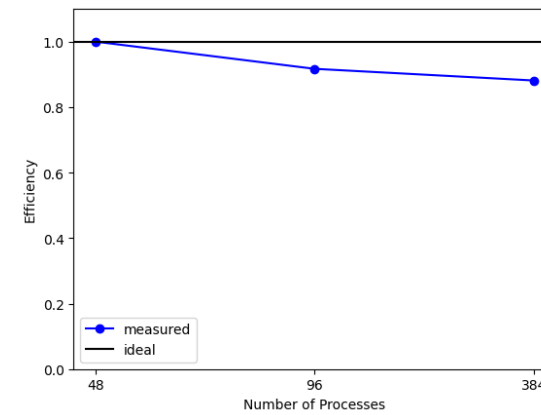
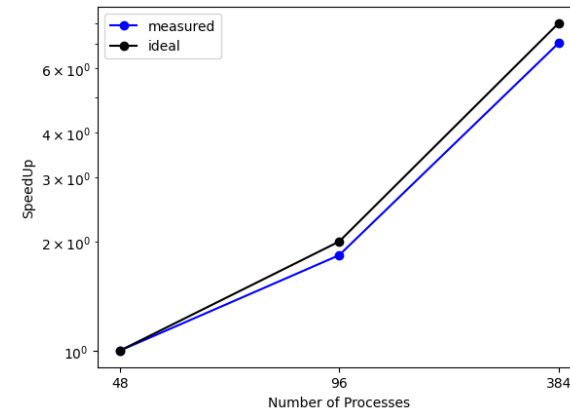
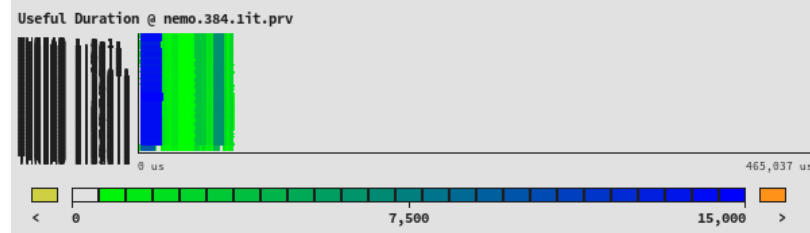
48



96



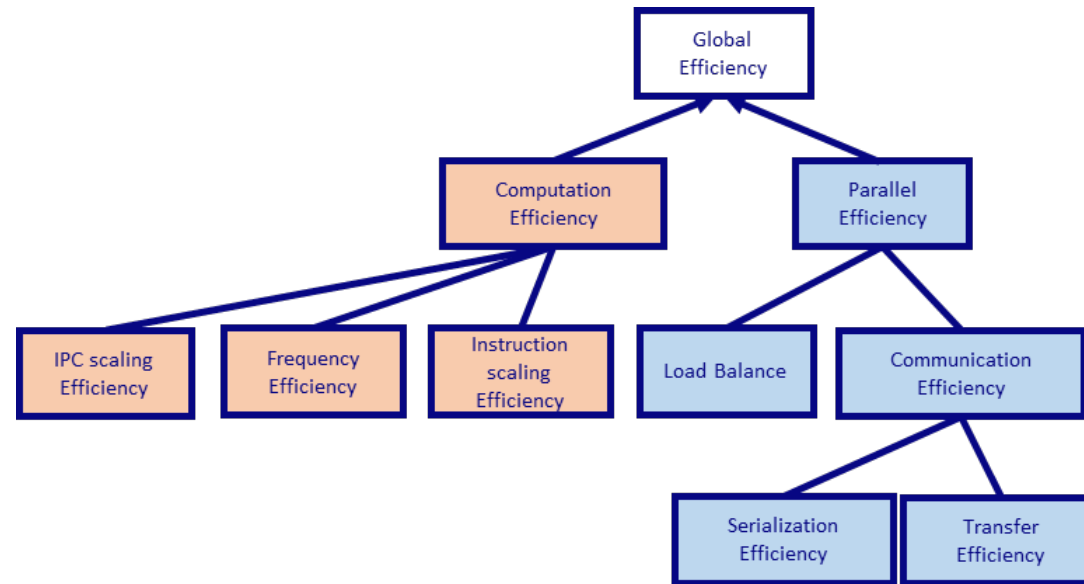
384



Hierarchical Performance Model



Efficiencies: $\sim (0,1]$
Multiplicative model



$$CompEff = Ieff * IPCeff * Feff$$

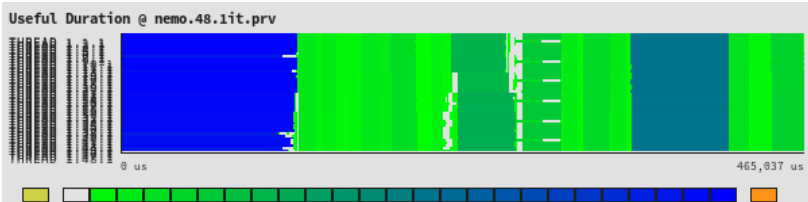
$$\eta_{||} = LB * \overset{CommEff}{Ser} * Trf$$



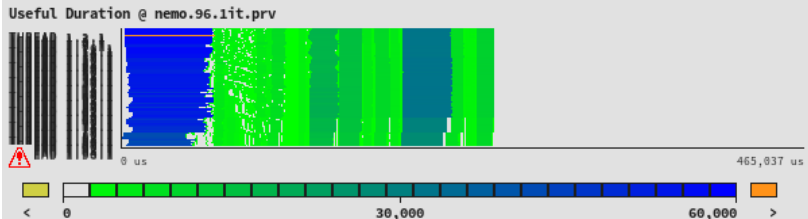
Efficiency model



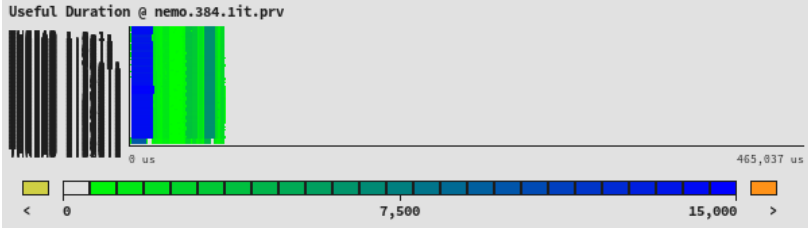
48



96



384



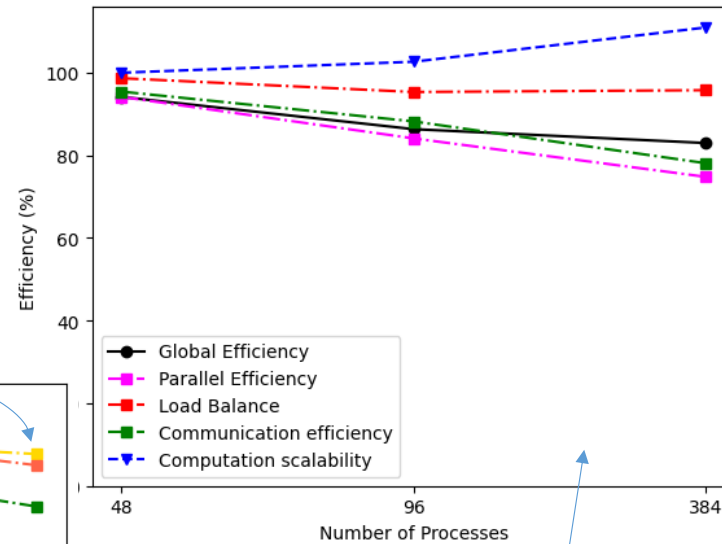
	48	96	384	
Global efficiency	94.16	86.35	82.99	
-- Parallel efficiency	94.16	84.11	74.80	
-- Load balance	98.67	95.34	95.77	
-- Communication efficiency	95.43	88.22	78.10	
-- Serialization efficiency	97.86	93.05	89.61	
-- Transfer efficiency	97.52	94.81	87.15	
-- Computation scalability	100.00	102.66	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

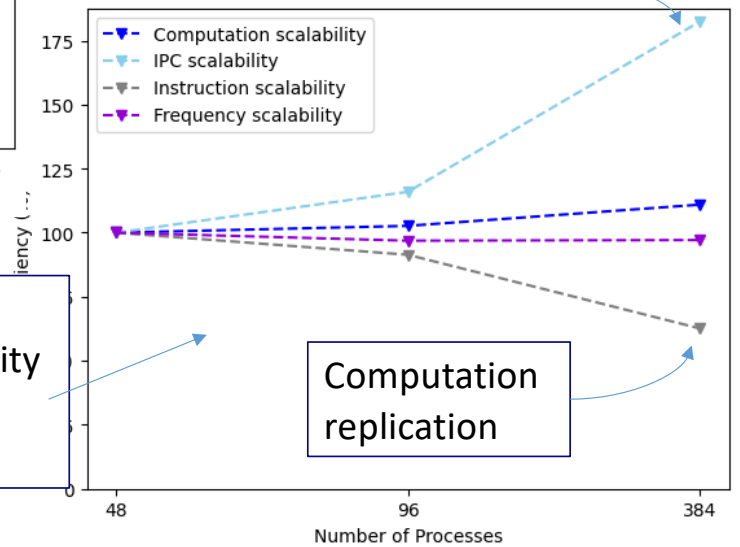
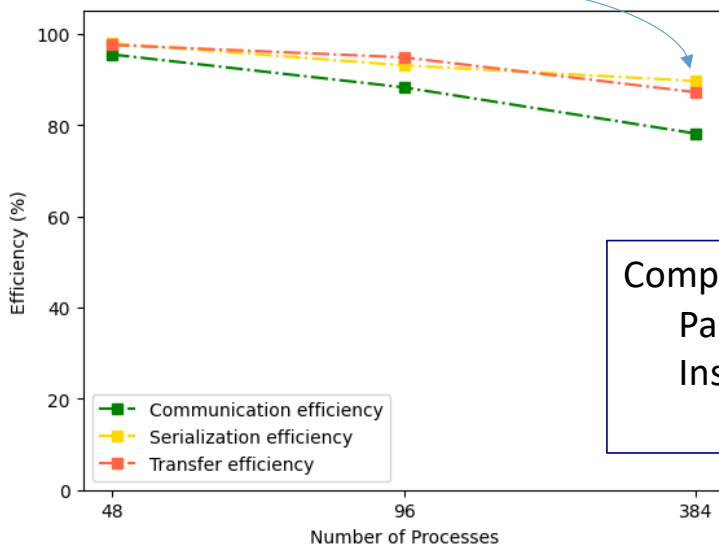


Efficiency model



Very good IPC scaling
 → probably very bad starting value
 → will not go on forever

Transfer and Serialization
 Transfer scaling worse



Compensation effects
 Parallel efficiency ↔ Computation scalability
 Instruction scalability (replication) ↔ IPC

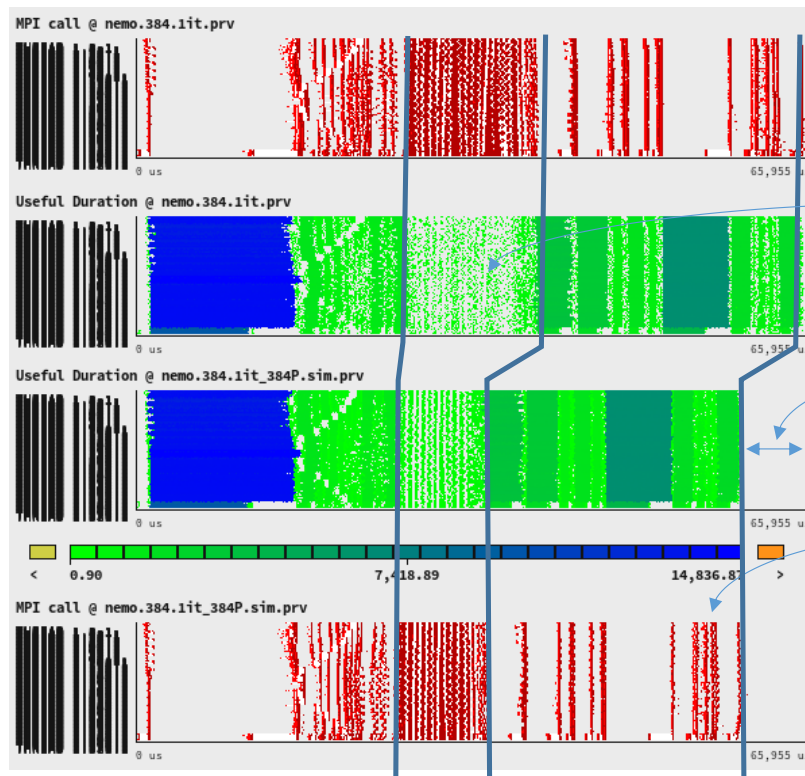
Computation
 replication



Communication analysis



Actual run



Fine granularity region sensitive to Interconnect latency and bandwidth

Transfer

Ideal Network

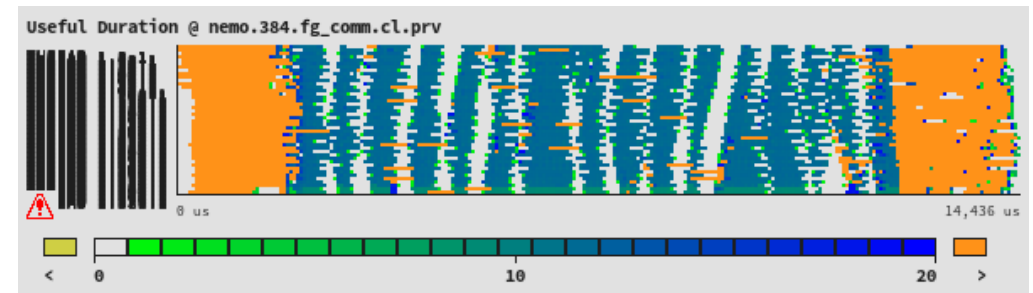
MPI does not disappear ... Elapsed time for the region does not shrink



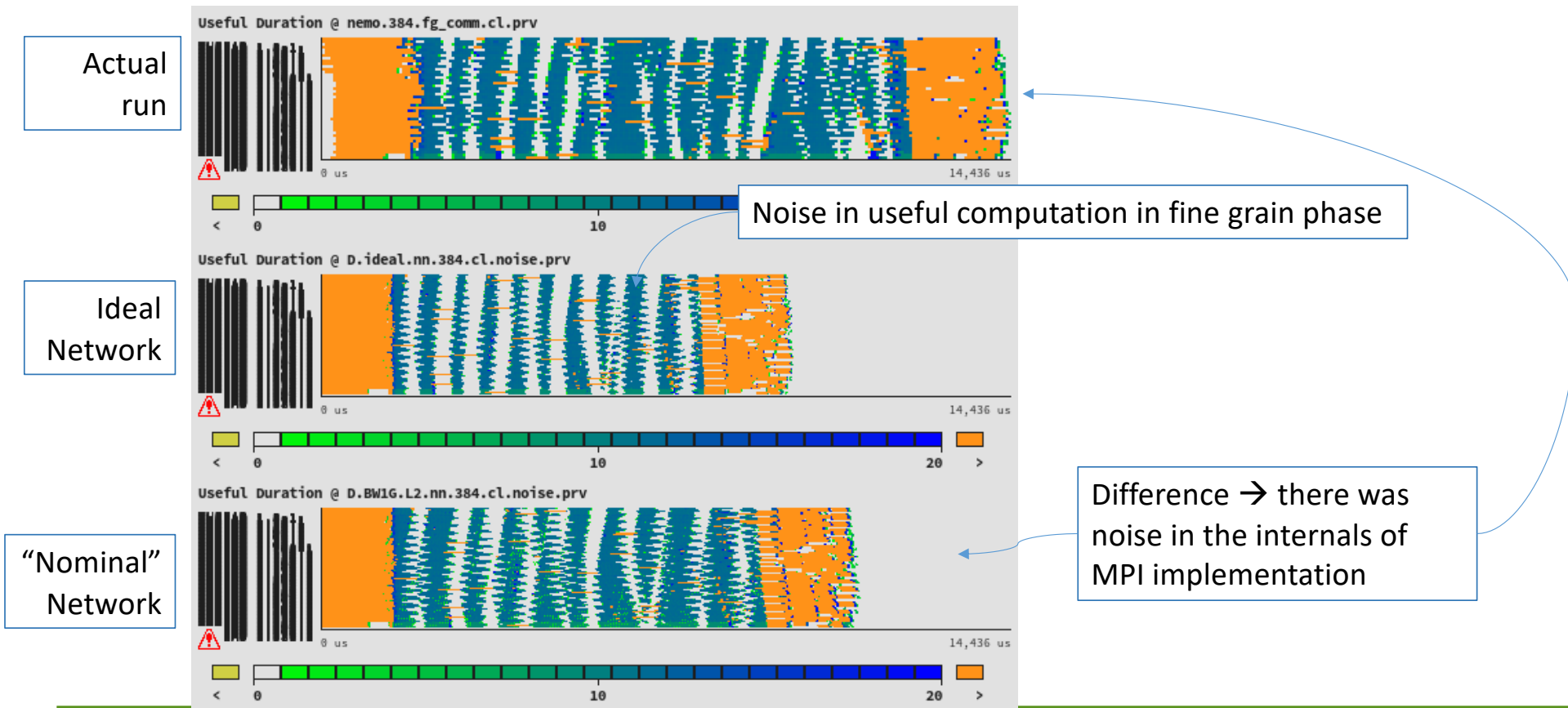
Communication analysis



- Focus on Inner fine grain communication phase



Communication analysis

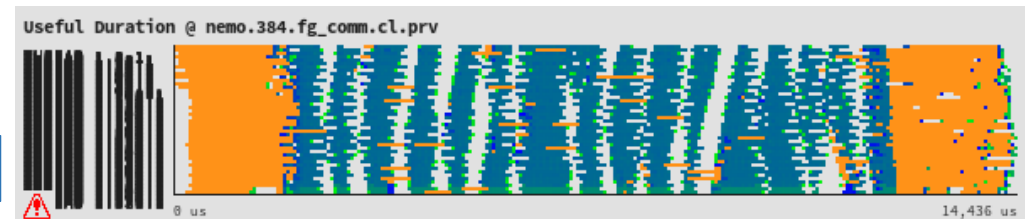


Communication analysis



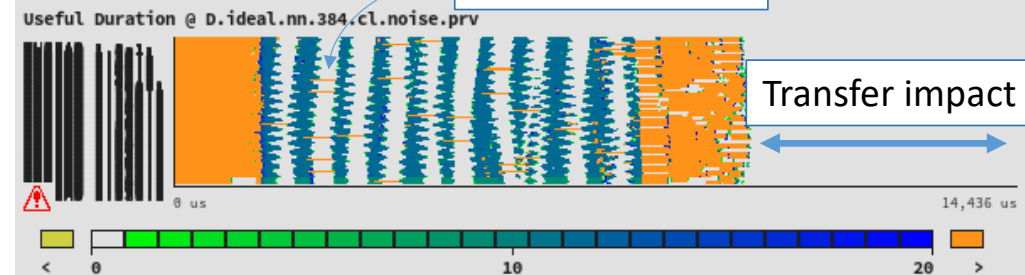
- Serialization caused by OS Noise

Actual run



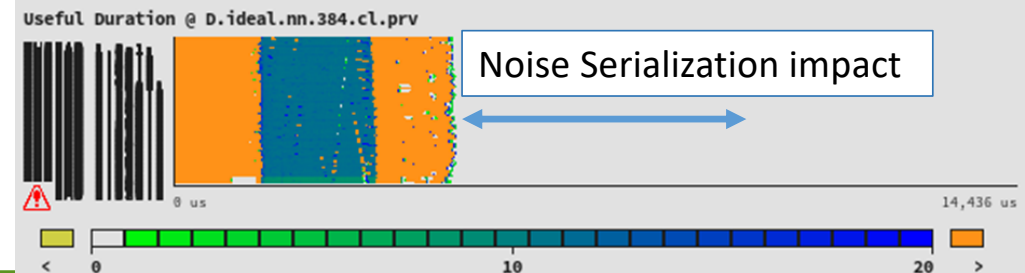
Noise in useful

Ideal Network



Transfer impact

Ideal Network
"Eliminating" OS noise



Noise Serialization impact



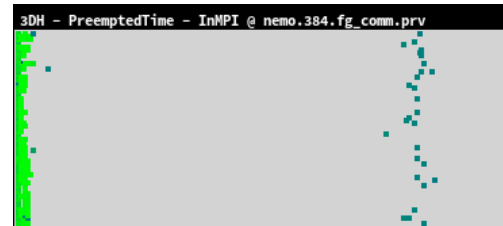
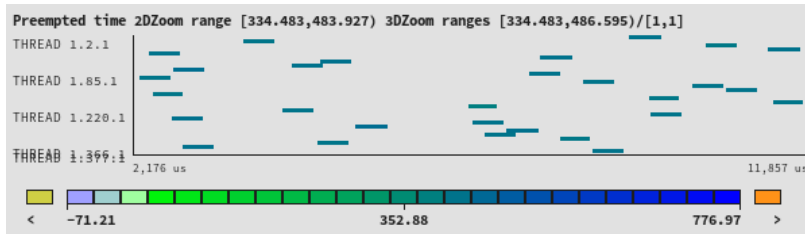
Close look at noise



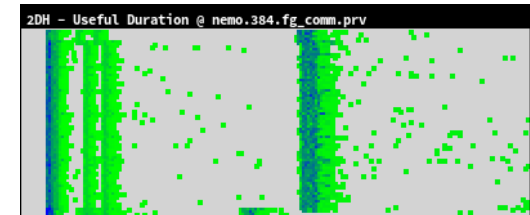
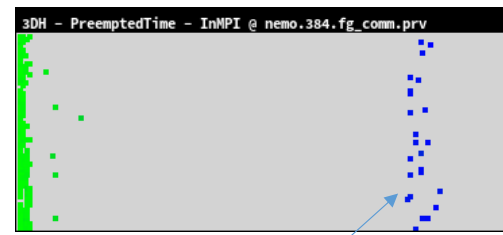
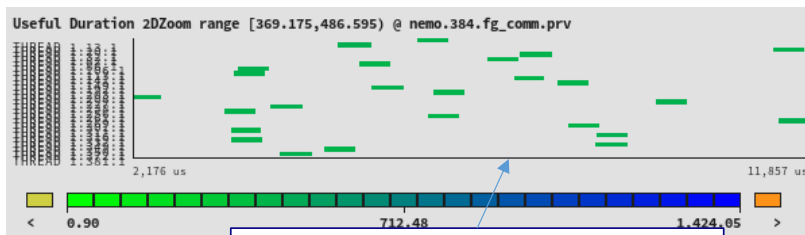
“Preemptions”

“Preempted” time

In MPI



In Useful



Scattered in space and time
Both in MPI and useful

Mode ~400 us

Compared to ~11 us
(and less) computations

Hybrid MPI+OpenMP with relatively dynamic scheduling would be a way to reduce the impact of noise

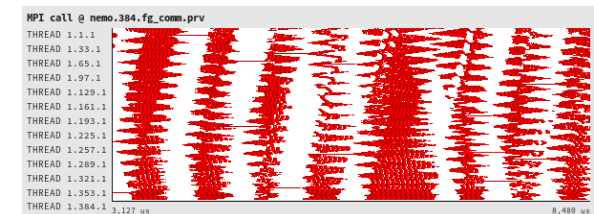
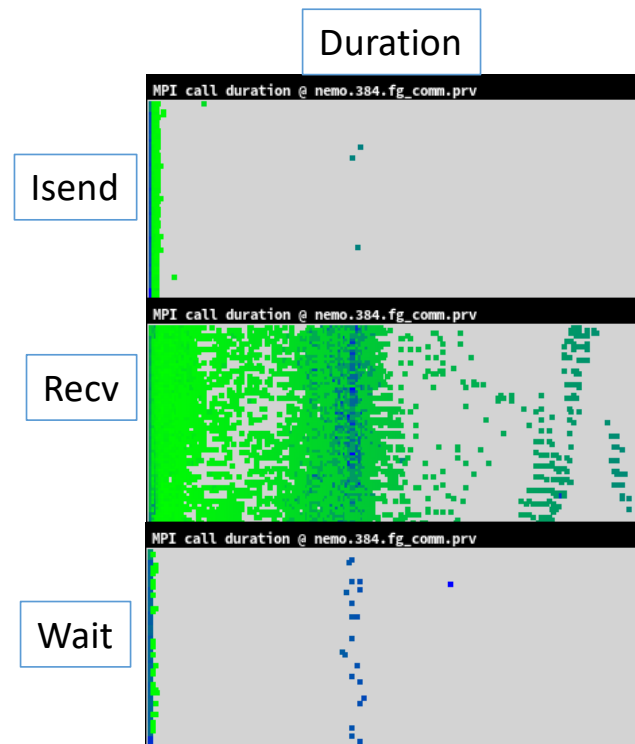
“Noise” cause ?
Cant fight noise, learn to live with it



Close look at noise



- MPI call durations affected & impacted by noise



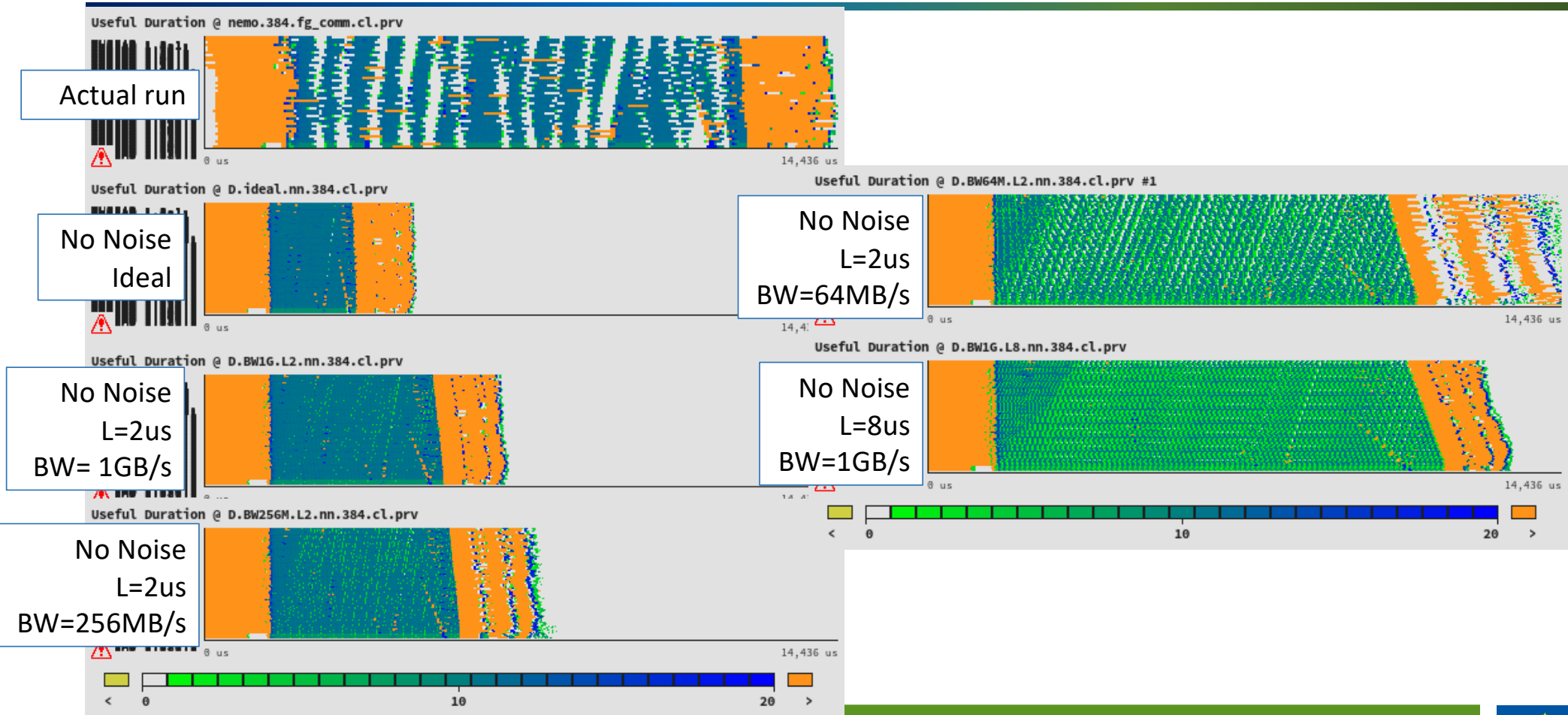
Noise within MPI
Requirement to vendor to implement noise toleration mechanisms within the node in their MPI implementation

“Noise” cause ?

Mode ~400 us



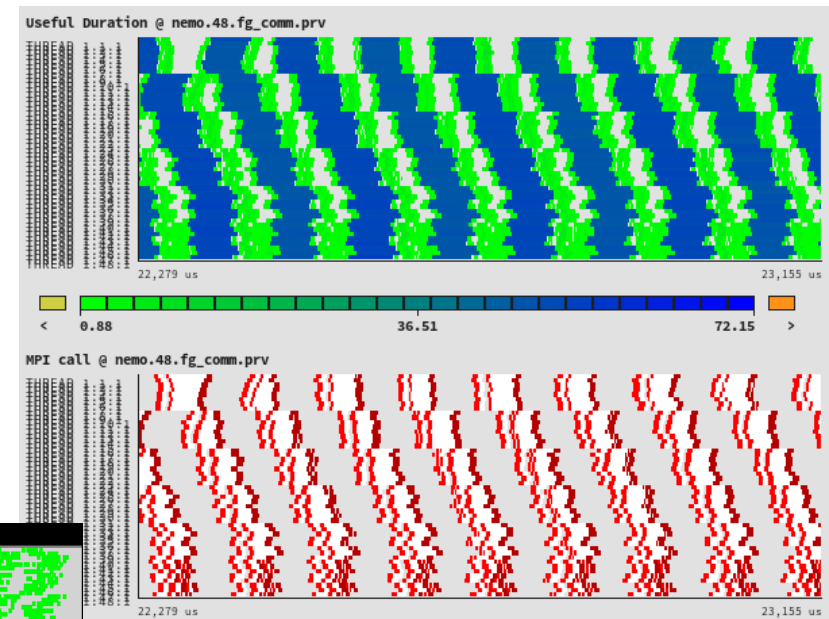
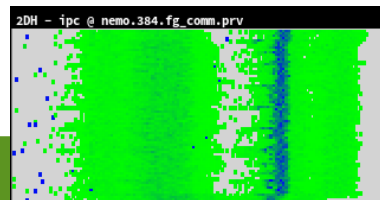
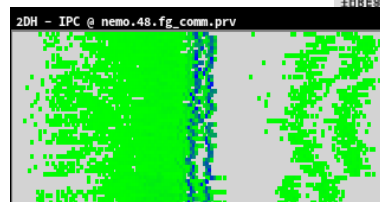
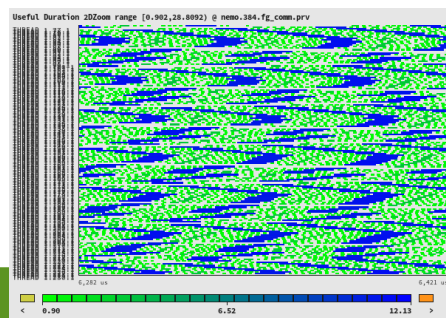
Transfer analysis



The skew



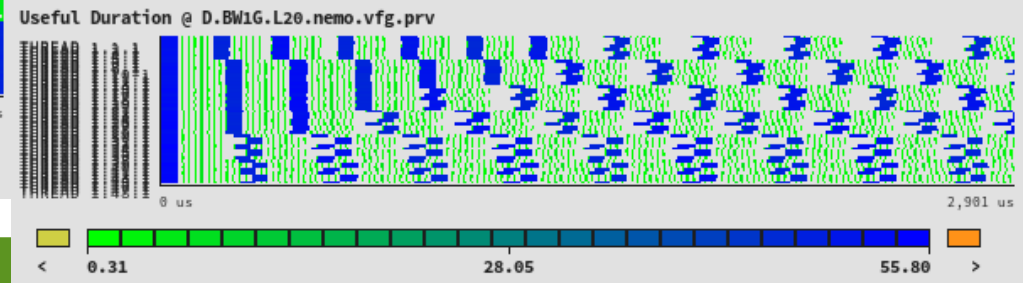
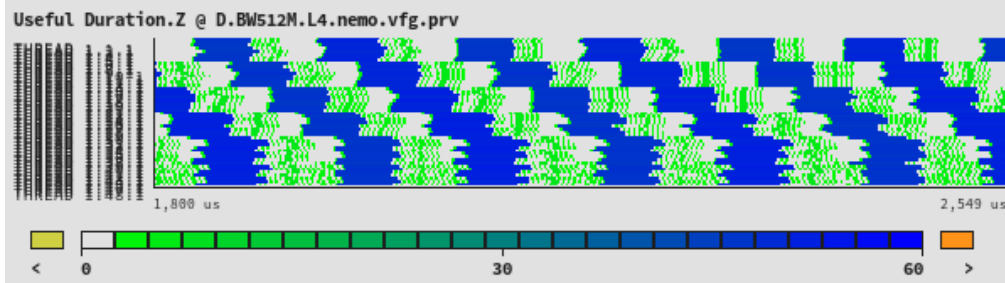
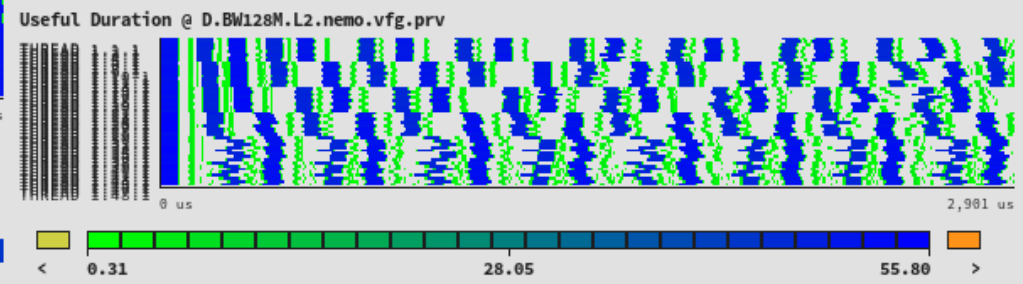
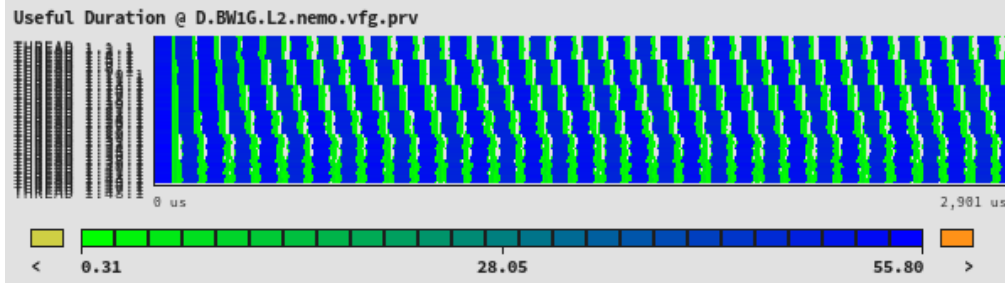
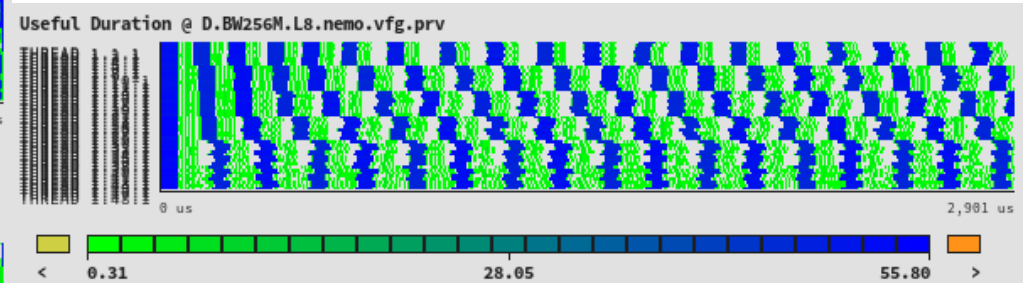
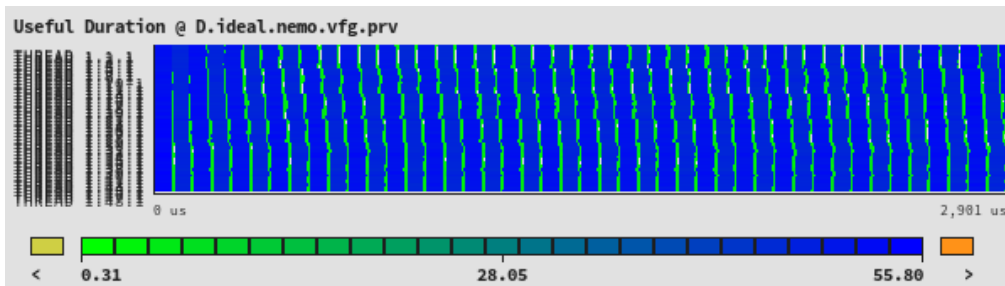
- Pattern in fine grain communication phase
- Efficiency loss on 48 processes (0.66)
- Significant impact at 384
 - Parallel efficiency ~ 0.5
 - Less overlap of computations & strong scaling \rightarrow better IPC
 - Compensating effects



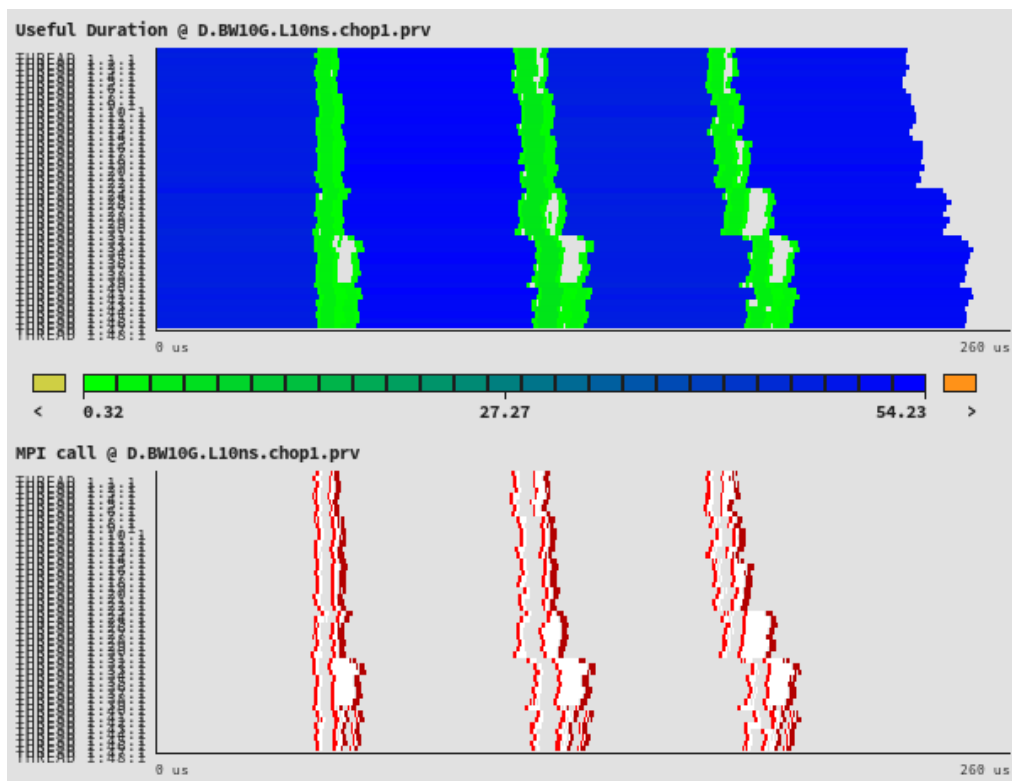
IPC: 0.5 -- 3.5



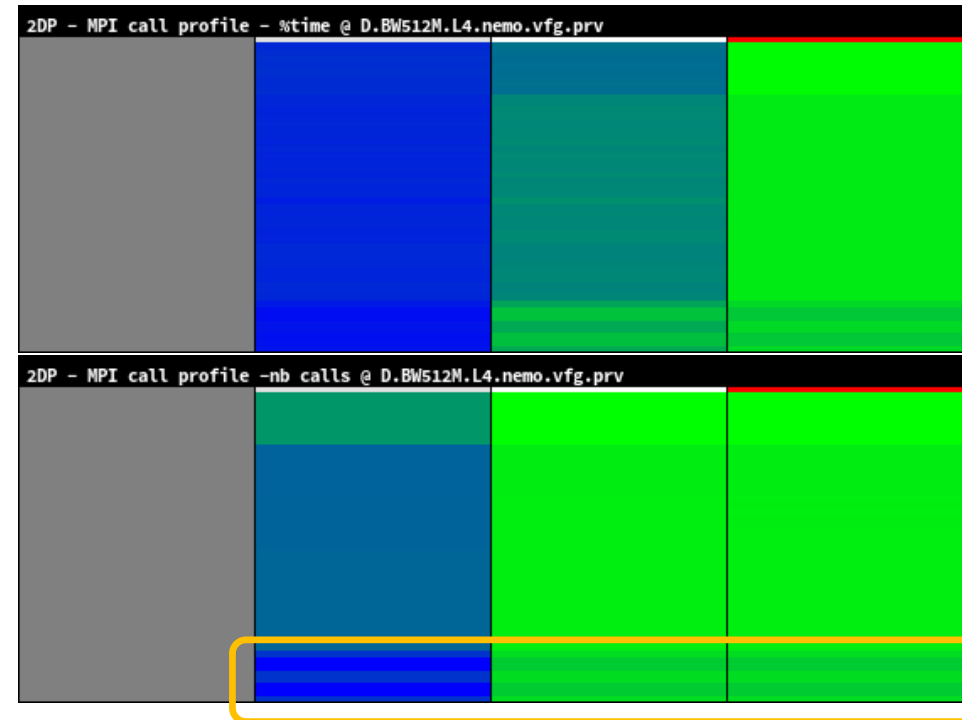
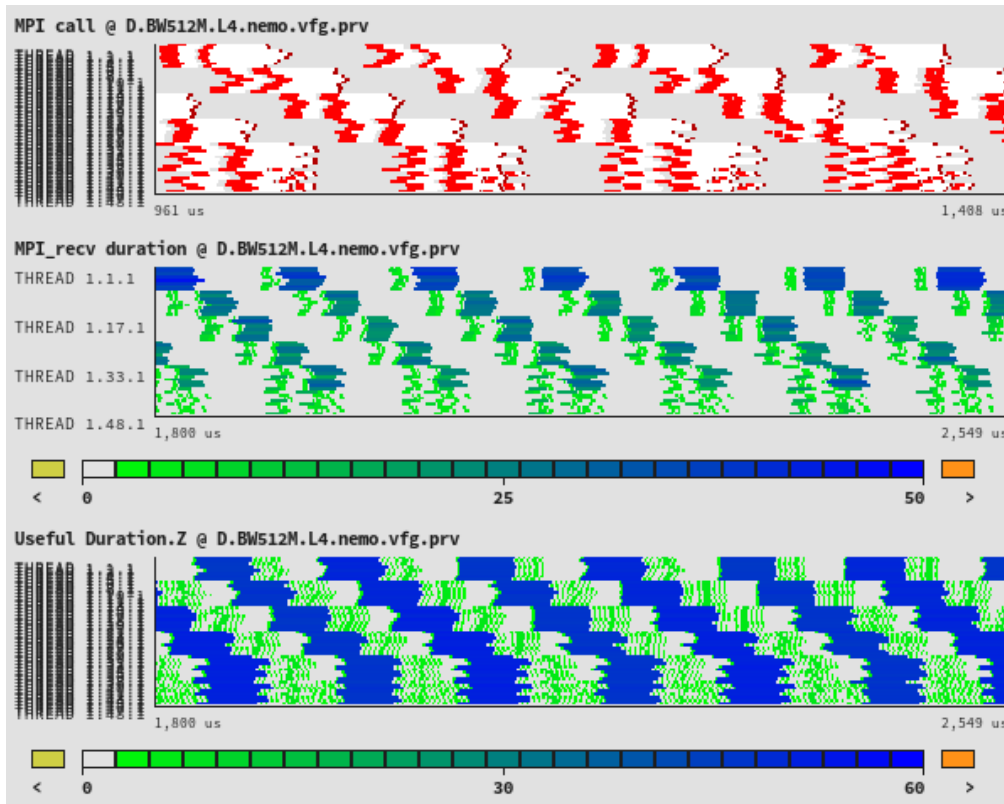
Impact of latency and BW on the skew



Genesis of the skew



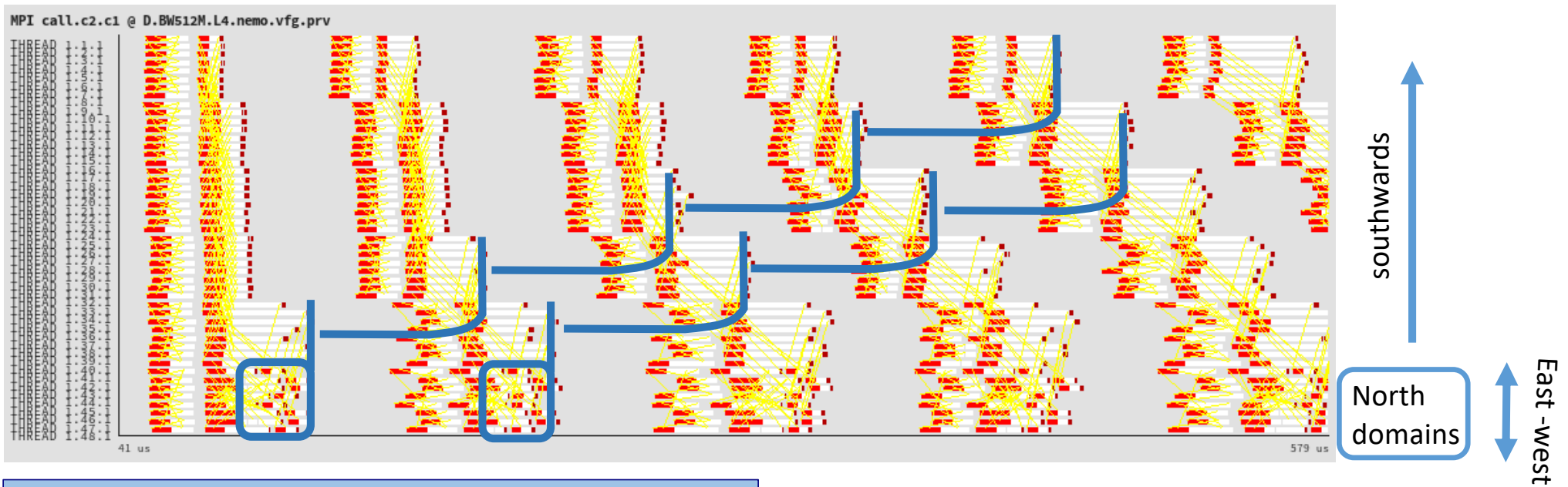
Genesis of the skew



Larger number of communication of processes
in the north domain



Genesis of the skew



Compensate skew → Assign less load to last processes ??

Reorder communication from below and computation may “alleviate” propagation ??

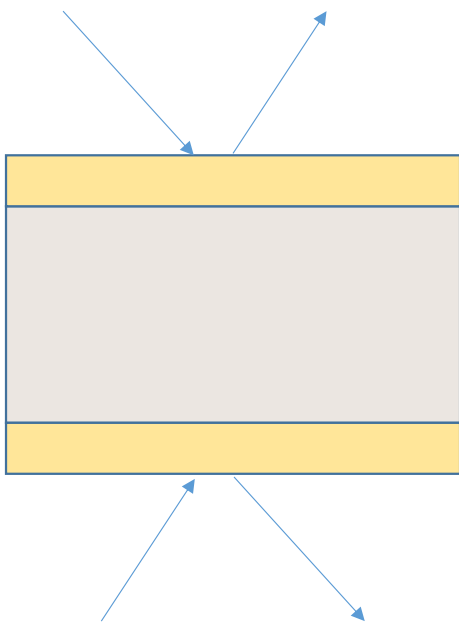
all isends before any receive ? Just on the pole?



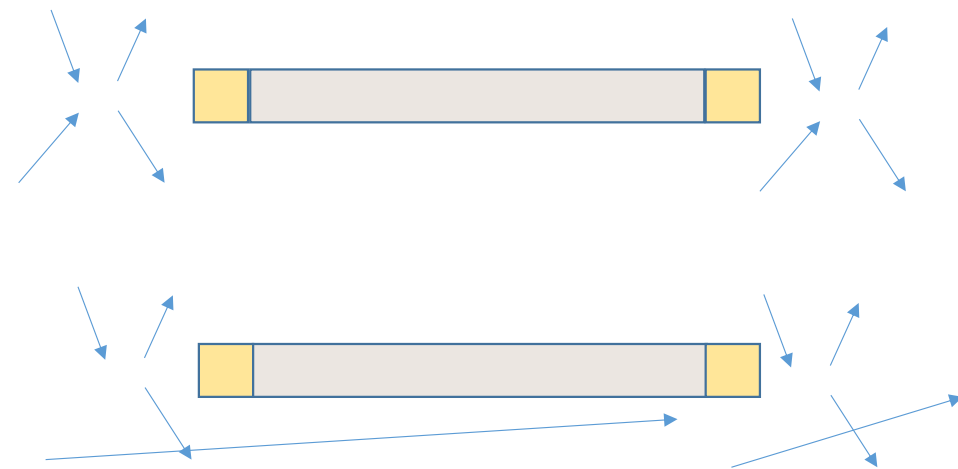
Overlap communication - computation?



“Physical/algorithmic” structure



“computational” structure



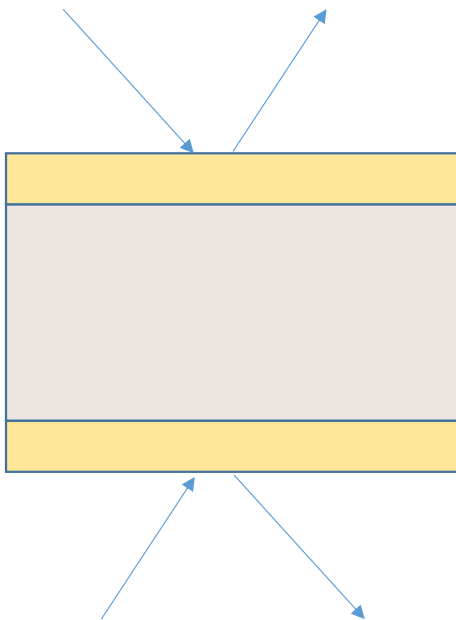
Can be reordered?
Source code reordering (irecv, wait)
tasks and dependences
(more dynamic, adaptive, overhead?)



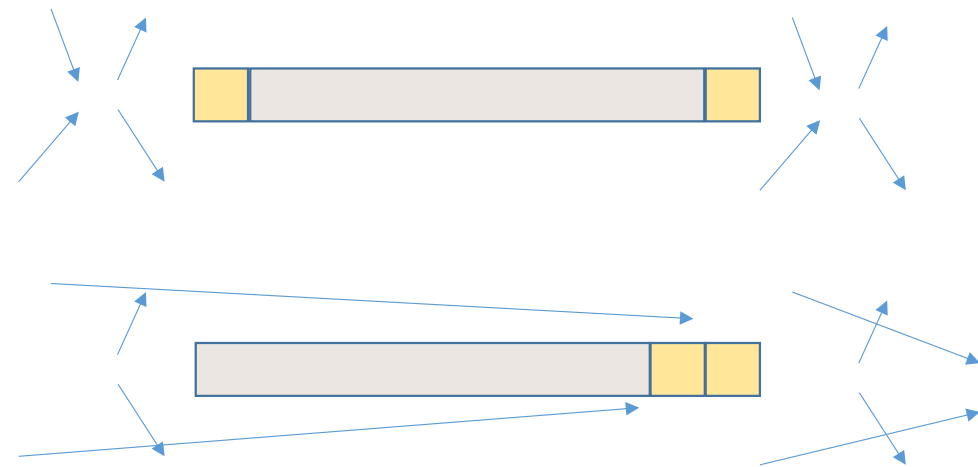
Reorder computation ?



“Physical/algorithmic” structure



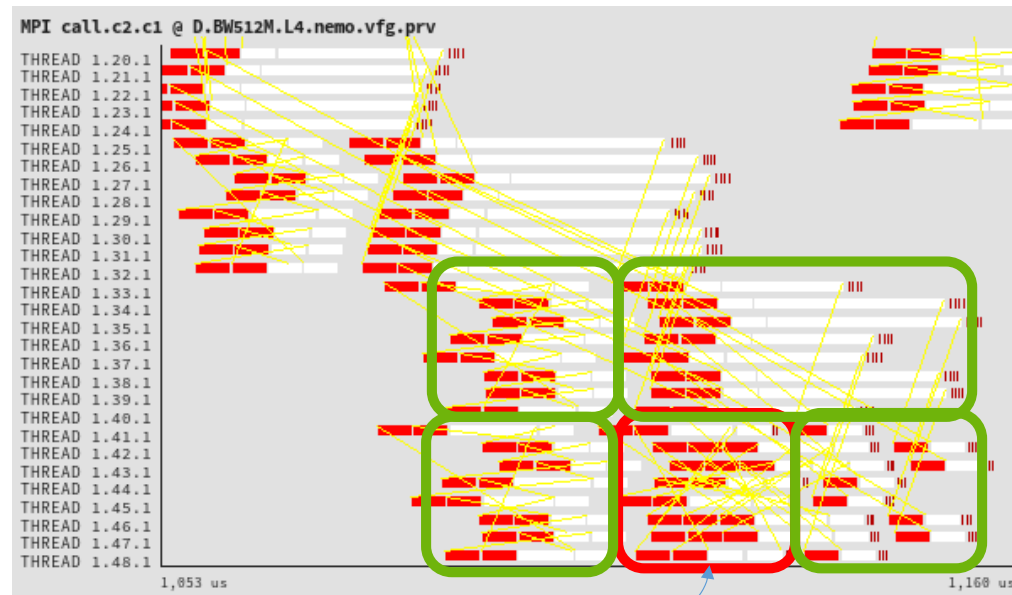
“computational” structure



Can be reordered? Dependences?
Source code reordering
tasks and dependences
(more dynamic, adaptive, overhead?)



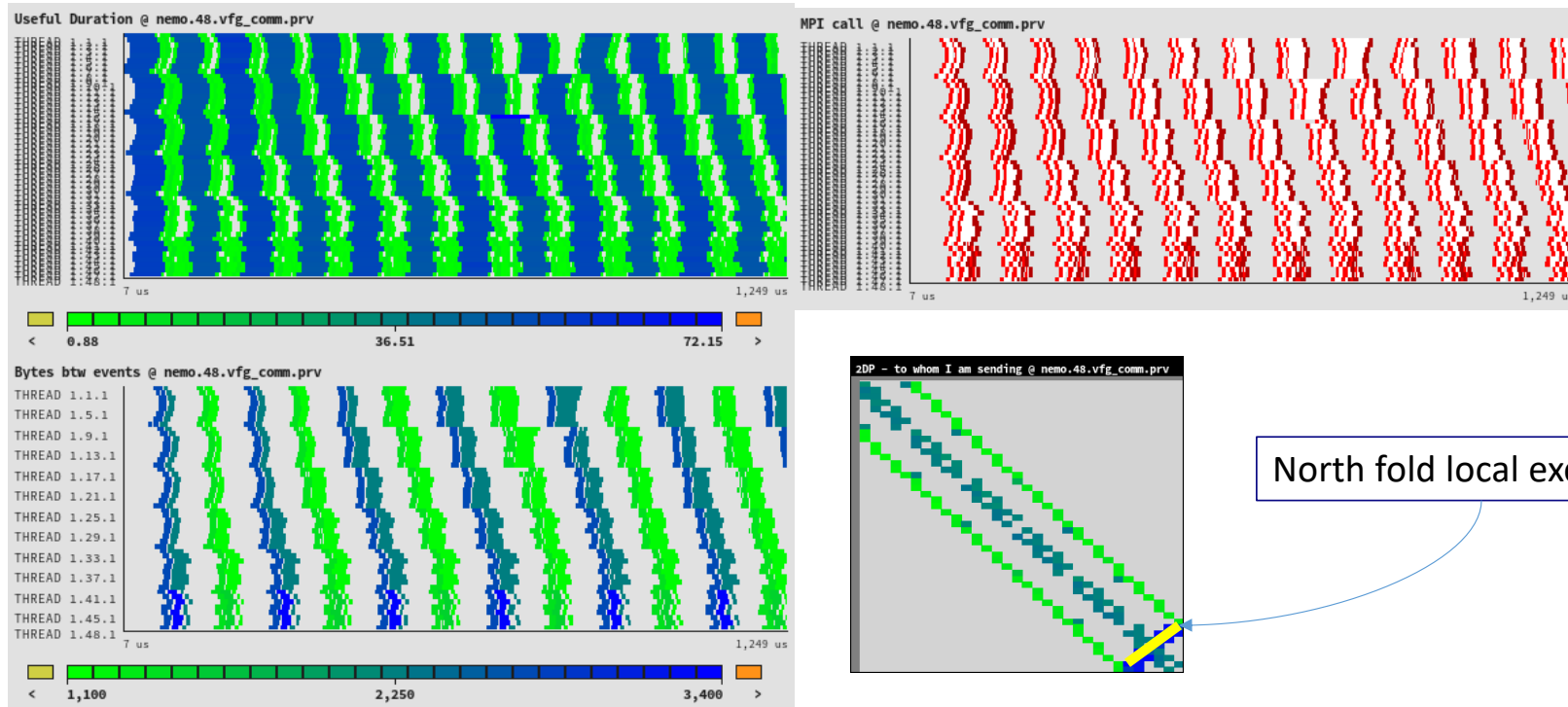
Genesis of the skew



North fold ?



North fold

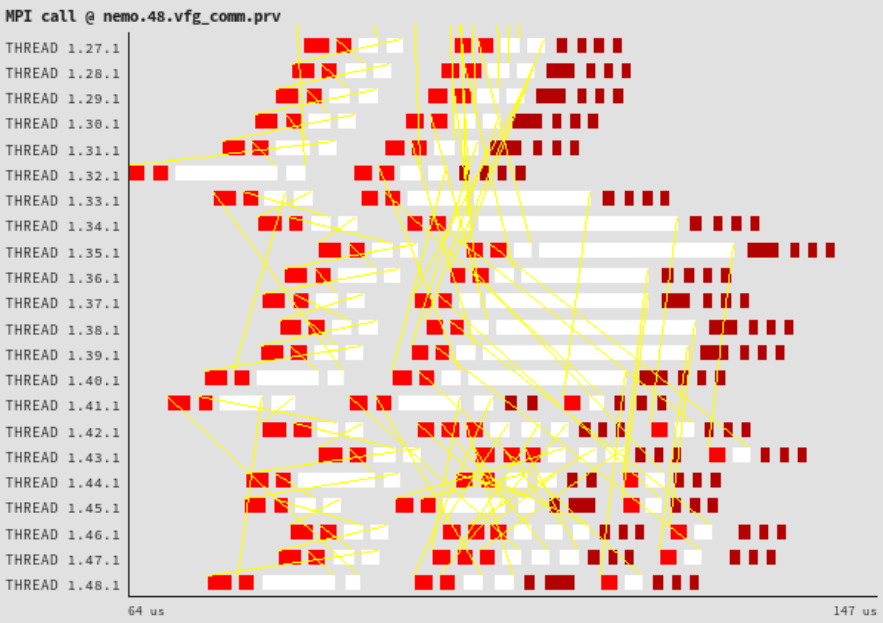


North fold local exchange

Substructure:
are the two phases dependent (potential to overlap)?

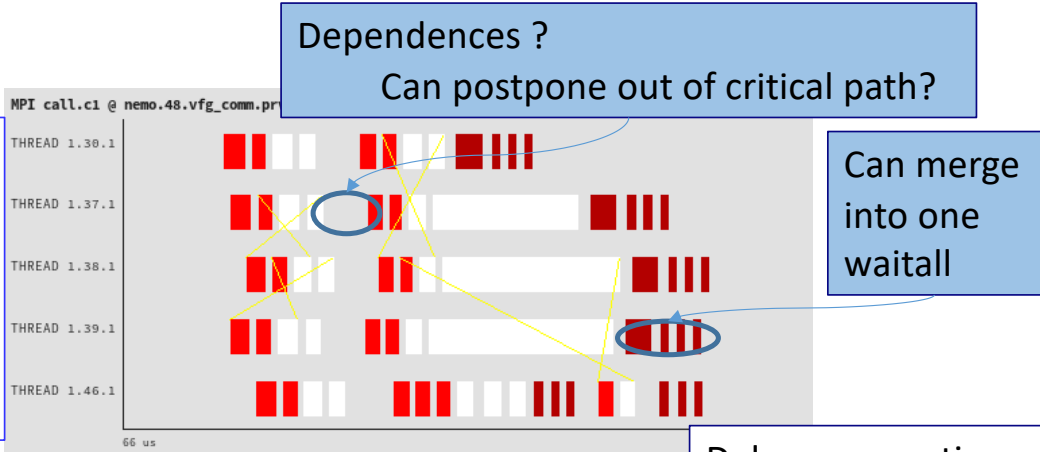


Detailed MPI call sequence

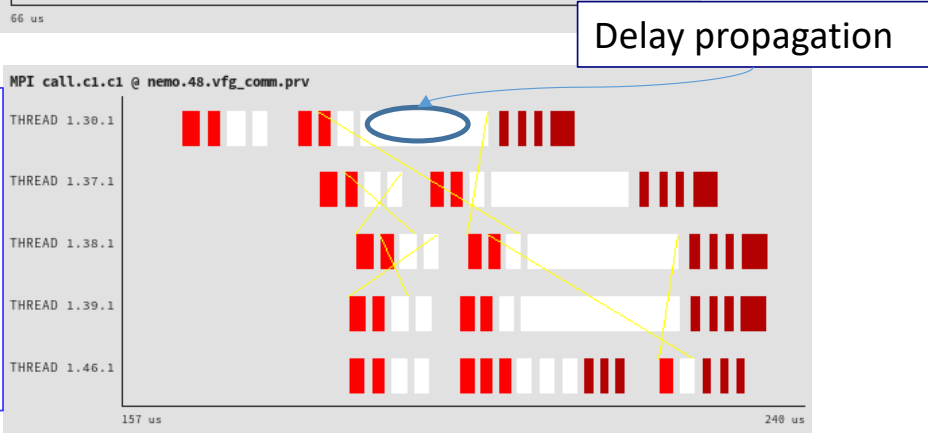


- Outside MPI
- MPI_Recv
- MPI_Isend
- MPI_Wait

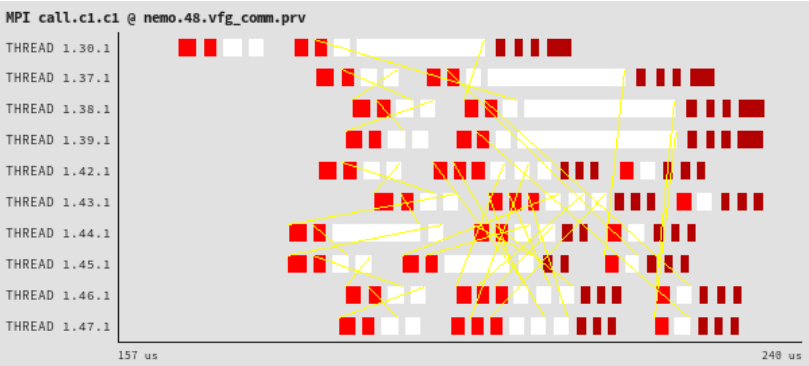
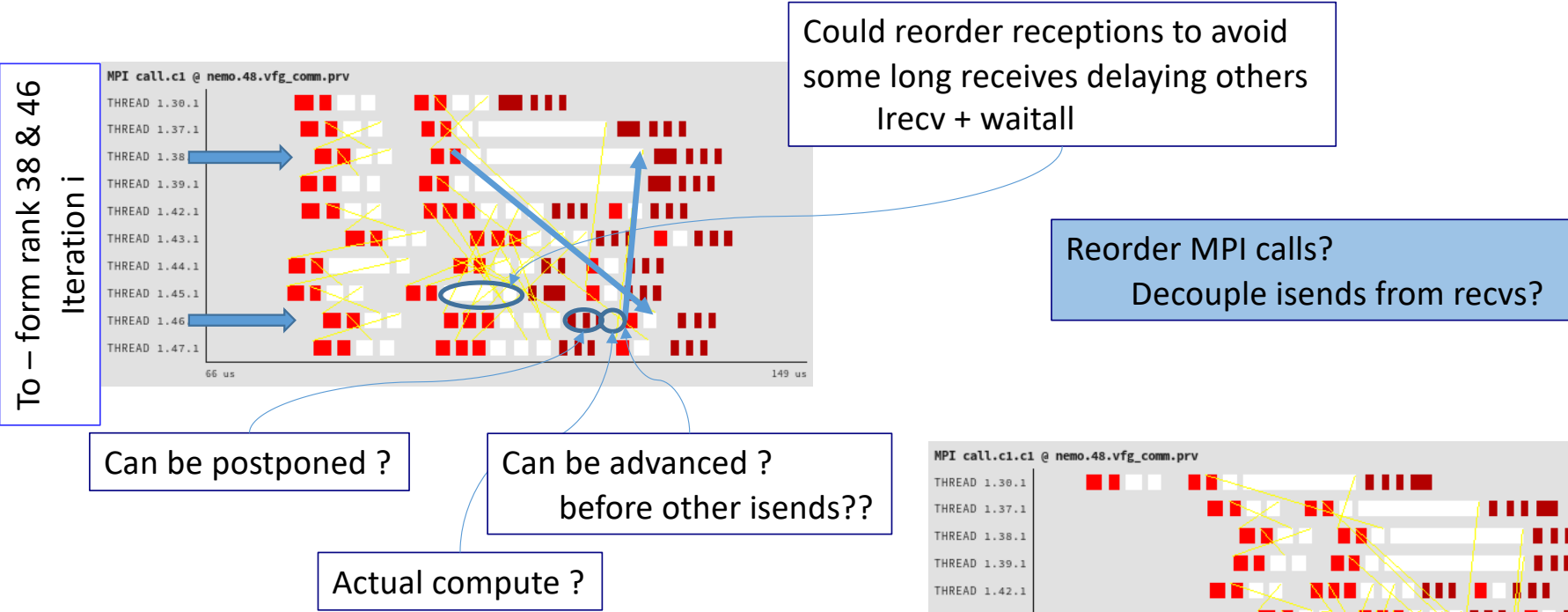
To – form rank 38
Iteration i



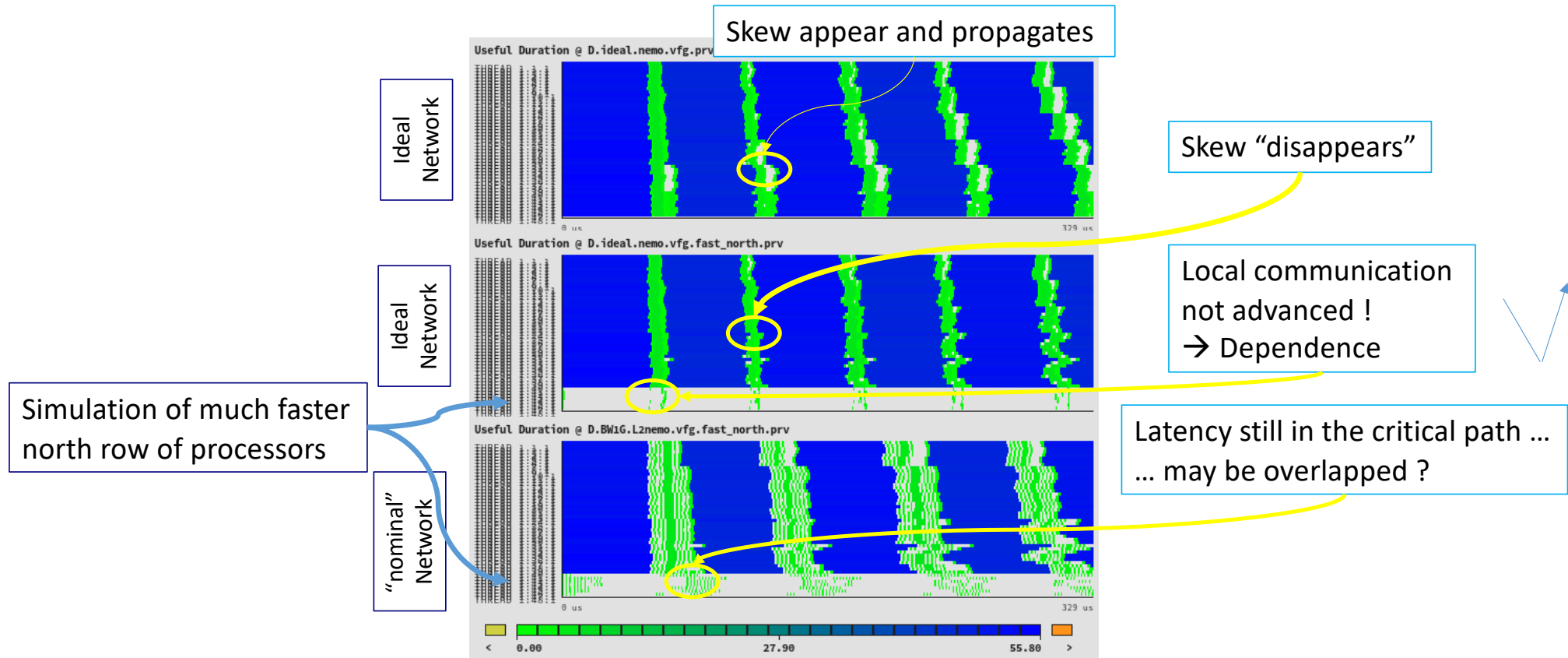
To – form rank 38
Iteration i+1



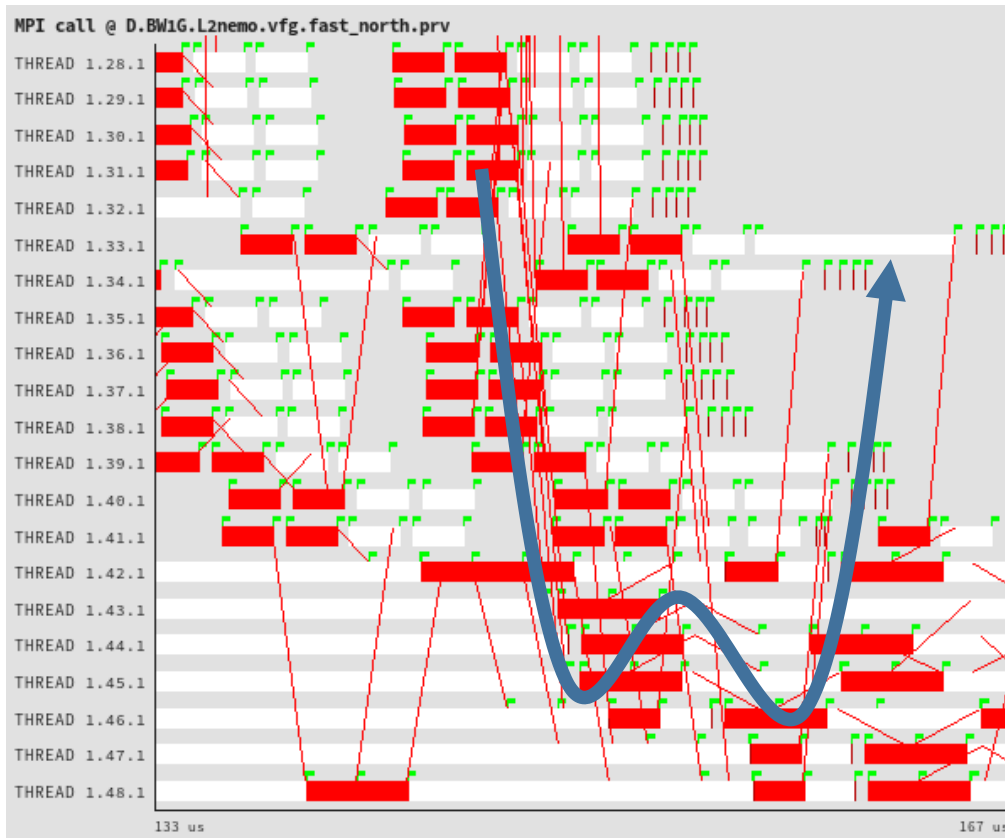
Detailed MPI call sequence



Potential of imbalance?



Dependence chain



Dependence chain?
Mostly communications?
What computation?

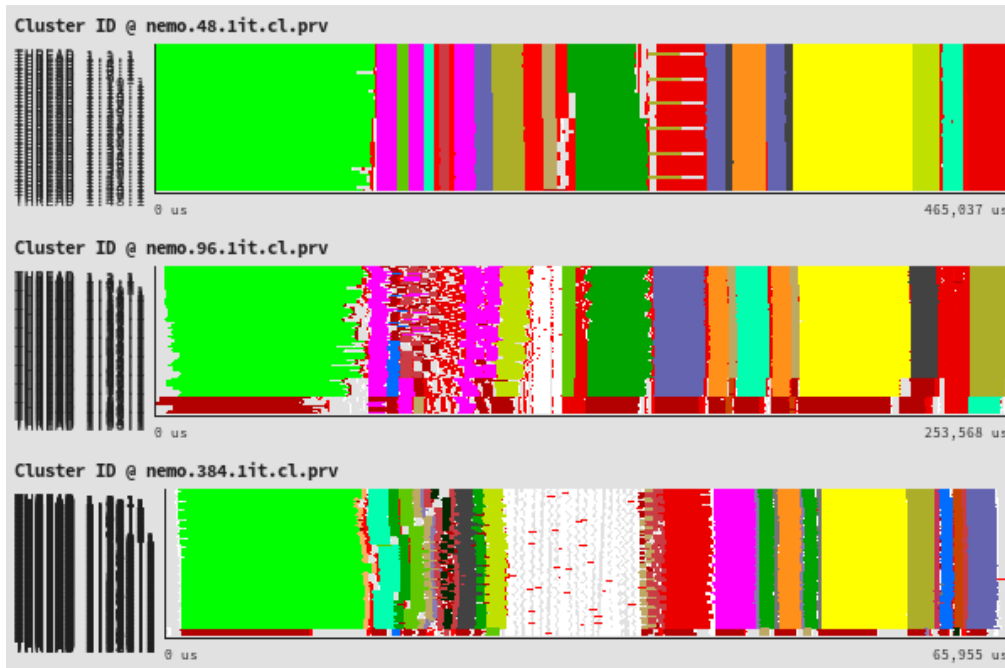
Can be embedded within less processes ?

Can reorder north computations/communications ?

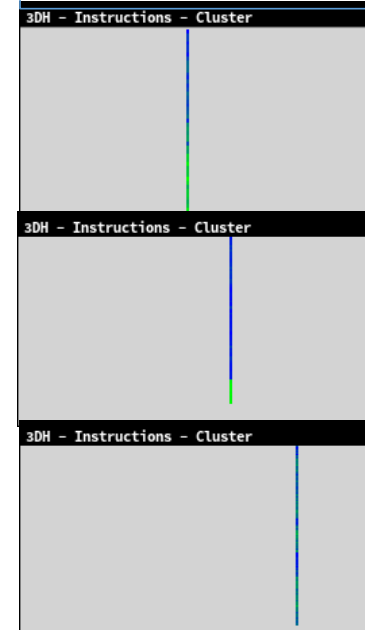
Can reduce #north processes? Only one?



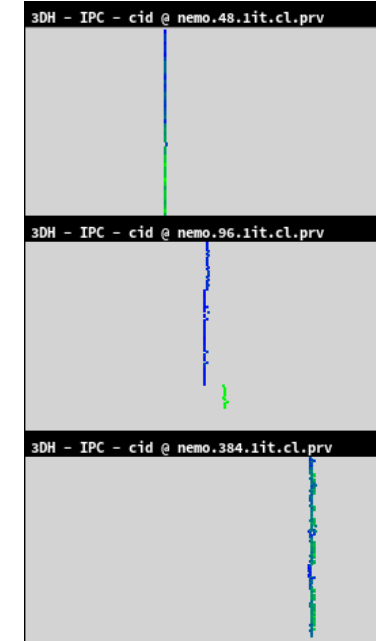
Computation scaling



Yellow Instructions
"strong scaled"



Yellow IPC



Caused by
locality at
L1, L2, L3
level !

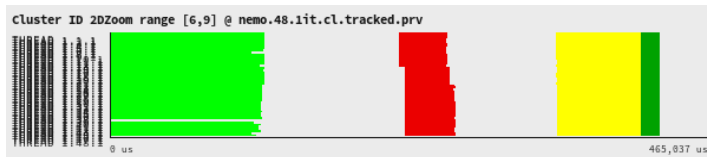
Compensating effects.
Interesting to understand individually why !



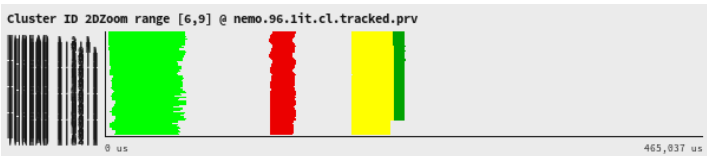
Computation Scaling



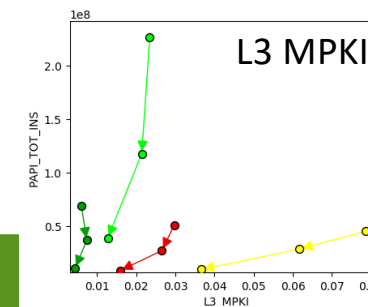
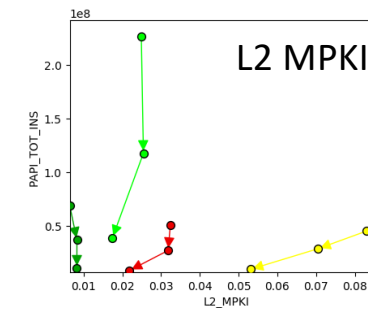
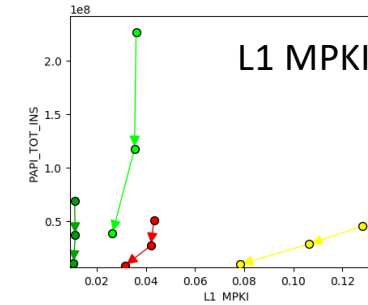
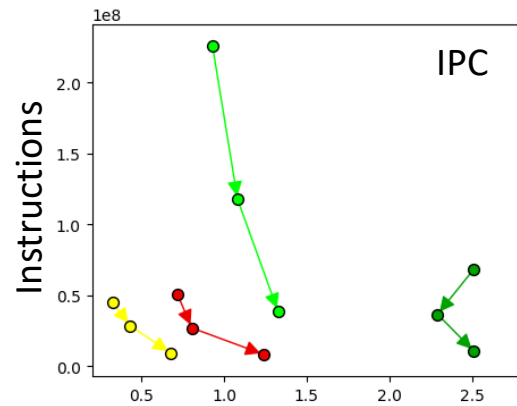
48



96



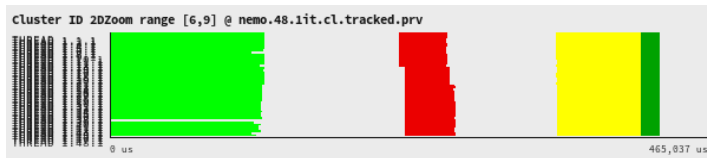
384



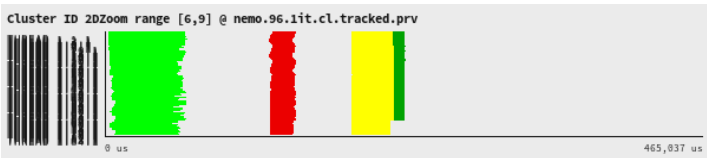
Computation Scaling



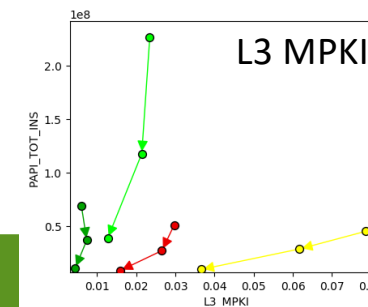
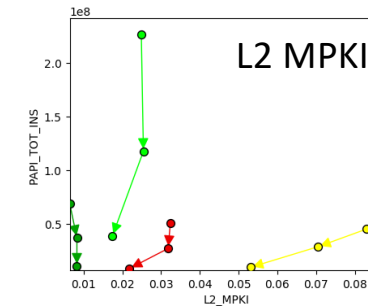
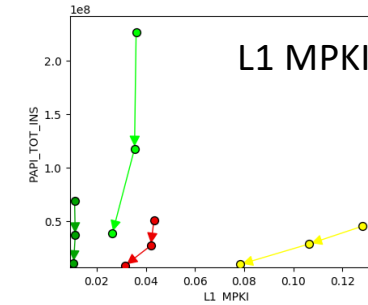
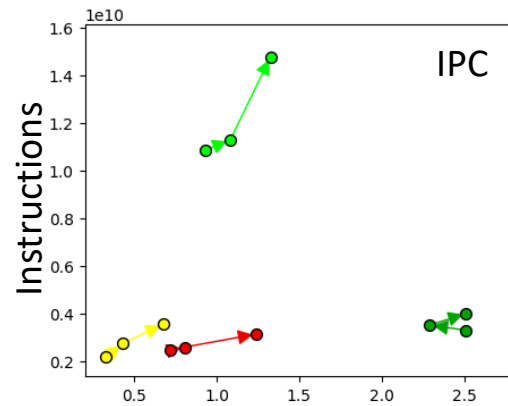
48



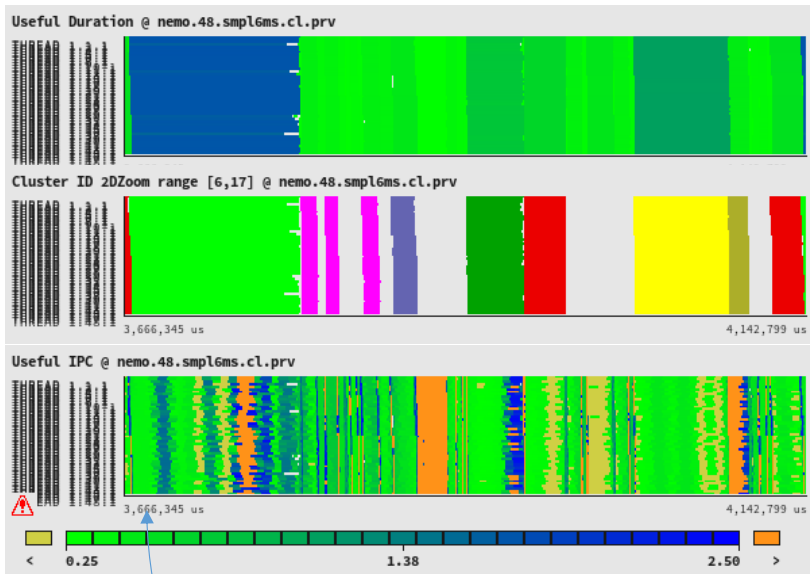
96



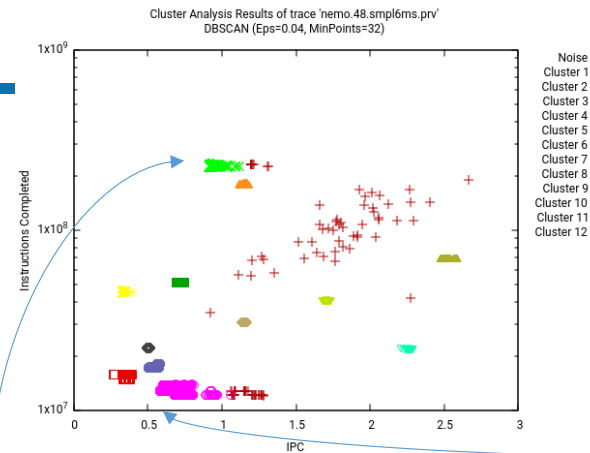
384



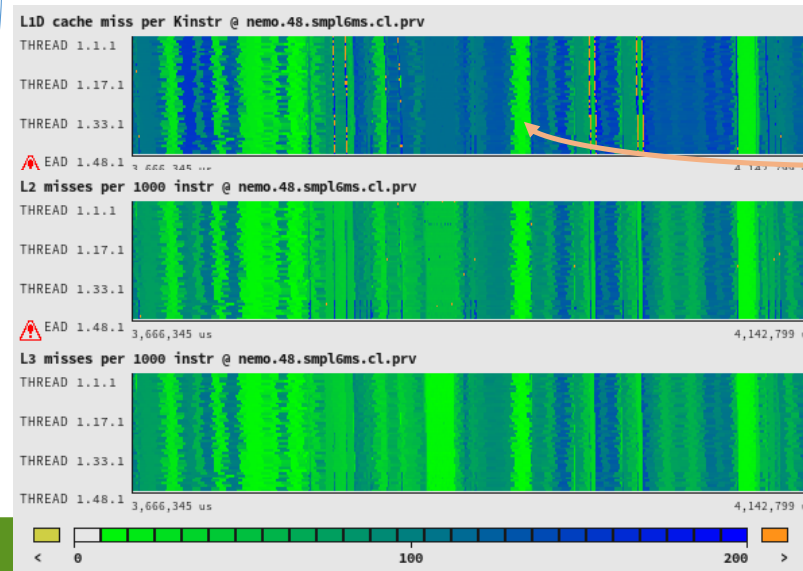
Sampled traces



Very poor IPC sub regions within region of moderate average IPC



Regions with poor IPC

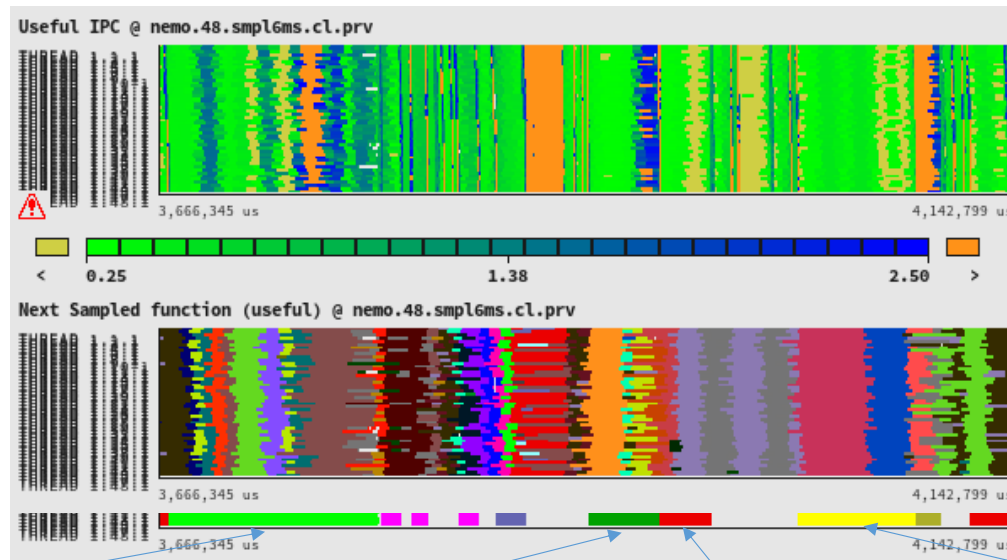


Cache miss ratios "explaining" IPC

Limited benefit of L3



Sampled traces



- End
- domvvl_m..nterpol_ [domvvl_mp_dom_vvl_interpol_]
- eosbn2_mp_bn2_
- zdf_tke_m..tke_tke_ [zdf_tke_mp_tke_tke_]
- zdf_tke_m..tke_avn_ [zdf_tke_mp_tke_avn_]
- zdf_iwm_m..zdf_iwm_ [zdf_iwm_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_]

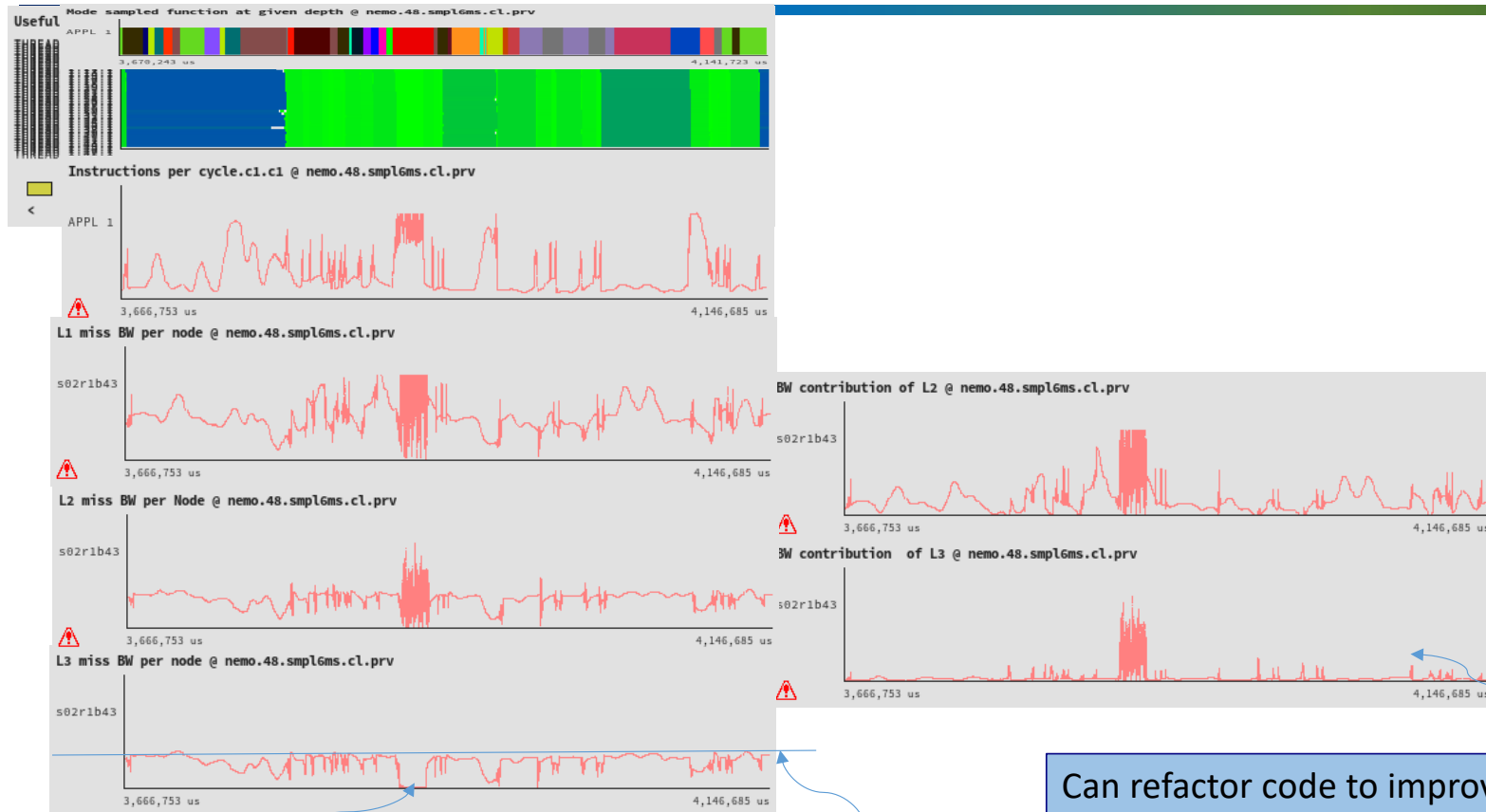
- End
- lbcInk_m.._3d_ptr_ [lbcInk_mp_mpp_lnk_3d_ptr_]
- traadv_f..adv_fct_ [traadv_fct_mp_tra_adv_fct_]
- traldf_i..ldf_iso_ [traldf_iso_mp_tra_ldf_iso_]
- trazdf_m..zdf_imp_ [trazdf_mp_tra_zdf_imp_]
- tra_nxt_vvl



Computation behavior



48 processes



L3 not useful elsewhere ?

Can refactor code to improve L3 use? Blocking?

Working out of L3 in fine grain phase

Node BW limit

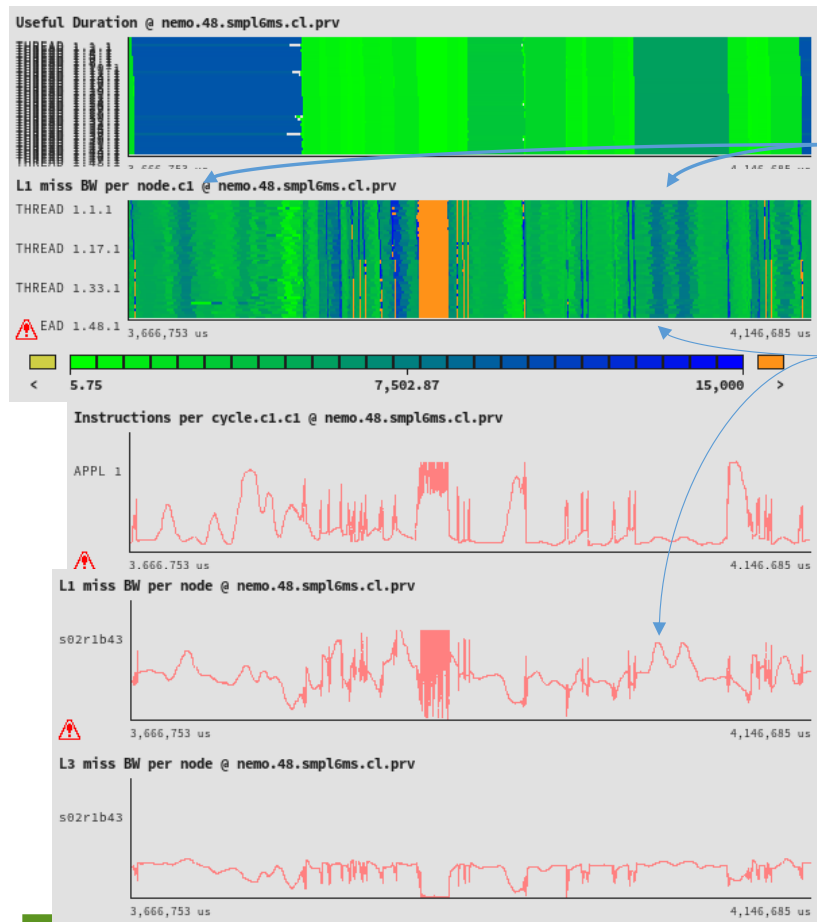
Co-design: no need of L3 ?



Computation behavior



48 processes



Sub substructure

Synchronized BW demands.

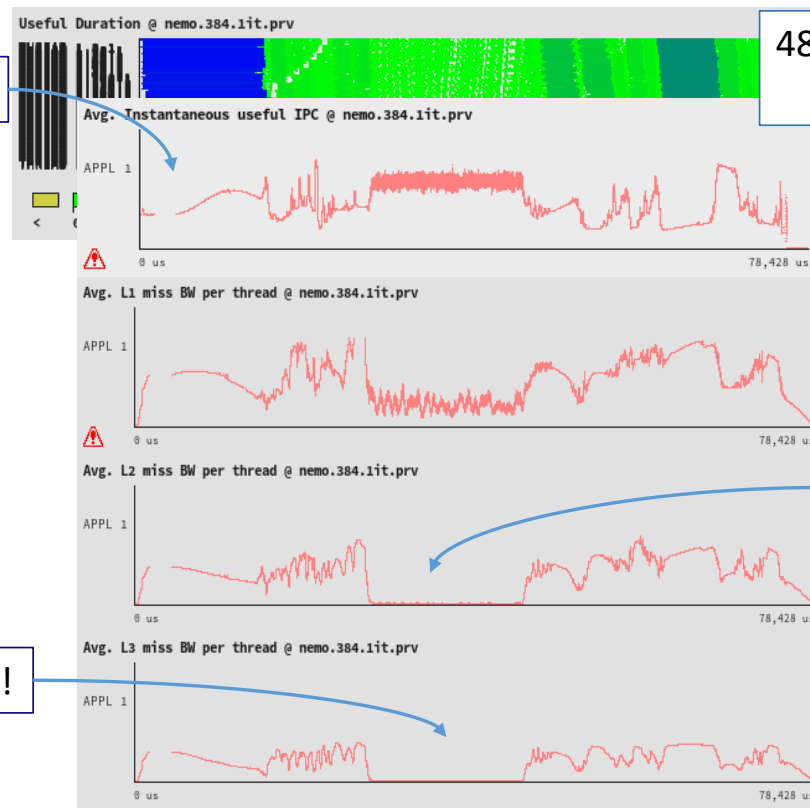
Possibility to unsynchronize?



Comparative Computation behavior



Loss of substructure?



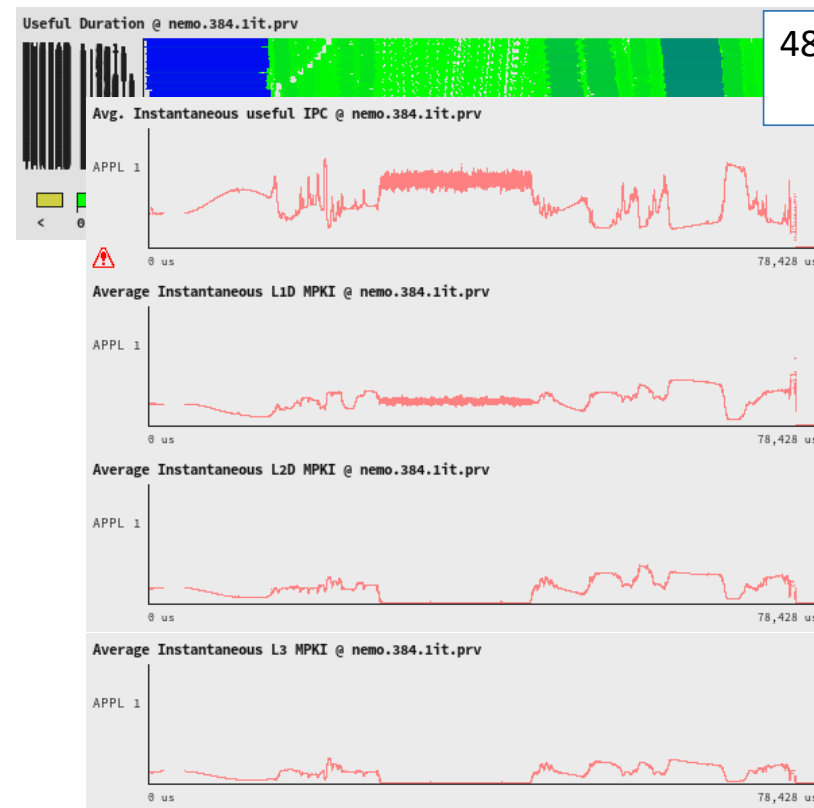
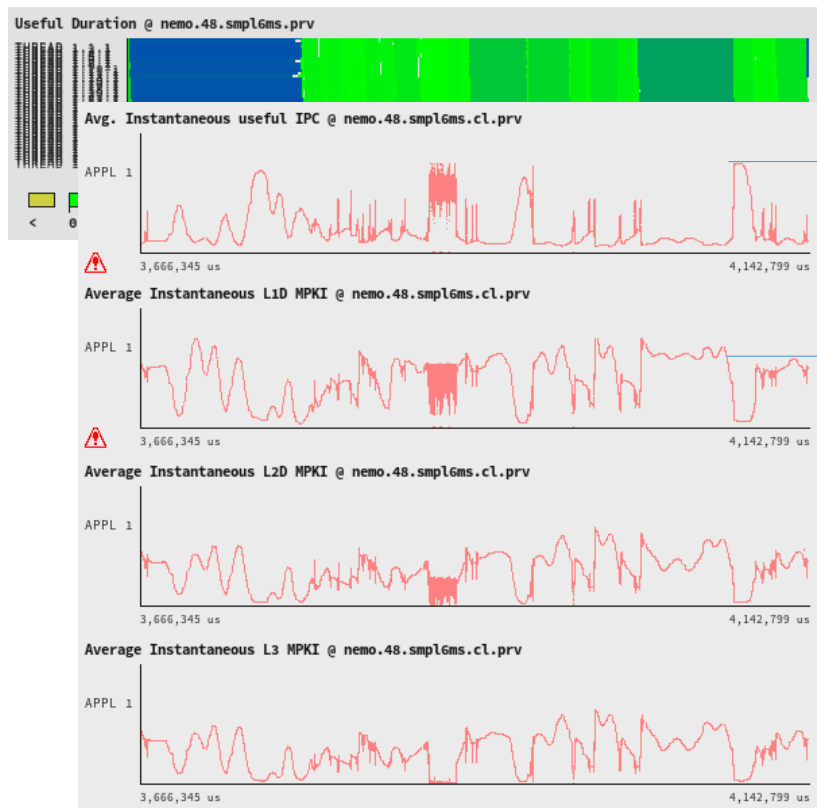
48 vs 384 processes
IPC - BW

Data fits in L2!

Data fits in L3!



Comparative Computation behavior



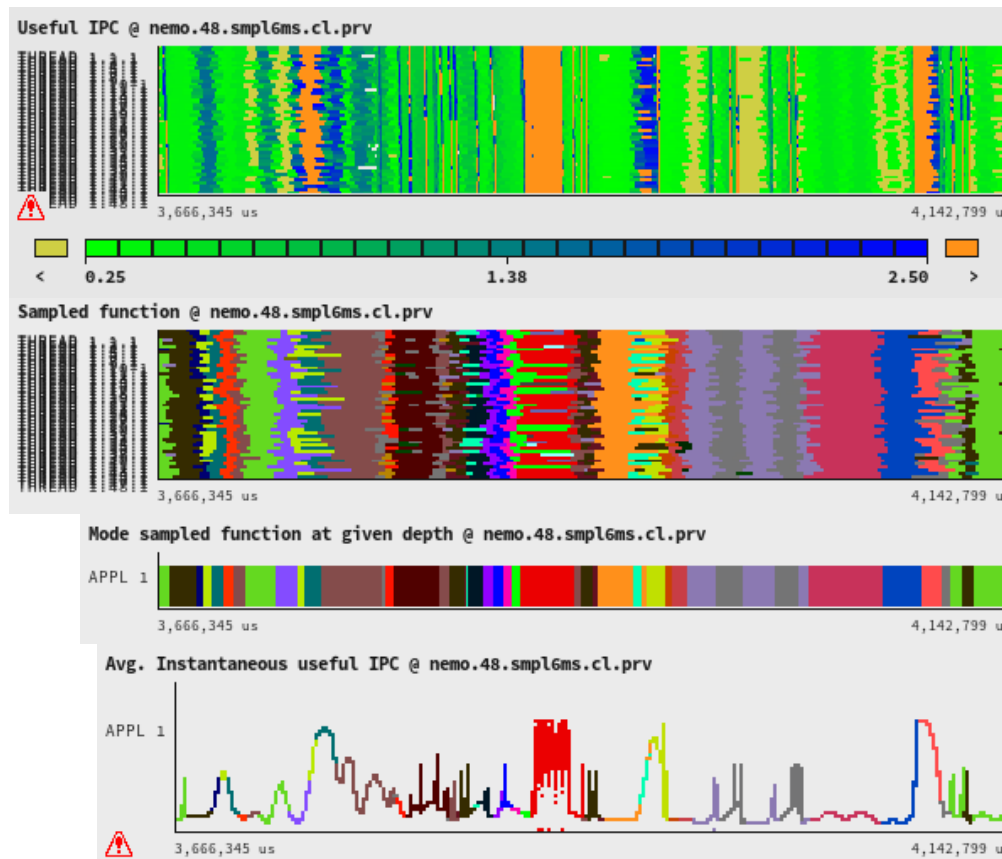
48 vs 384 processes
IPC – Miss Ratios

Overall
improvement
in Miss ratios

*Same scales for all miss ratios, same scales for all IPCs



Link to Source



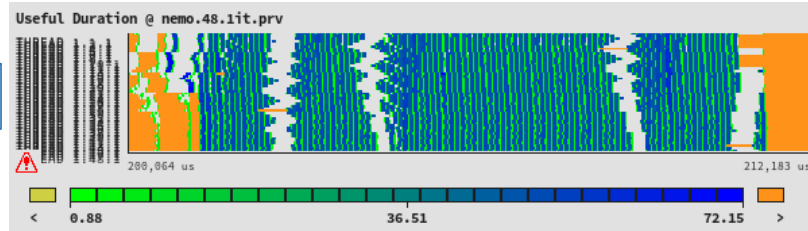
- lbcInk_m.._3d_ptr_ [lbcInk_mp_mpp_lnk_3d_ptr_]
- domvvl_m..nterpol_ [domvvl_mp_dom_vvl_interpol_]
- divhor_m..div_hor_ [divhor_mp_div_hor_]
- eosbn2_mp_bn2_
- zdfcke_m..tke_tke_ [zdfcke_mp_tke_tke_]
- zdfcke_m..tke_avn_ [zdfcke_mp_tke_avn_]
- zdfevd_m..zdf_evd_ [zdfevd_mp_zdf_evd_]
- zdfddm_m..zdf_ddm_ [zdfddm_mp_zdf_ddm_]
- zdfiwm_m..zdf_iwm_ [zdfiwm_mp_zdf_iwm_]
- ldfslp_m..ldf_slp_ [ldfslp_mp_ldf_slp_]
- domvvl_m.._sf_nxt_ [domvvl_mp_dom_vvl_sf_nxt_]
- dynkeg_m..dyn_keg_ [dynkeg_mp_dyn_keg_]
- dynzad_m..dyn_zad_ [dynzad_mp_dyn_zad_]
- dynvor_m..vor_eeen_ [dynvor_mp_vor_eeen_]
- dynldf_l..ldf_lap_ [dynldf_lap_blp_mp_dyn_ldf_lap_]
- dynhpg_m..hpg_sco_ [dynhpg_mp_hpg_sco_]
- dynspg_t.._spg_ts_ [dynspg_ts_mp_dyn_spg_ts_]
- dynzdf_m..dyn_zdf_ [dynzdf_mp_dyn_zdf_]
- sshwzv_mp_wzv_
- traqsr_m..tra_qsr_ [traqsr_mp_tra_qsr_]
- ldftra_m..eiv_trp_ [ldftra_mp_ldf_eiv_trp_]
- traadv_m..tra_adv_ [traadv_mp_tra_adv_]
- traadv_f..adv_fct_ [traadv_fct_mp_tra_adv_fct_]
- traadv_f.._nonosc_ [traadv_fct_mp_nonosc_]
- traldf_i..ldf_iso_ [traldf_iso_mp_tra_ldf_iso_]
- trazdf_m..zdf_imp_ [trazdf_mp_tra_zdf_imp_]
- dynnxt_m..dyn_nxt_ [dynnxt_mp_dyn_nxt_]
- domvvl_m.._sf_swp_ [domvvl_mp_dom_vvl_sf_swp_]
- stpctl_m..stp_ctl_ [stpctl_mp_stp_ctl_]
- eosbn2_mp_rab_3d_
- zdfsh2_m..zdf_sh2_ [zdfsh2_mp_zdf_sh2_]
- tra_nxt_vvl
- eosbn2_m.._insitu_ [eosbn2_mp_eos_insitu_]
- eosbn2_m..itu_pot_ [eosbn2_mp_eos_insitu_pot_]



Solver Computational Scaling



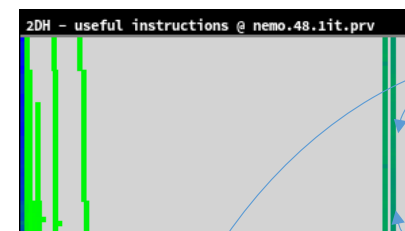
48



Duration Histogram

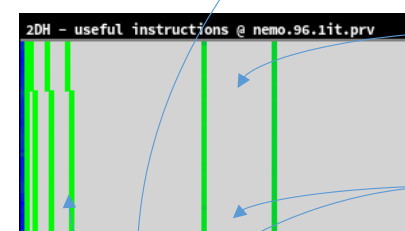
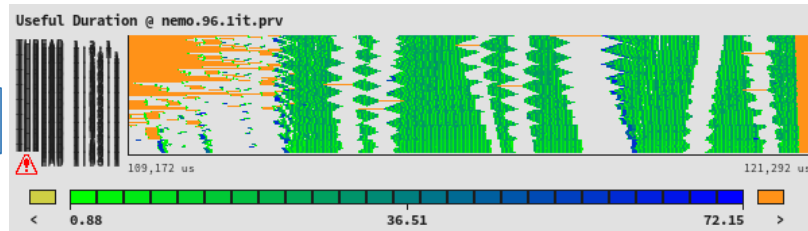


Instructions Histogram



Poor scaling

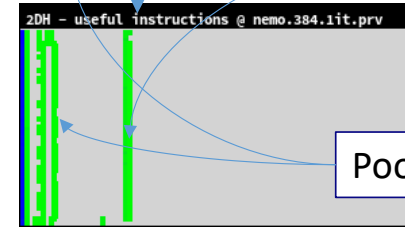
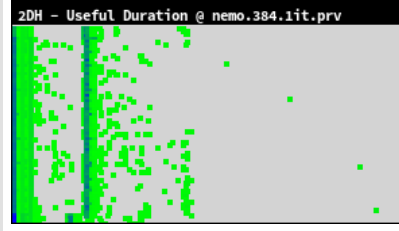
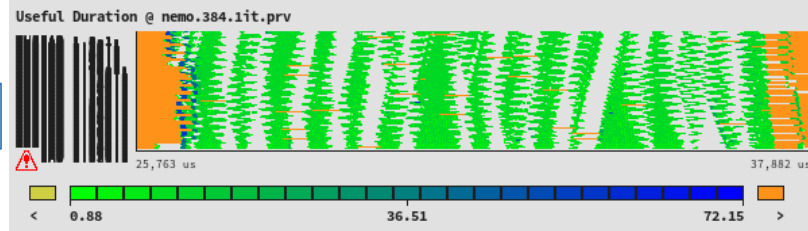
96



Strange scaling

Strange scaling

384



Poor scaling

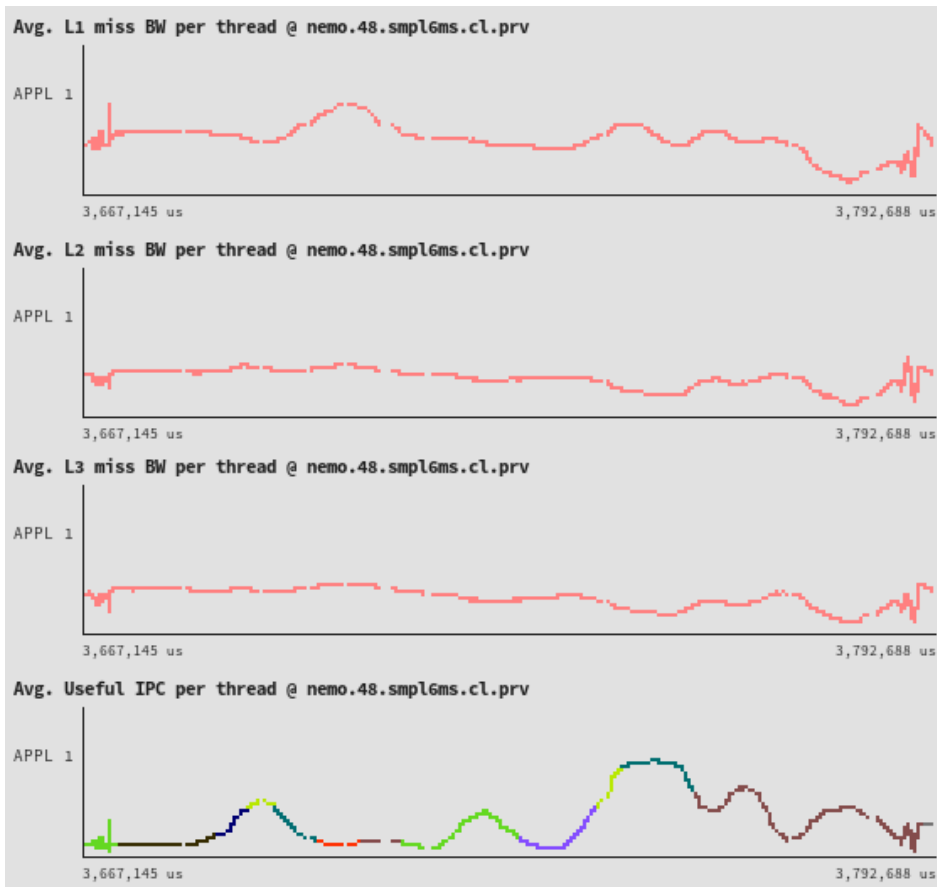
Same time scale
~same duration 😞

<~ fair burst
duration scaling

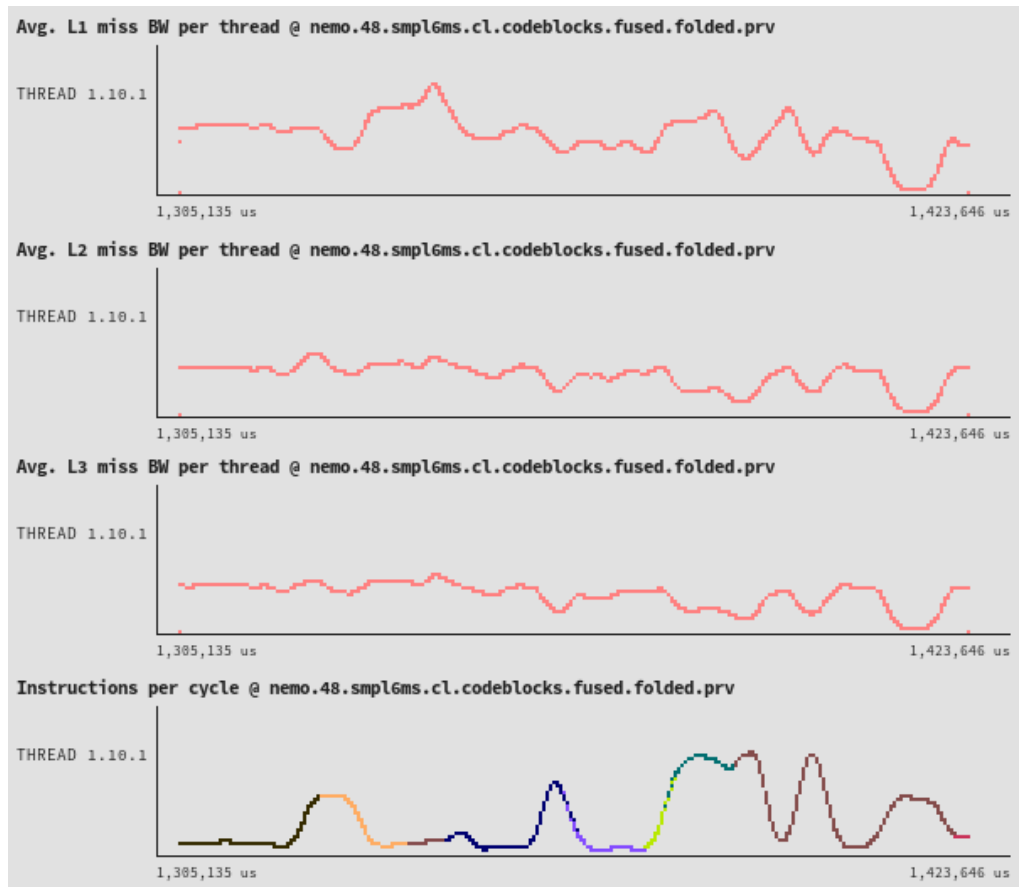
~ poor/strange
instruction scaling



Increasing accuracy ?



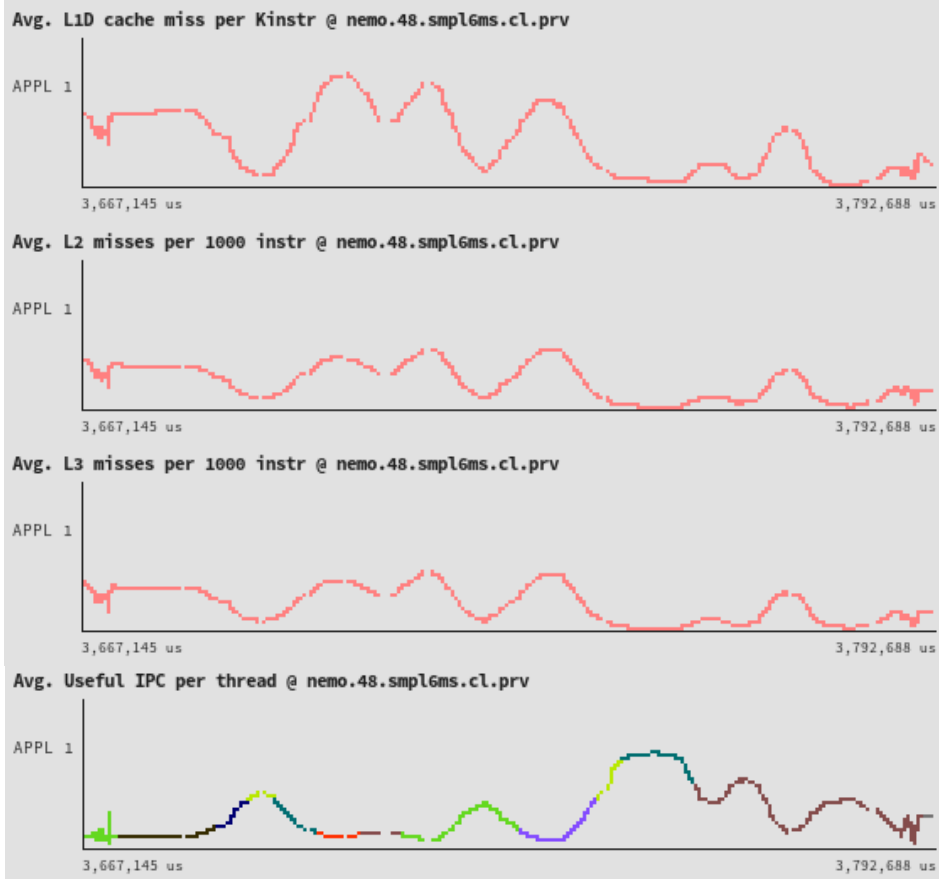
Aggregated Sampling



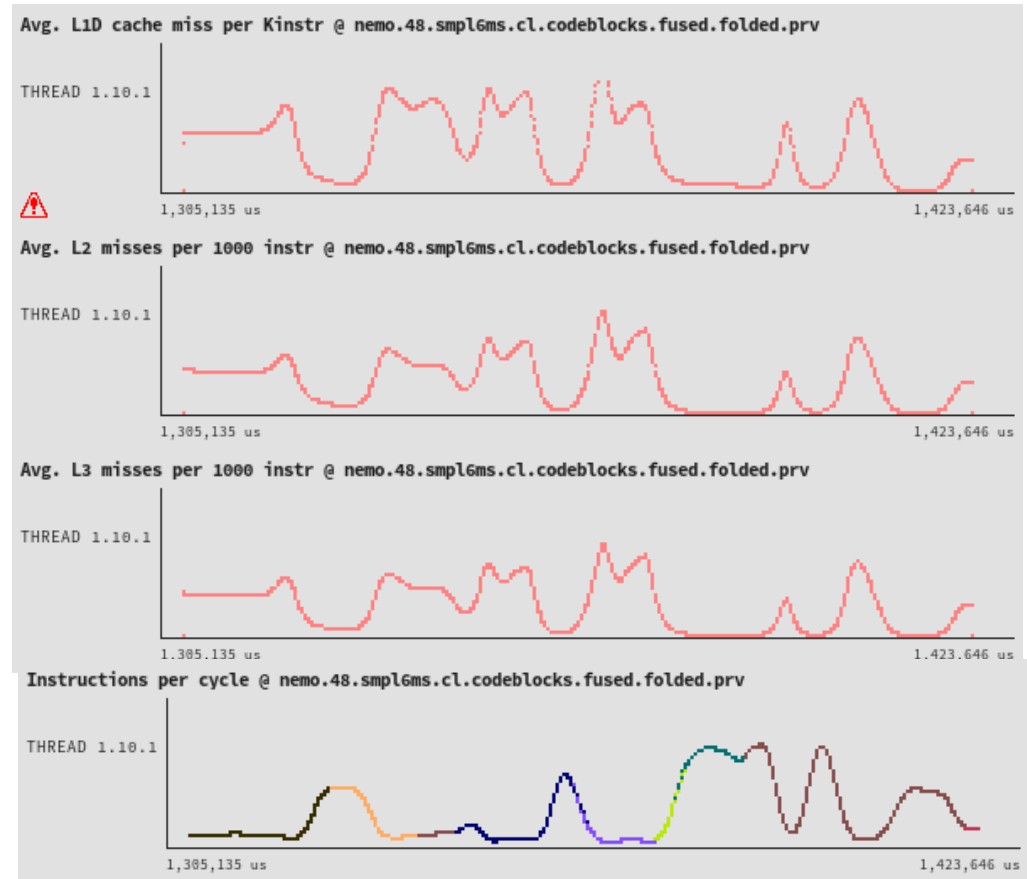
Folding Sampling



Increasing accuracy ?



Aggregated Sampling



Folding Sampling
Cluster 1



Recommendations



- Merge sequence of waits for Isends to waitall
- Reorder north: advance isends postpone recs (if possible)
- Assign less load to north fold to take internal communication out of the critical path (combined with reordering of comms)
- Reduce #north processes?
 - → only one ?
- Increase granularity in solver
 - Gather to single rank per node → Solve → Scatter
- OpenMP in solver
 - → Tasks with dependences between communications and computations for out of order execution and noise tolerance

Possibility for a Mockup POC



Recommendations



- Try to improve locality → better IPC
 - L3 usage: Blocking?
- Convergence of numerical method: Jacobi vs. Gauss-Seidel ?





Performance Optimisation and Productivity

A Centre of Excellence in Computing Applications

Contact:

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

<mailto:pop@bsc.es>



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 676553.

