



# Experiences of an HPC analyst in AI land

ISC workshop. June 26<sup>th</sup>, 2026  
Jesus Labarta(jesus.labarta@bsc.es) , BSC

HORIZON-EUROHPC-JU-2023-COE



**EuroHPC**  
Joint Undertaking

Grant Agreement No 101143931

1 January 2024– 31 December 2026

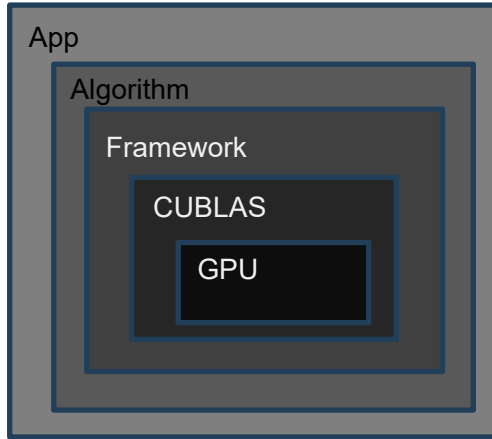
# POP: From HPC to AI ?



- HPC and AI ... speak a different language
  - Task, architecture, performance, ... same words, DIFFERENT meaning
  - Mainly strong scaling vs. mainly weak scaling
- But still, we should be able to ...
  - ... build on HPC performance analysis tools and methodologies ...
  - ... towards seamless analysis of AI apps from very large to very small scale

How efficiently are we using our resources?

- Black matryoshkas



- “no” visibility
  - What are we really doing?
  - How are we performing ?

- Just leave it to us !!
  - Our libraries, frameworks will do the BEST for you.
  - We do not offer mechanisms for you to take/steer our decision
- Why would you want explicit control of ...
  - Kernel variant selected , #blocks, threads
  - In GEMMs, NCCLs, ...

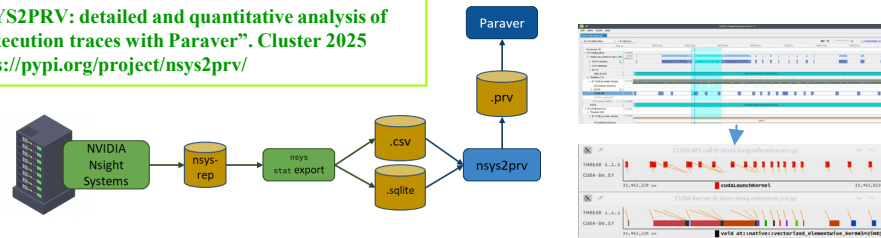
vendor, ...

## Assessments

- Important to understand ...
- ... and may be optimize our use case?
  - A bit more informed decisions than exhaustive parameter exploration

- Can we leverage Nsight Systems's acquisition capability ...
- and use Paraver to better squeeze the information in those traces ?
  - Scalability: dynamic range from very coarse to extremely fine grain
  - Analytics: timelines, histograms, correlations beyond profiles

M. Clascá et al. "NSYS2PRV: detailed and quantitative analysis of large-scale GPU execution traces with Paraver". Cluster 2025  
<https://pypi.org/project/nsys2prv/>



Towards microscopic insight ...

... from single GPU to "large" clusters



---

**A few tales ...**

# A few tales

---



- Long, long time ago, ...
- Common wisdom, real wisdom?
- Load imbalance in AI ?
- How important is communication ?
- What scares AI application engineers?
- You do not want that !!



---

# Long, long time ago

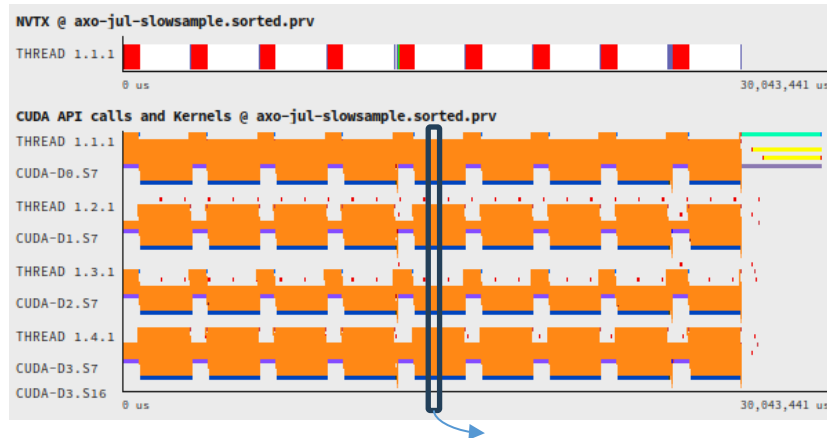
(24 months)

# in a far, far ...

# My first Paraver view ...



- ... of an LLM training (fine tuning)
- @ 4 GPUs



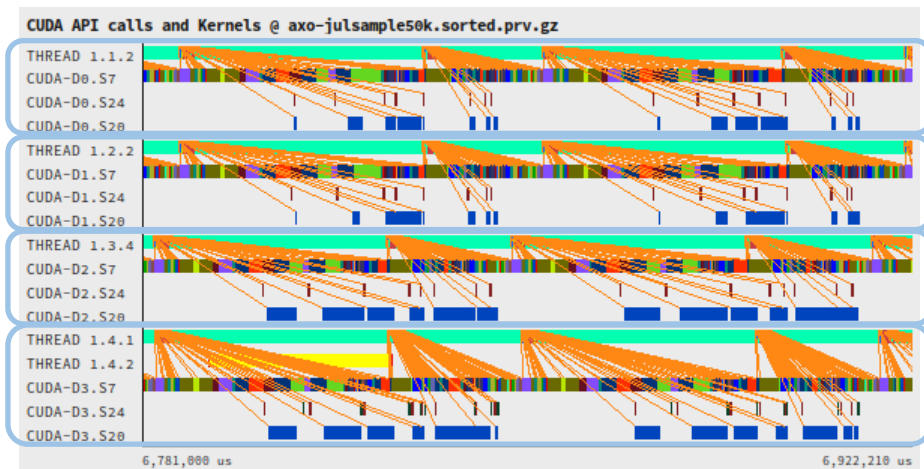
- End
- :DeepSpeedEngine.backward
- :DeepSpeedEngine.forward
- :step10
- :step11
- NCCL:nccLGroupStart

# API calls and kernels



- Kernels, CUDA calls, Launch lines

GPUs & task based models?  
Sparsity in performance data



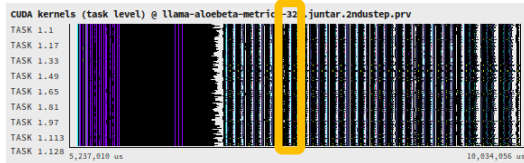
- End
- cudaEventDestroy
- cudaEventQuery
- cudaLaunchKernel
- cudaMemcpyAsync
- cudaMemsetAsync
- cudaStreamIsCapturing\_v10000
- cudaStreamSynchronize
- cudaStreamWaitEvent

- End
- ncclDevKernel\_AllReduce\_Sum\_bf16\_RING\_LL(ncclDevComm \*, unsigned long, ncclWork \*)
- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_nn\_n\_tilsize128x128x64\_warpgroupsz1x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas
- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_nt\_n\_tilsize256x128x64\_warpgroupsz2x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas
- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_nt\_n\_tilsize256x128x64\_warpgroupsz2x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas
- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_tn\_n\_tilsize128x128x64\_warpgroupsz1x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas
- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_tn\_n\_tilsize256x128x64\_warpgroupsz2x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas

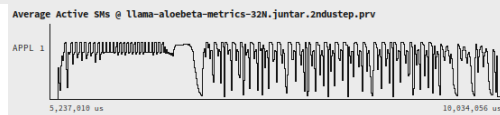
# @ 128 GPUs ?



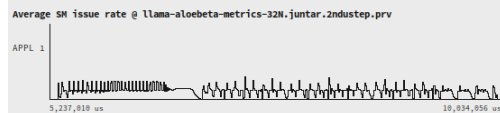
Kernels



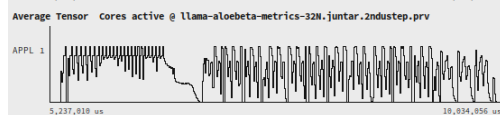
SM active



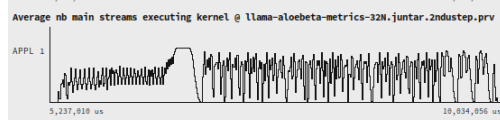
SM instr. Issue rate



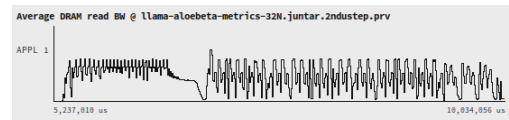
Tensor core active



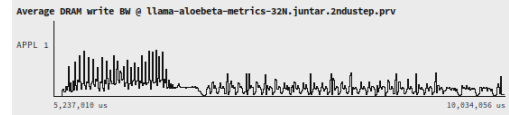
% active kernels



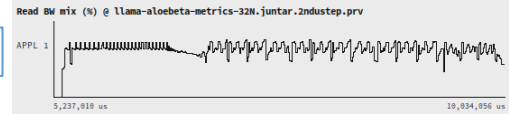
Read BW



Write BW



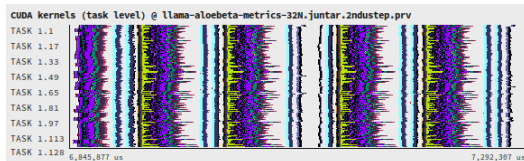
Read mix



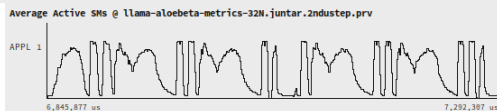
# @ 128 GPUs ?



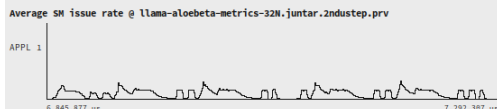
Kernels



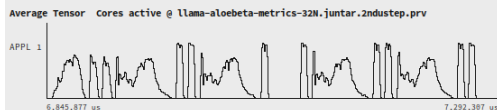
SM active



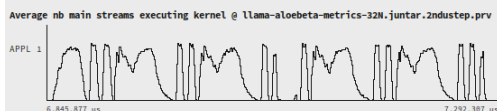
SM instr.  
Issue rate



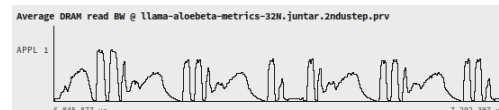
Tensor  
core active



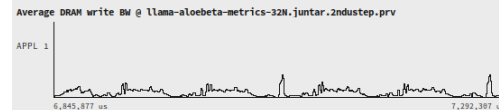
% active  
kernels



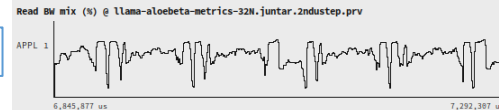
Read BW



Write BW



Read mix



Lots of opportunities ... who to “blame”?

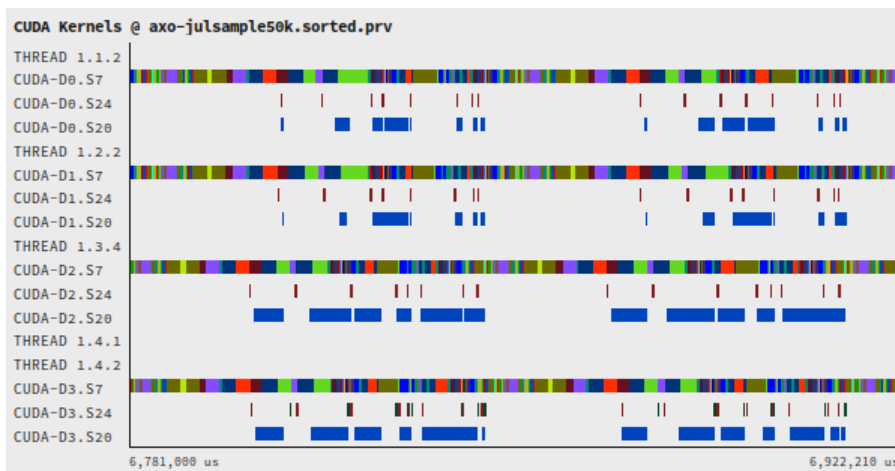


# Common wisdom ... real wisdom?

# Common wisdom (I)



- “AI == communication bound” ?



End

ncc1DevKernel\_AllReduce\_Sum\_bf16\_RING\_LL(ncclDevComm \*, unsigned long, ncclWork \*)

sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_nn\_n\_tilsize128x128x64\_warpgroupsize1x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas

sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_nt\_n\_tilsize256x128x64\_warpgroupsize2x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas

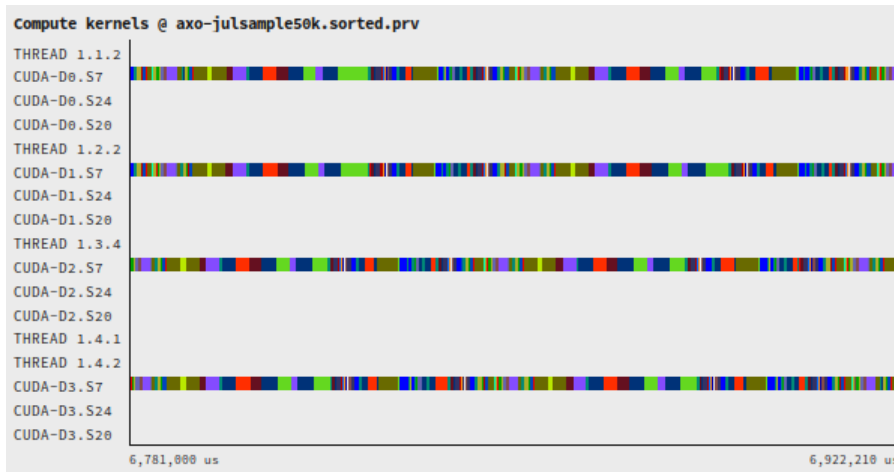
sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_tn\_n\_tilsize128x128x64\_warpgroupsize1x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas

sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_tn\_n\_tilsize256x128x64\_warpgroupsize2x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas

# Common wisdom (I)



- “AI == communication bound”
  - ... but this training case is certainly NOT communication bound

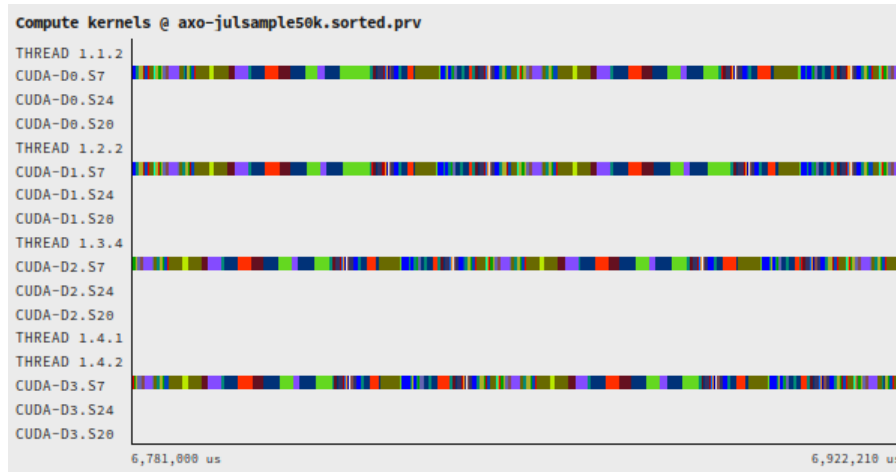


Idle %: 3.56

# Common wisdom (II)



- “AI == MxM”

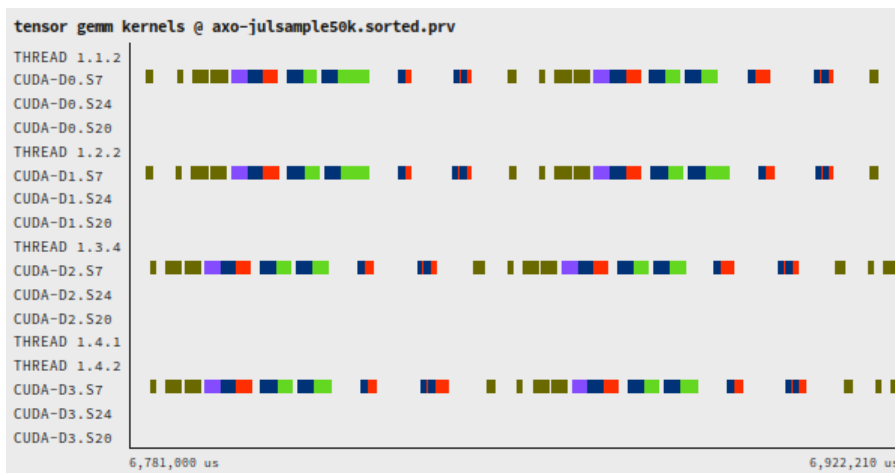


Idle %: 3.56

# Common wisdom (II)



- “AI == MxM”
  - ... at least in current platforms NOT the only/bottleneck



Idle %: 3.56

Gemm %: 49.95

- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_nn\_n\_tilesize128x128x64\_warpgroupsize1x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas
- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_nn\_n\_tilesize256x128x64\_warpgroupsize2x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas
- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_nt\_n\_tilesize256x128x64\_warpgroupsize2x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas
- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_tn\_n\_tilesize128x128x64\_warpgroupsize1x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas
- sm90\_xmma\_gemm\_bf16bf16\_bf16f32\_f32\_tn\_n\_tilesize256x128x64\_warpgroupsize2x1x1\_execute\_segment\_k\_off\_kernel\_\_5x\_cublas

# Common wisdom (II)



- “AI == MxM”
  - ... at least in current platforms NOT the only/bottleneck



Idle %: 3.56

Gemm %: 49.95

Non Gemm %: 46.49

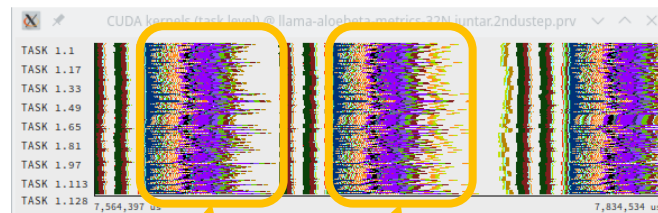


# Load imbalance in AI?

# @ 128 GPUs ?



- Kernels and duration



Load imbalance !!!  
Repetitive pattern

Load balance:  
77%

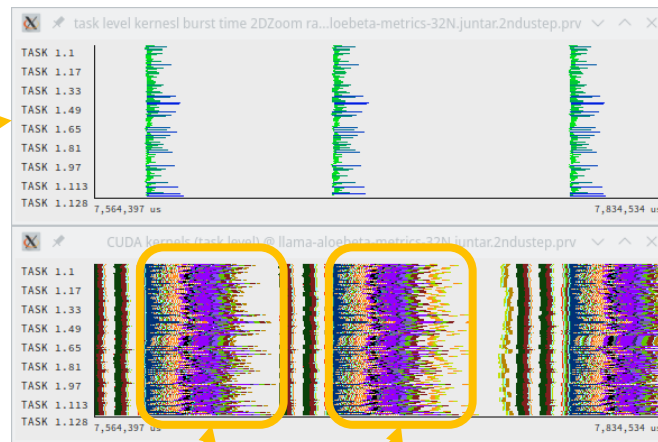
# @ 128 GPUs ?



- Kernels and duration

flash\_bwd\_dq... kernel duration  
Repetitive pattern

Load imbalance !!!  
Repetitive pattern



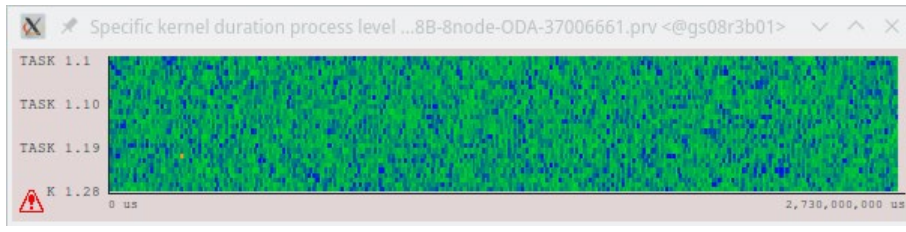
# Packing impact along time



- flash\_bwd\_... kernel duration

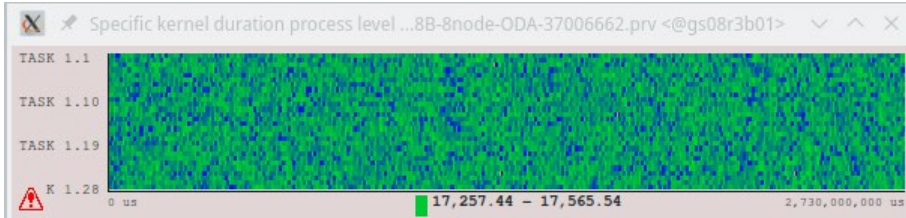
Sequence length 32K; 360 steps; 15 GB

Sequential



Stalls: 1.12%

MP default



0

flash\_bwd duration

2730 s

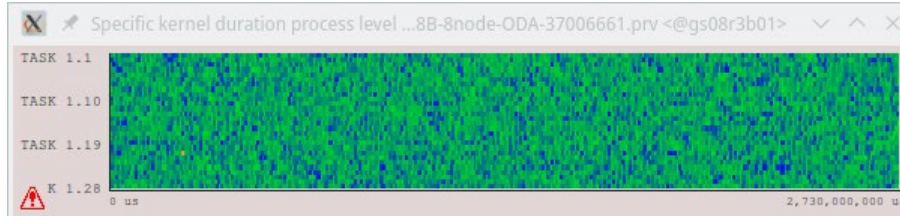
# Packing impact along time



- flash\_bwd\_... kernel duration

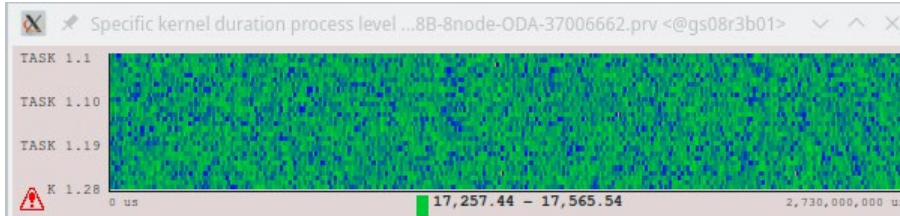
Sequence length 32K; 360 steps; 15 GB

Sequential

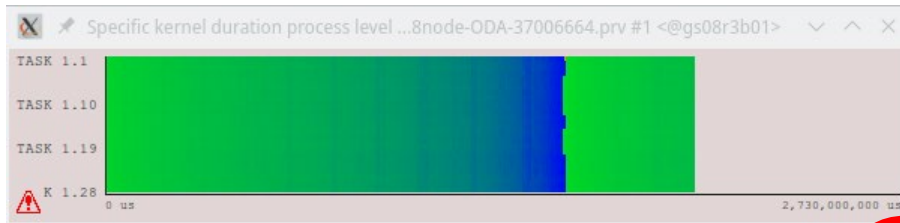


Stalls: 1.12%

MP default



Cycle MP Sort



Stalls: 2.56%

0

flash\_bwd duration

2730 s

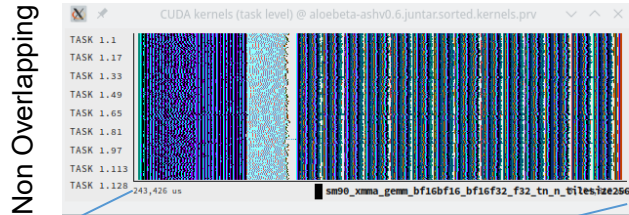


How important is “communication”?

# To overlap or not to overlap



- Different microscopic behaviors ... similar macroscopic duration !!!!

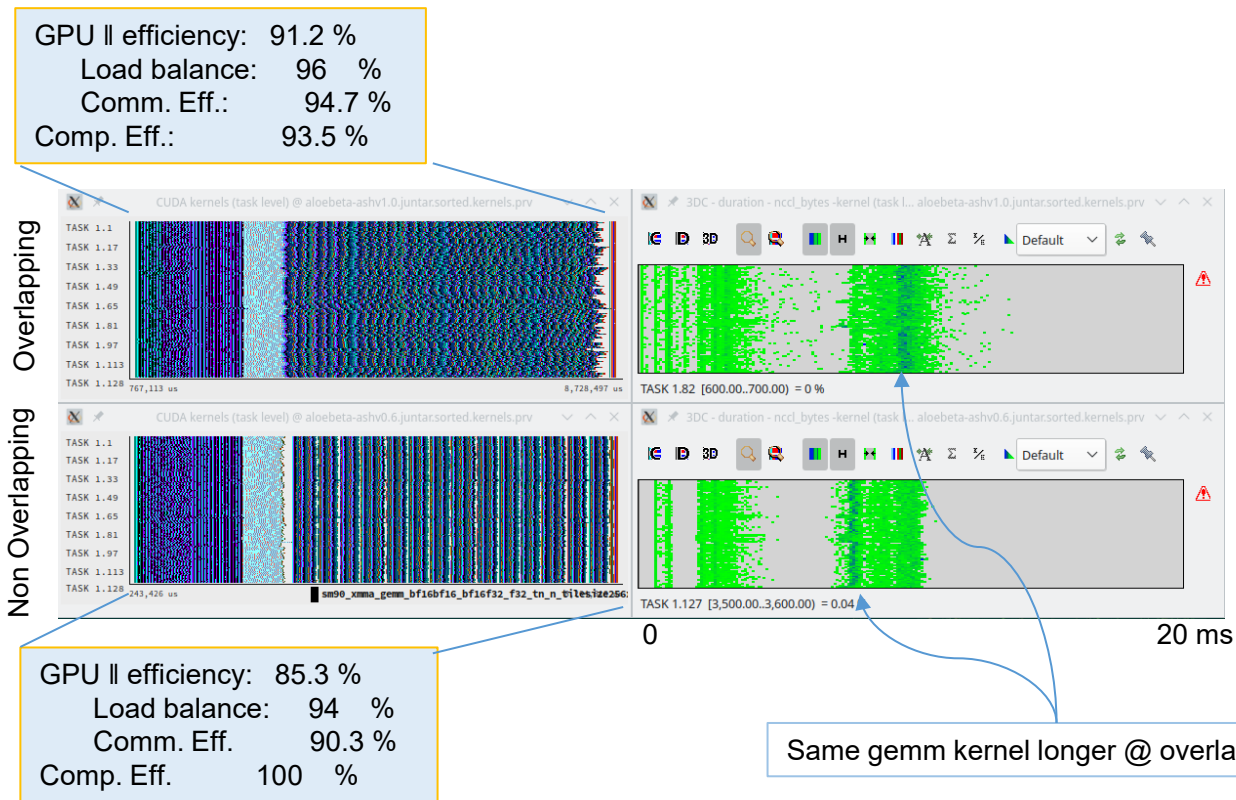


GPU II efficiency: 85.3 %  
Load balance: 94 %  
Comm. Eff. 90.3 %  
Comp. Eff. 100 %

# To overlap or not to overlap



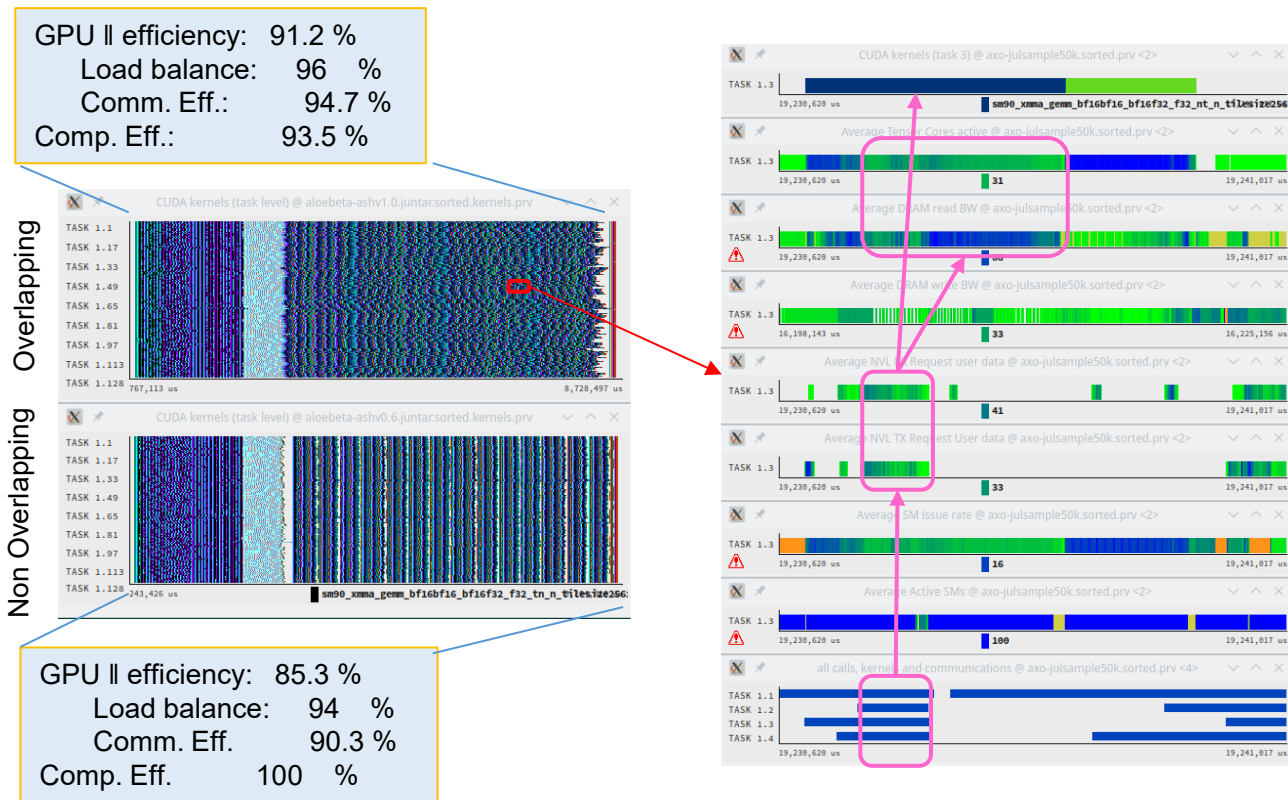
- Different microscopic behaviors ... similar macroscopic duration !!!!



# To overlap or not to overlap



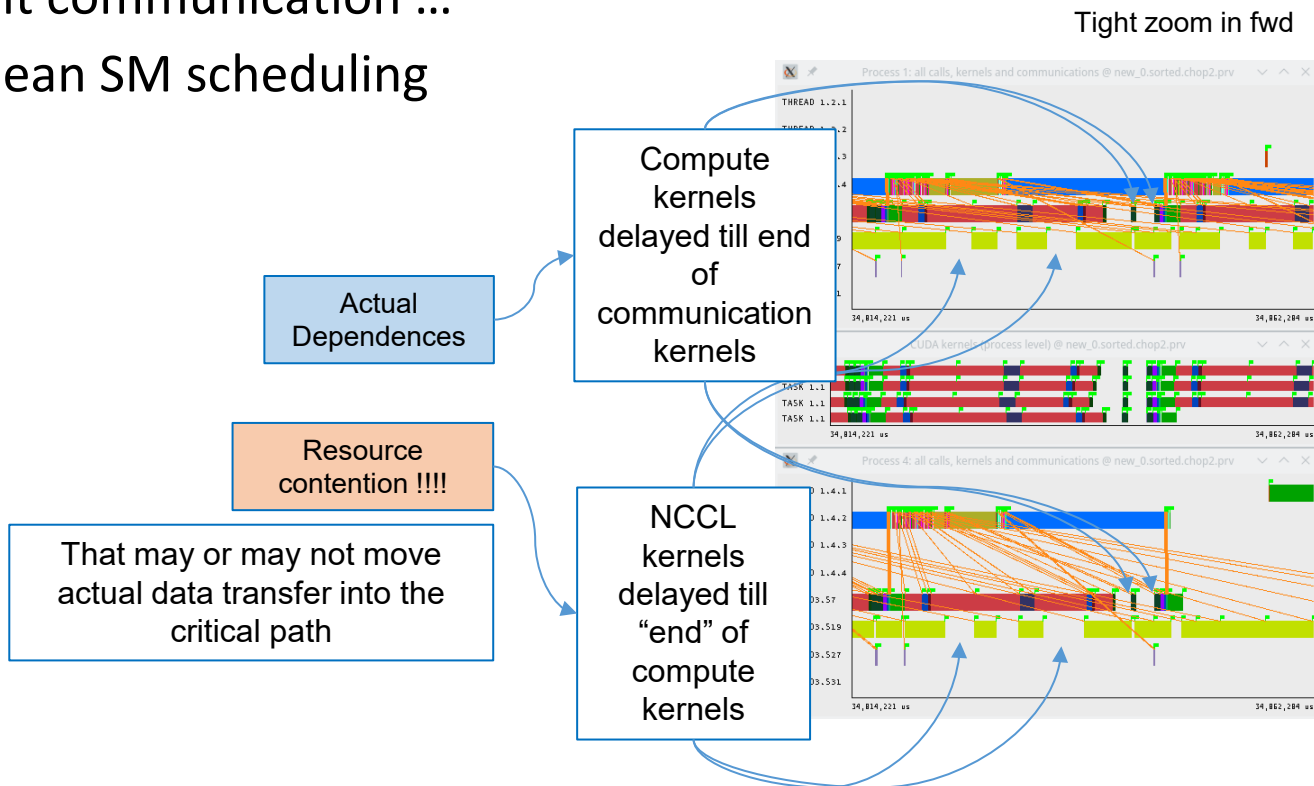
- Different microscopic behaviors ... similar macroscopic duration !!!



# Comp – Comm interferences ?



- Call it communication ...
- ... mean SM scheduling





To conclude ...

# Life in AI times

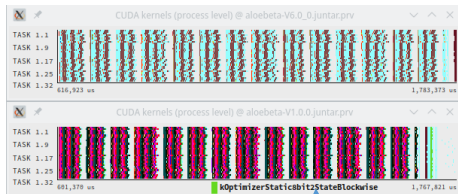


“My dear, here we must run as fast as we can, just to stay in place.  
And if you wish to go anywhere you must run twice as fast as that.”

Lewis Carroll. Through the Looking-Glass



From when you say is the framework version?  
Three months ago? ... ☹ ... update !!



Very different set of kernels

Some global gain

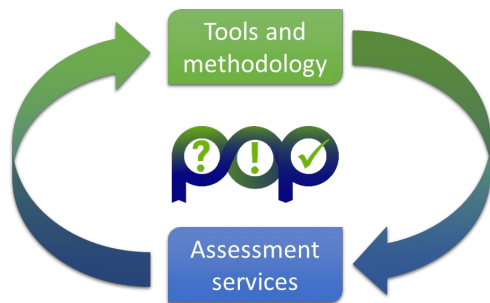
But same fundamental behavior

But ...

...we can leverage HPC  
Performance analysis tools,  
methodology and background  
to get external assessment of  
the behavior of AI apps !



- Performance optimization and productivity COE



## FREE Services provided by the CoE



### Parallel Application Performance Assessment

- Primary service
- Initial analysis measuring a [range of performance metrics](#) to assess quality of performance and identify the issues affecting performance (at customer site)
- If needed, undertakes further performance evaluations to identify the root causes of the issues found and qualify and quantify approaches to address them (recommendations)

### Second Level Services

- Second level services may follow after conclusion of an initial performance assessment:
  - Proof-of-concept:** explore the potential benefit of proposed optimisations by applying them to selected regions of the applications
  - Correctness study:** evaluate the correctness of hybrid MPI + OpenMP applications
  - Energy study:** investigate improvements of energy consumption or efficiency
  - Implementation consultancy:** providing consultancy for customers that choose to implement proposed optimisations

our experts and customer!

10



HORIZON-EUROHPC-JU-2023



Grant Agreement 10101743931

1 January 2024 – 31 December 2026

Give POP a try

Performance Optimisation and Productivity  
A Centre of Excellence in HPC

Services

Our experts provide performance optimisation and productivity services for (your?) academic AND industrial code(s) in all domains! We offer a portfolio of services designed to help our users optimise parallel software and understand performance issues. While our primary customers are code developers & owners our services are also available to code users and infrastructure & service centres.

The services are free of charge to academic, research or commercial organisations in the EU!

To discuss these services please get in touch with us via email ([pop@isc.eu](mailto:pop@isc.eu)).  
In order to request our services, use our [Request Service Form](#).

Services provided by POP

- Parallel Application Performance Assessment

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



# Performance Optimisation and Productivity 3

A Centre of Excellence in HPC

## Contact:



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



[pop@bsc.es](mailto:pop@bsc.es)



[@POP\\_HPC](#)



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

