



# Selected performance assessments using Scalasca/Score-P/CUBE

Brian Wylie (Jülich Supercomputing Centre)

HORIZON-EUROHPC-JU-2023-COE



**EuroHPC**  
Joint Undertaking

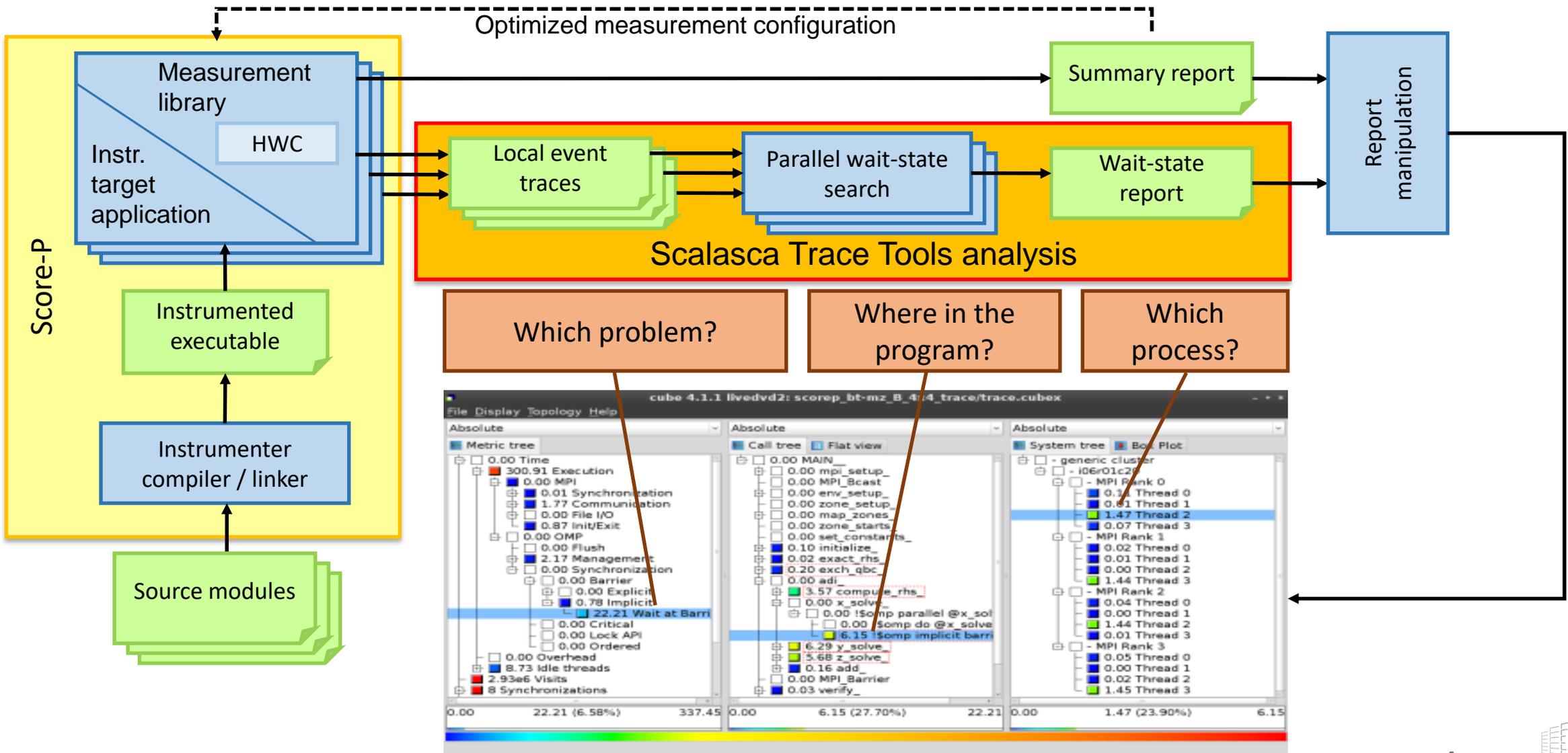
1 January 2024– 31 December 2026

Grant Agreement No 101143931

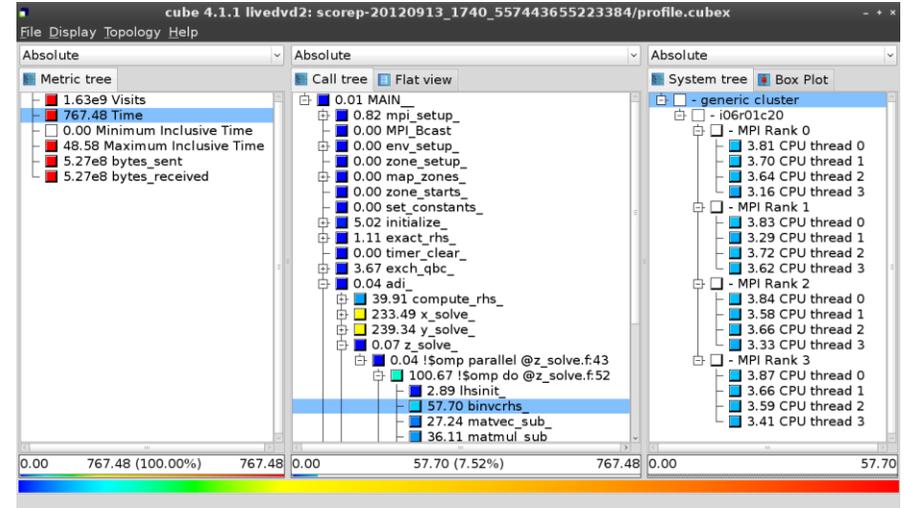
- Collection of trace-based performance tools
  - Specifically designed for large-scale systems
  - Features automatic trace analyzer providing wait-state, critical-path & delay analysis
  - Supports MPI, OpenMP, POSIX threads, and hybrid MPI+OpenMP/Pthreads
  - Uses Score-P instrumentation & measurement infrastructure and CUBE analysis report infrastructure
- Available under 3-clause BSD open-source license
- Documentation & sources:
  - <https://www.scalasca.org>
- Contact:
  - mailto: [scalasca@fz-juelich.de](mailto:scalasca@fz-juelich.de)

- Infrastructure for instrumentation and performance measurements
- Instrumented application can be used to produce several results:
  - Call-path profiling: CUBE4 data format used for data exchange
  - Event-based tracing: OTF2 data format used for data exchange
- Supported parallel paradigms:
  - Multi-process: MPI, SHMEM
  - Thread-parallel: OpenMP, Pthreads
  - Accelerator-based: CUDA, HIP, OpenCL, OpenACC, Kokkos
- Open Source; portable and scalable to all major HPC systems
- Initial project funded by BMBF
- Further developed in multiple third-party funded projects
- Documentation & sources: <https://www.score-p.org>

# Scalasca workflow



- Parallel program analysis report exploration tools
  - Libraries for Cube report reading & writing
  - Algebra utilities for report processing
  - GUI for interactive analysis exploration
- Available under 3-clause BSD open-source license
- Documentation & sources:
  - <http://www.scalasca.org>
- User guide also part of installation:
  - <prefix>/share/doc/CubeGuide.pdf
- Contact:
  - mailto: [scalasca@fz-juelich.de](mailto:scalasca@fz-juelich.de)



Assessment of execution efficiency factors using POP model

The screenshot displays the CubeGUI-4.5.0 interface for analyzing performance. The main window shows a hierarchical tree of metrics under 'Absolute' and 'Metric selection percent'. A 'POP Assessment' window is overlaid, providing a detailed breakdown of efficiency factors.

**POP Assessment : SimulationMaster::RunSimulation**

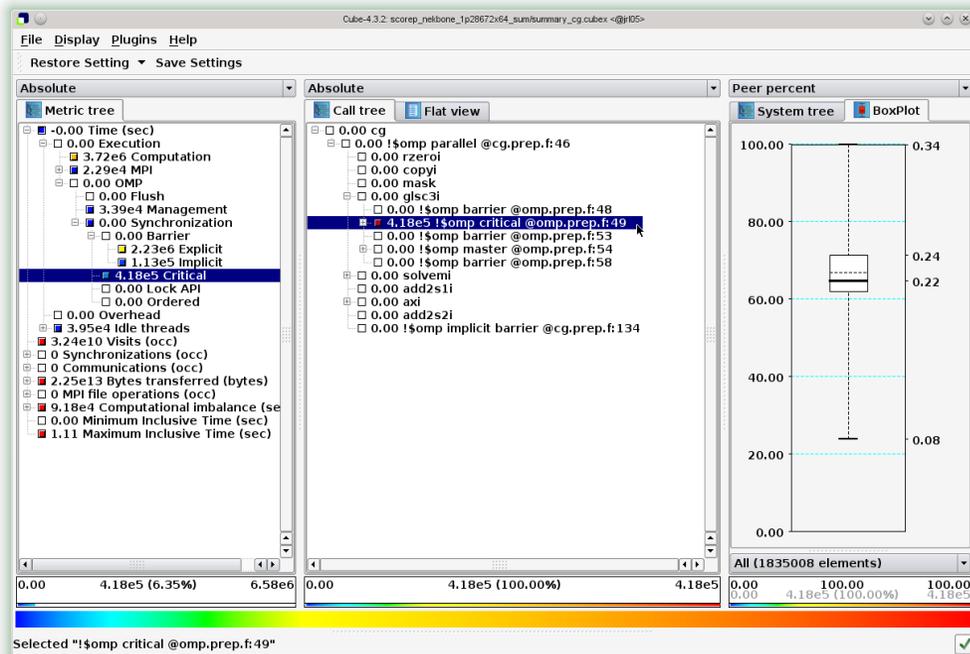
Parallel Efficiency	49.0 %		Fair	
Load balance	49.0 %		Fair	
Communication Efficiency	100.0 %		Very good	
Serialisation Efficiency	100.0 %		Excellent	
Transfer efficiency	100.0 %		Very good	
Stalled Resources	0.0 %		Excellent	
IPC	0.0		Value	
Instructions (only computation)				
Computation time	10377876.2		Value	

On the right, a 'Statistic...jureca' window shows a violin plot of data with a mean of 757.04 and a range from 518.59 to 1544.56. Below the plot, summary statistics are provided: 8.589 750.684 +/- 54.924 1544.563.

# Scalasca exa-scale readiness

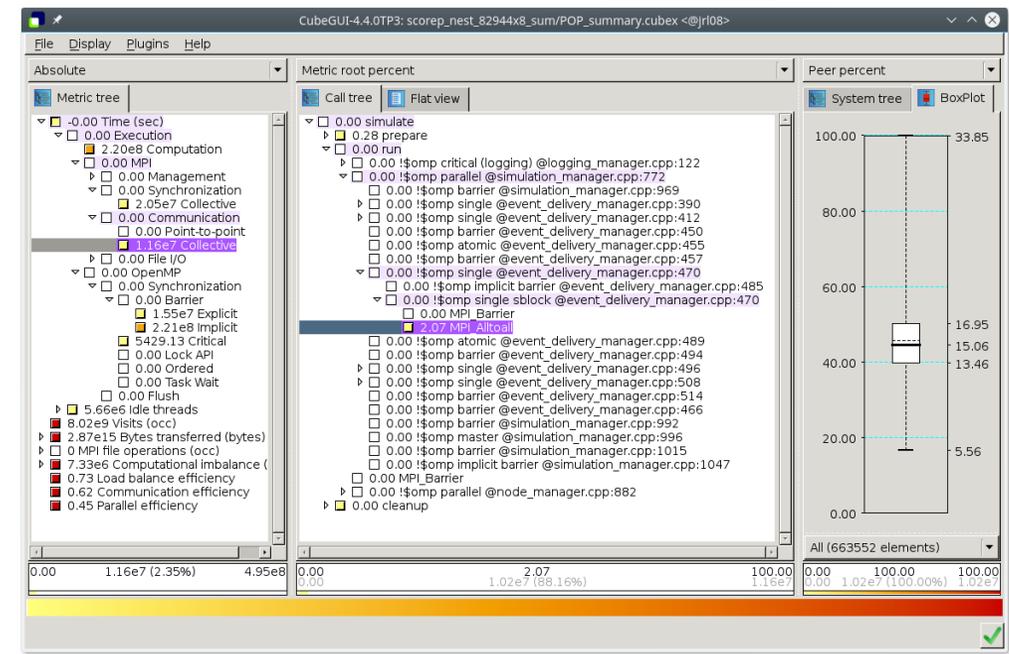


- Largest experiment
  - Application: NekBone
  - System: *JUQUEEN* IBM BG/Q
  - $28,672 \times 64 = 1,835,008$  threads



<https://co-design.pop-coe.eu/reports/POP-AR-112-Nekbone.html>

- Largest experiment by user
  - Application: NEST
  - System: *K* computer
  - $82,944 \times 8 = 663,552$  threads



[https://co-design.pop-coe.eu/reports/POP-AR-111-Nest\\_5g.html](https://co-design.pop-coe.eu/reports/POP-AR-111-Nest_5g.html)

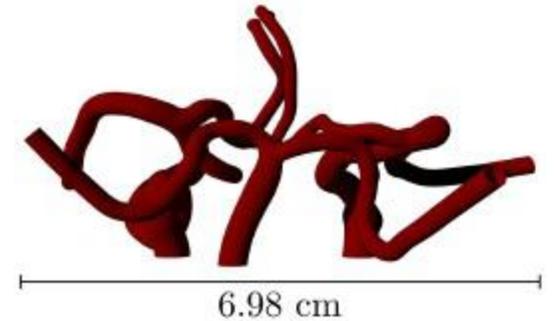
# Example performance assessments



- HemeLB (MPI) on *SuperMUC-NG*
  - also previously assessed on *ARCHER* Cray XC30 & *Blue Waters* Cray XE6
  - accelerated prototype (MPI+CUDA) assessed on *JUWELS-Volta* V100 GPUs
- SPECFEM3D (MPI+CUDA) on *Leonardo-B*
  - MPI version previously assessed on *Joloit-Curie*

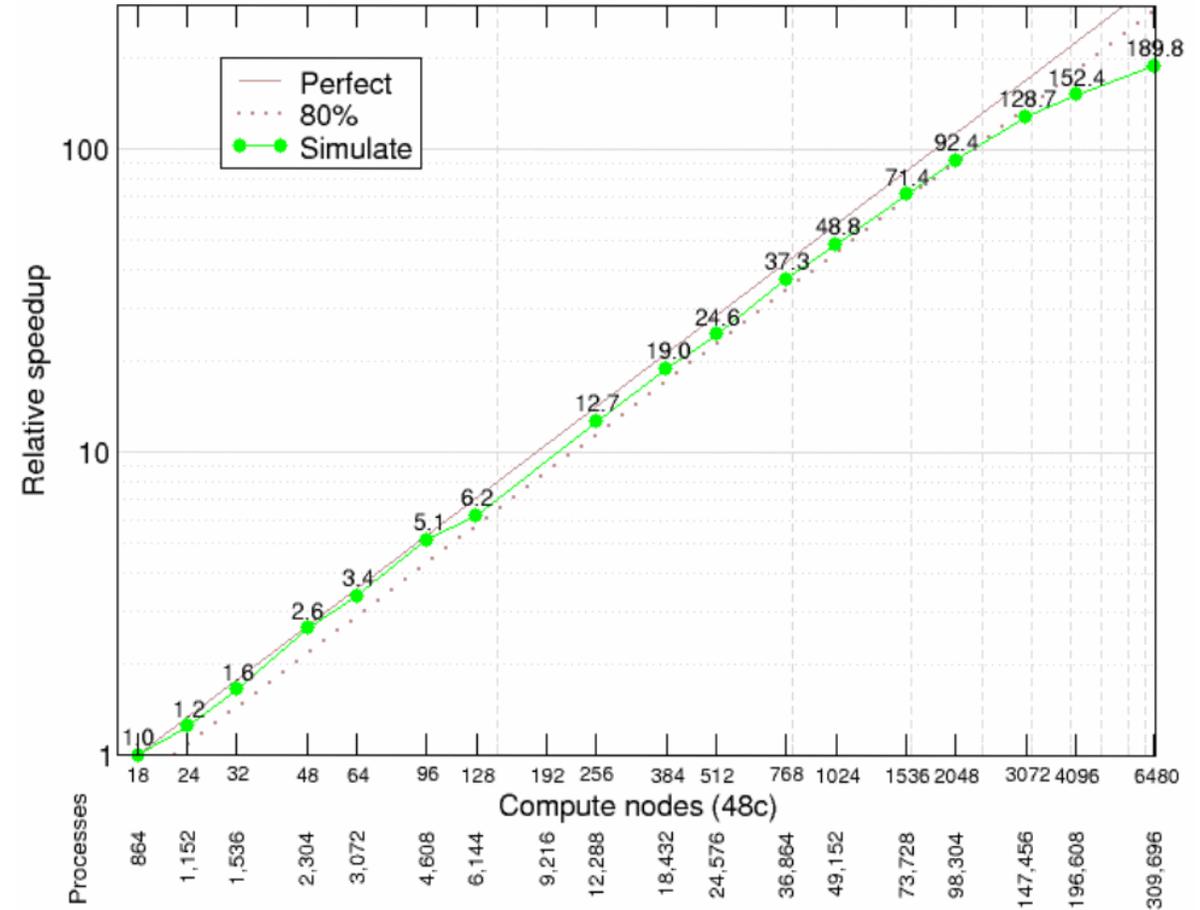
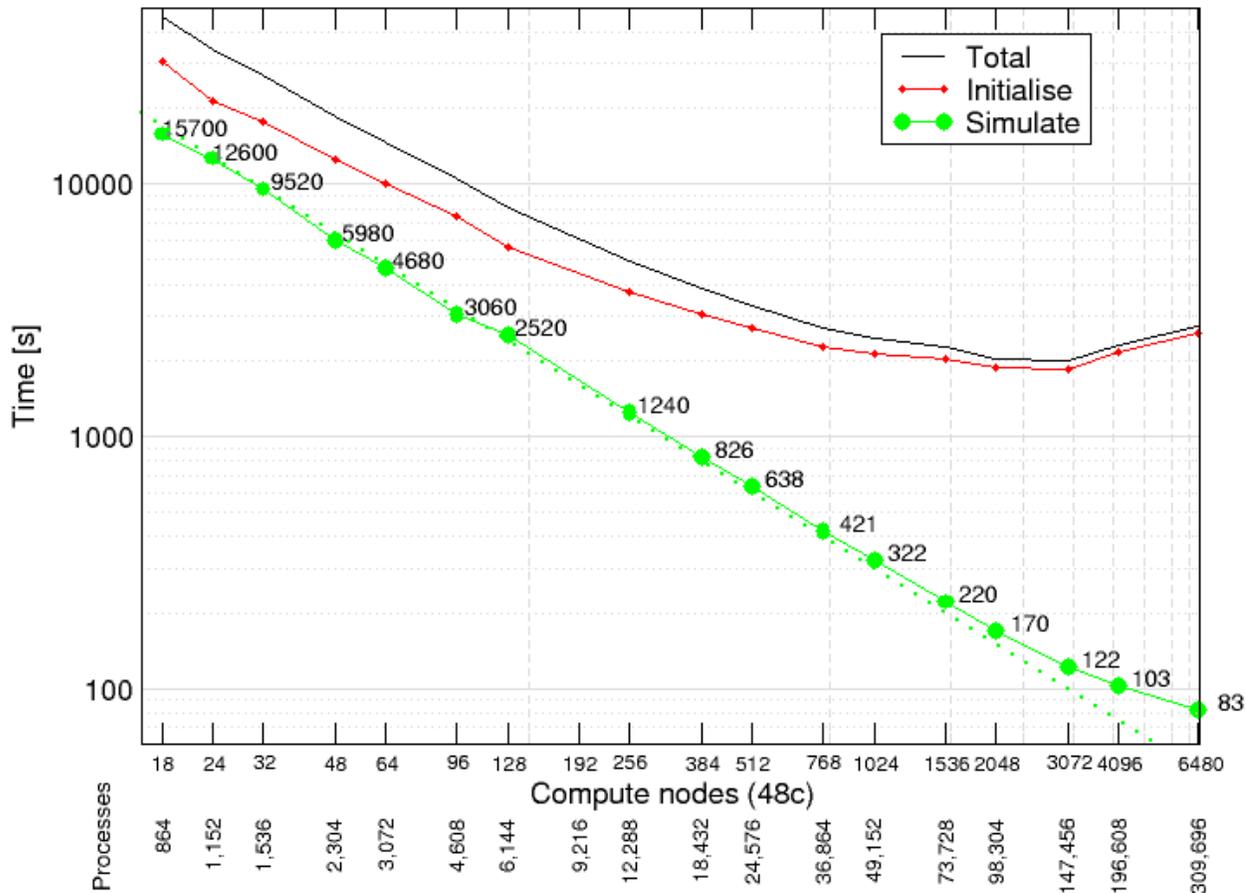


- 3D macroscopic blood flow in human arterial system developed by UC London (UK)
  - lattice-Boltzmann method tracking fluid particles on a lattice grid with complex boundary conditions
  - exascale flagship application of EU H2020 HPC Centre of Excellence for Computational Biomedicine (CompBioMed)
- HemeLB open-source code and test case: [www.hemelb.org](http://www.hemelb.org)
  - C++ parallelized with MPI
    - Intel Studio 2019u4 compiler and MPI library (v19.0.4.243)
    - configured with 2 ‘reader’ processes (intermediate MPI file writing disabled)
    - MPI-3 shared-memory model employed within compute nodes to reduce memory requirements when distributing lattice blocks from reader processes
  - Focus of analysis 5,000 time-step (500 $\mu$ s) simulation of cerebrovascular “circle of Willis” geometry
    - 6.4 $\mu$ m lattice resolution (21.15 GiB): 10,154,448,502 lattice sites
- Executed on *SuperMUC-NG* Lenovo ThinkSystem SD650 (LRZ):
  - 2x 24-core Intel Xeon Platinum 8174 (‘Skylake’) @ 3.1GHz
  - 48 MPI processes/node, 6452 (of 6480) compute nodes: 309,696 MPI processes
  - 190x speed-up from 864 cores: 80% scaling efficiency to over 100,000 cores



***⇒ Identification & quantification of impact of load balance and its variation***

# HemeLB@SNG strong scaling



[Execution of 9,216 processes on 192 compute nodes not possible due to insufficient compute nodes with adequate memory in 'fat' partition (768 GiB vs. regular 96 GiB node memory)]



# Advisor: POP efficiency assessment



CubeGUI-4.5.0: scorep\_hemehack\_13824\_sum.IPCX\summary.cubex

File Display Plugins Help

POP Assessment Runtime threshold: 2.5 % runtime

Absolute Metric tree

- 0.00 Time (sec)
  - 0.00 Execution
    - 6.10e7 Computation
      - 0.00 MPI
        - 60347.38 Management
        - 0.00 Synchronization
        - 0.00 Communication
          - 1.02e7 Point-to-point
            - 59875.44 Collective
            - 0.00 One-sided
          - 0.00 File I/O
        - 0.00 Overhead
      - 2.04e9 Visits (occ)
    - 0 Bytes transferred (bytes)
      - 0 Point-to-point
        - 2.70e14 Sent
        - 2.70e14 Received
      - 3.45e7 Collective
      - 0 Remote Memory Access
    - 0 MPI file operations (occ)
    - 553183.46 Computational imbalance (sec)
      - 2.72e17 PAPI\_TOT\_INS (#)
      - 1.63e17 PAPI\_TOT\_CYC (#)
      - 2.82e16 PAPI\_RES\_STL (#)
      - 5.43e16 CYCLE\_ACTIVITY:STALLS\_TOTAL (#)
      - 2.11e16 CYCLE\_ACTIVITY:STALLS\_MEM\_ANY (#)

Absolute Call tree Flat view

- 9.20e8 hemehack.scorep-trace
  - 9.57e8 main
    - 6.53e7 MPI\_Initialized
    - 1.64e11 MPI\_Init
    - 1.89e16 MEASUREMENT OFF
    - 2.22e15 RunSimulation
    - 5.32e10 ~SimulationMaster
    - 3.16e12 MPI\_Finalize

Advisor Score-P Configuration Source

Recalculate automatic

POP Assessment : SimulationMaster::RunSimulation

Parallel Efficiency	50%	Fair	0.50	?
Load balance	50%	Fair	0.50	?
Communication Efficiency	99%	Very good	1.00	?
Serialisation Efficiency	0%			!
Transfer efficiency	0%			!
Stalled Resources	57%	Fair	0.58	?
IPC		Value	1.56	?
Instructions (only computation)		Value	3.72e16	?
Computation time		Value	1.04e7	?

Candidates

Callpath	Issue
----------	-------

0 2.11e16 (100.00%) 2.11e16 0 2.22e15 (10.50%) 2.11e16

Ready



# HemeLB\_GPU (JUWELS-Volta)



[https://co-design.pop-coe.eu/reports/POP2-AR-065-HemeLB\\_GPU.html](https://co-design.pop-coe.eu/reports/POP2-AR-065-HemeLB_GPU.html)

DOI 10.5281/zenodo.4117942

DOI 10.5281/zenodo.4081080

- 3D macroscopic blood flow in human arterial system developed by UC London (UK)
  - lattice-Boltzmann method tracking fluid particles on a lattice grid with complex boundary conditions
  - exascale flagship application of EU H2020 HPC Centre of Excellence for Computational Biomedicine
- HemeLB open-source code and test case: [www.hemelb.org](http://www.hemelb.org)
  - C++ parallelized with MPI + CUDA (in development)
    - GCC/8.3.0 compiler, CUDA/10.1.105 and ParaStationMPI/5.4 library
    - configured with 2 'reader' processes and intermediate MPI file writing
    - rank 0 'monitor' process doesn't participate in simulation
  - Focus of analysis 2,000 time-step (each 100 $\mu$ s) simulation of CBM2019\_Arteries\_patched geometry
    - 1.78 GiB: 66,401,494 lattice sites, 1+38 iolets
- Executed on *JUWELS-Volta* (@JSC):
  - 2x 20-core Intel Xeon Platinum 8168 ('Skylake') CPUs + 4 Nvidia V100 'Volta' GPUs
  - 4\* MPI processes/node (one per GPU), 32 (of 56) compute nodes: 129 MPI processes



***⇒ Identification & quantification of impact of CPU/GPU load balance and its variation***

# Time for asynch. CUDA kernels



CubeGUI-4.5.0: scorep\_hemepure\_gpu\_20a+IO\_129\_trace\summary.cubex

File Display Plugins Help

**Absolute**

Metric tree

- 14974.45 Time (sec)
- 1.51e7 Visits (occ)
- 5.10e11 Bytes transferred (bytes)
- 202268 MPI file operations (occ)
- 269.69 Computational imbalance (sec)
- 1.86e9 io\_bytes\_read (bytes)
- 7.44e9 io\_bytes\_written (bytes)

**Absolute**

Call tree Flat view

- 0.83 hemepure\_gpu\_20a.scorep
  - 0.02 main
    - 0.00 MPI\_Initialized
    - 122.79 MPI\_Init
    - 12890.80 SimulationMaster
      - 1324.46 RunSimulation
      - 193.65 ~SimulationMaster
      - 8.91 MPI\_Finalize
    - 0.40 BUFFER\_FLUSH
    - 396.53 hemelb::GPU\_CollideStream\_mMidFluidCollision\_mWallCollision\_sBB
    - 13.12 hemelb::GPU\_CollideStream\_iolets\_NashZerothOrderPressure
    - 14.74 hemelb::GPU\_CollideStream\_wall\_sBB\_iolet\_Nash
    - 8.21 hemelb::GPU\_StreamReceivedDistr

**Absolute**

System tree Statistics

- machine Linux
  - node jwc09n033.adm09.juwels.fzj.de
    - MPI Rank 0
      - 11.73 Master thread
    - MPI Rank 1
      - 9.48 Master thread
        - 0.00 CUDA[0:7]
        - 0.03 CUDA[0:17]
        - 2.22 CUDA[0:23]
        - 0.41 CUDA[0:26]
        - 0.41 CUDA[0:28]
        - 0.08 CUDA[0:31]
    - MPI Rank 2
      - 9.48 Master thread
        - 0.00 CUDA[1:7]
        - 0.05 CUDA[1:17]
        - 2.34 CUDA[1:23]
        - 0.36 CUDA[1:26]
        - 0.35 CUDA[1:28]
        - 0.01 CUDA[1:31]
    - 11.99 MPI Rank 3
    - 11.92 MPI Rank 4
  - node jwc09n036.adm09.juwels.fzj.de
    - 12.48 MPI Rank 5
    - 11.23 MPI Rank 6
    - 11.35 MPI Rank 7
    - 14.40 MPI Rank 8
  - node jwc09n039.adm09.juwels.fzj.de
    - 12.60 MPI Rank 9
    - 13.32 MPI Rank 10
    - 11.54 MPI Rank 11
    - 14.03 MPI Rank 12

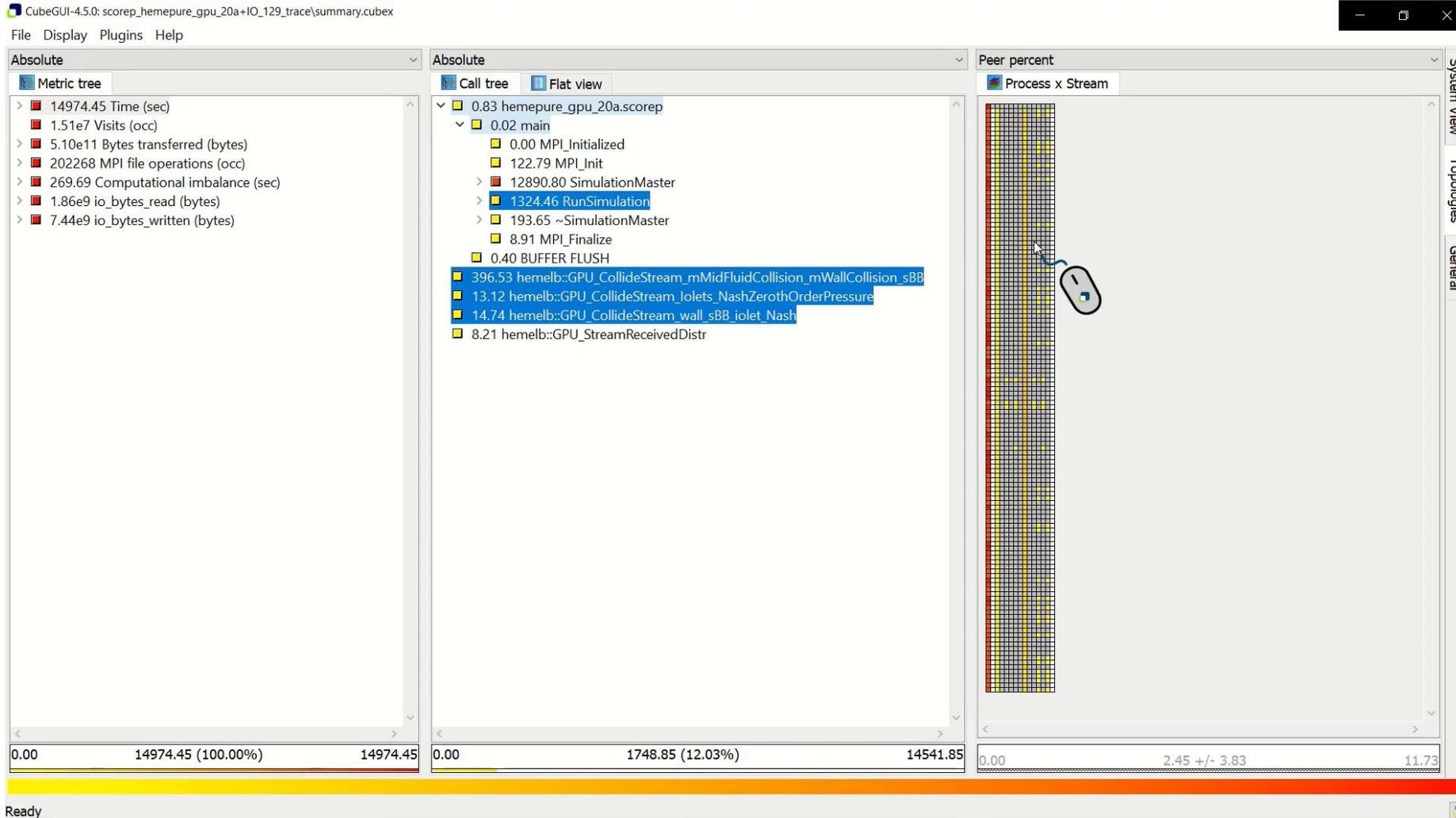
System View Topologies General

All (715 elements)

0.00	14974.45 (100.00%)	14974.45	0.00	1757.06 (12.08%)	14541.85	0.00	1757.06
------	--------------------	----------	------	------------------	----------	------	---------

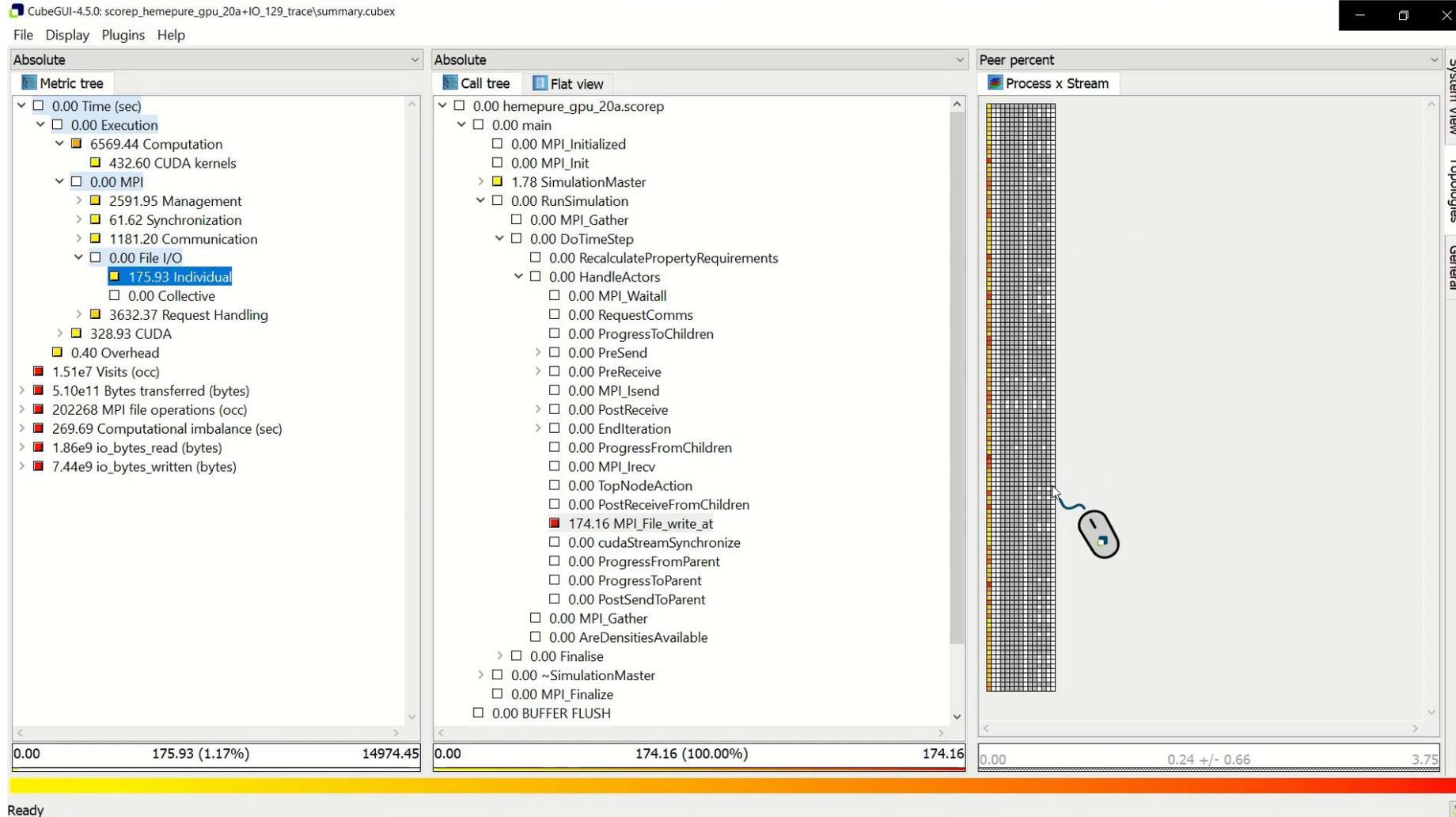
Ready

# Time for asynch. CUDA kernels



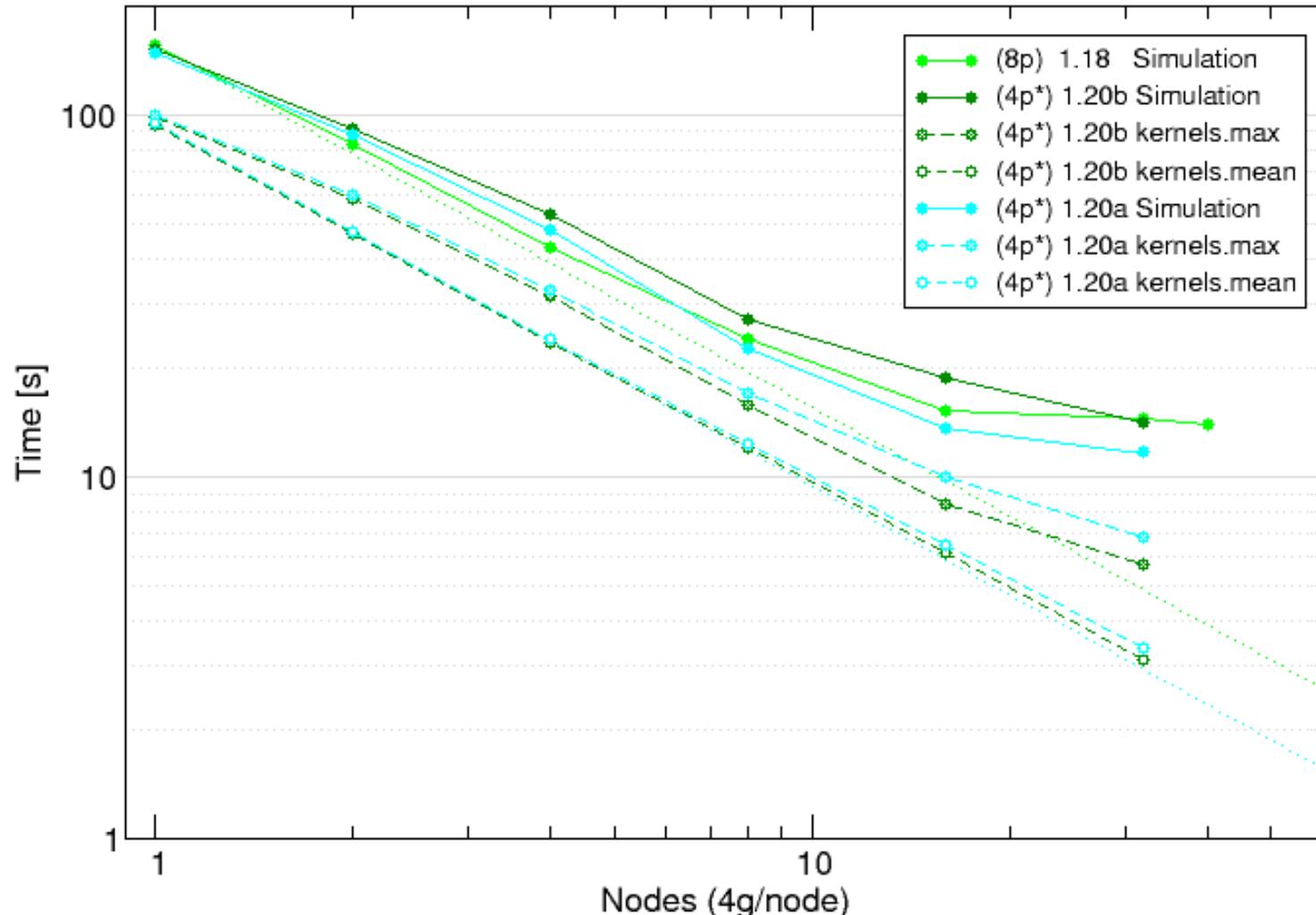


# Time for MPI file writing on CPU



Ready

# HemeLB@JUWELS-Volta strong scaling

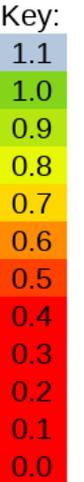


- Reference execution with 8ppn
  - multiple processes offloading GPU kernels generally unproductive
- Comparison of versions (4ppn)
  - v1.20a generally better
- Synchronous MPI file writing is the primary bottleneck
- CUDA kernels on GPUs
  - less than half of Simulation time (therefore GPUs mostly idle)
  - total kernel time scales very well (0.93 scaling efficiency)
  - load balance deteriorates (0.95 for single node, 0.50 for 32 nodes)

# HemeLB@JUWELS-V scaling efficiency



	1n 5p	2n 9p	4n 17p	8n 33p	16n 65p	32n 129p
Simulation time [s]	147.87	88.38	48.13	22.66	13.68	11.67
Global scaling efficiency	0.64	0.53	0.49	0.52	0.43	0.25
– Parallel efficiency	0.64	0.53	0.50	0.54	0.47	0.29
– – Load balance efficiency (GPU)	0.95	0.78	0.73	0.73	0.65	0.50
– – Communication efficiency (GPU)	0.67	0.68	0.68	0.75	0.73	0.58
– Computation scaling (GPU)	1.00	1.00	0.99	0.96	0.92	0.87

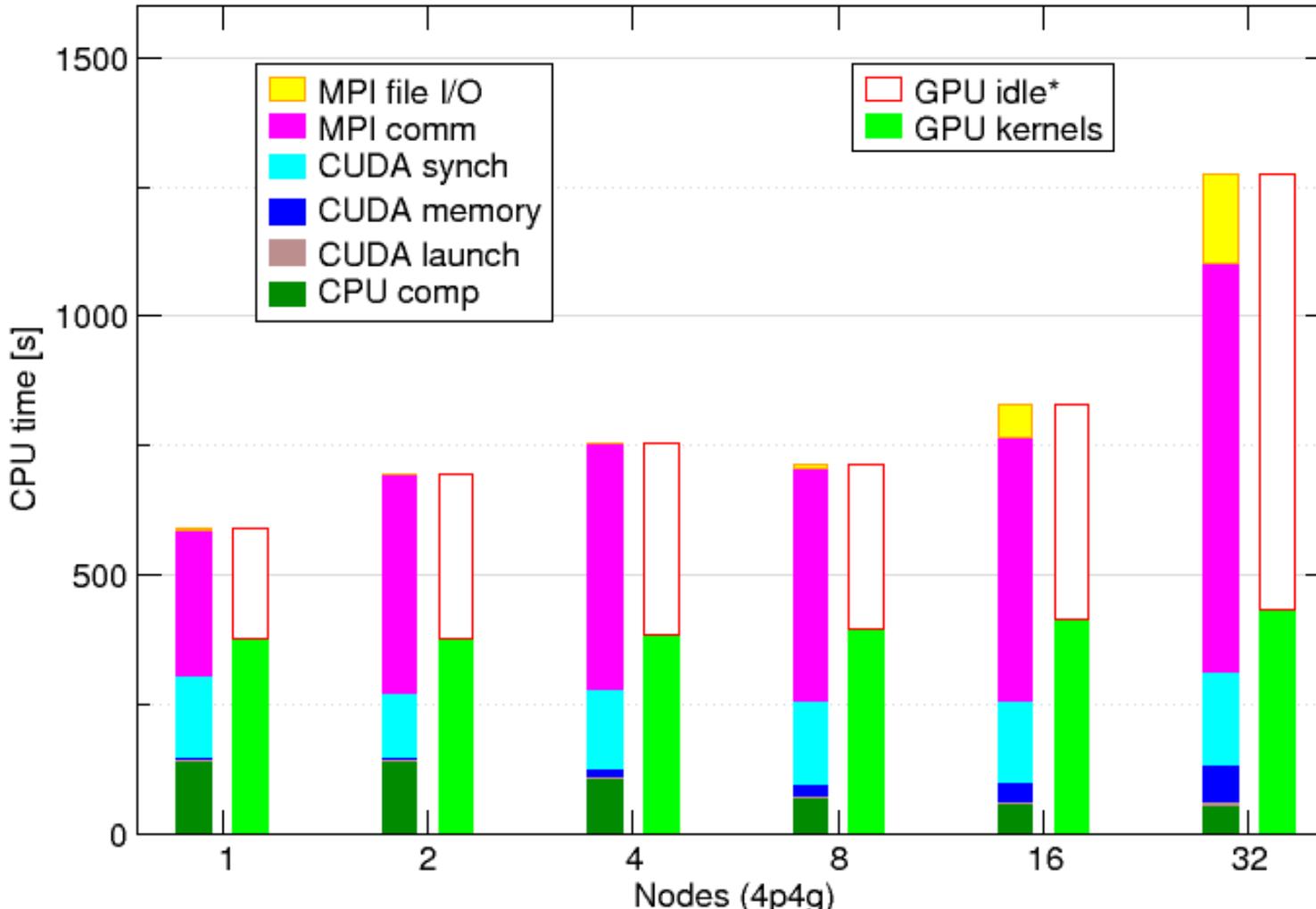


Only considering GPUs (ignoring all CPU cores, 90% of which are completely unused)

- Single (quad-GPU) node already suffers significant communication inefficiency
  - includes MPI file writing, but doesn't degrade much as additional nodes are included
- Load balance of GPUs deteriorates progressively
- GPU computation scaling remains reasonably good

[POP CoE scaling efficiency model: [www.pop-coe.eu](http://www.pop-coe.eu)]

# HemeLB@JUWELS-V strong scaling



- CPU+GPU time breakdown
- CUDA kernels on GPUs
  - less than half of Simulation time (therefore GPUs mostly idle)
  - total kernel time scales very well (0.87 scaling efficiency)
- MPI processes on CPUs
  - computation time decreases
  - CUDA synchronization time fairly constant, but time for memory management increases somewhat
  - MPI communication time dominates, with much more time for file writing with 16+ nodes

# Example performance assessments



- HemeLB on *SuperMUC-NG* (MPI)
  - also previously assessed on *ARCHER* Cray XC30 & *Blue Waters* Cray XE6
  - accelerated prototype (MPI+CUDA) assessed on *JUWELS-Volta* V100 GPUs
- SPECFEM3D on *Leonardo-B* (MPI+CUDA)
  - also previously assessed on *Joliot-Curie*



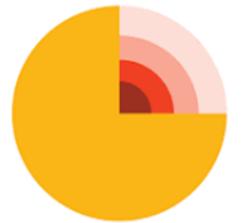
# SPECFEM3D (Leonardo-B)



<https://co-design.pop-coe.eu/reports/POP3-AR-002-SPECFEM3D.html>

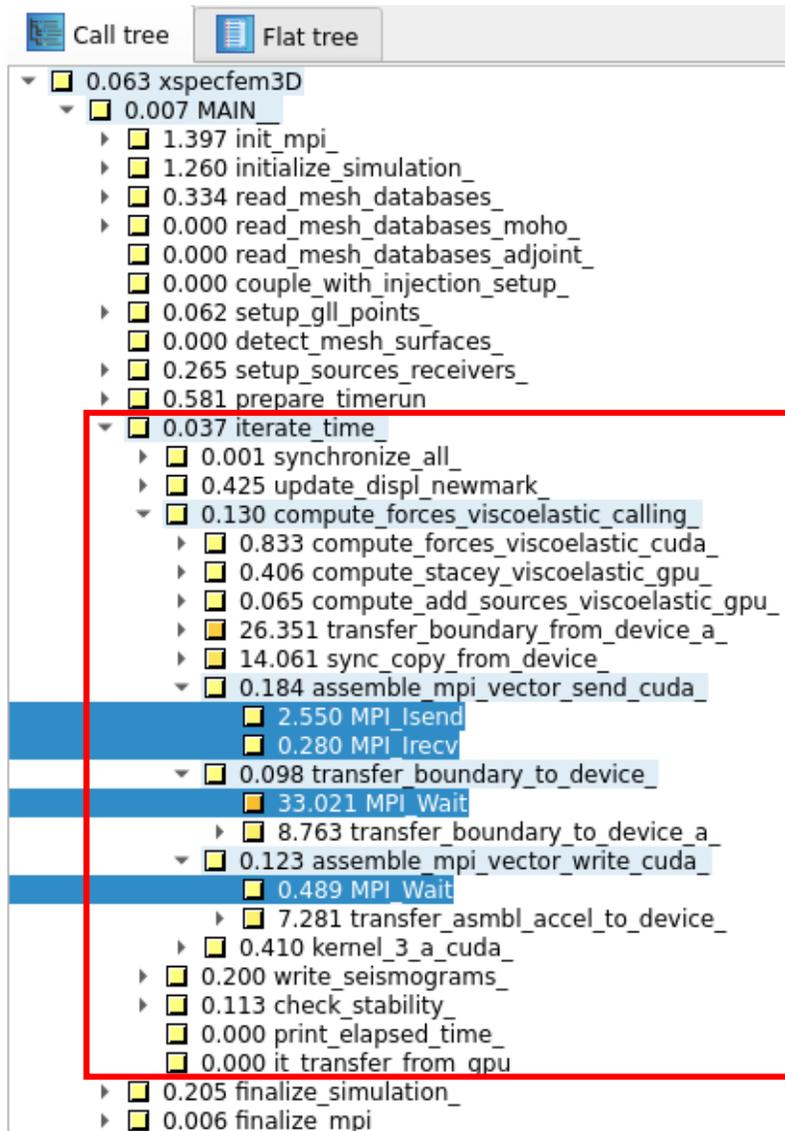
DOI 10.5281/zenodo.13643996

- SPECFEM3D
  - Software package for simulation of seismic wave propagation based on the spectral-element method
  - Assessment for HPC CoE for Exascale in Solid Earth (ChEESE)
  - Version 4.0.0 (release)
- Fortran90 (and some C) parallelized with MPI & CUDA: one MPI process per GPU
  - Intel oneAPI 2023.0.0 compilers and 2021.7.1 MPI libraries (not GPU-Aware)
- Testcase: 1 source in elastic domain; 4 seismic receiver stations
  - 48000 solver timesteps with intermediate writing disabled
  - weak scaling (22x 128x128 = 360,448 elastic elements per rank)
  - strong scaling (22x 1024x1024 = 23,068,672 elastic elements total)
- Executed on *Leonardo-Booster* Atos Bull Sequana XH21355 (CINECA)
  - 2345 compute nodes with 32-core Intel Xeon Platinum 8358 ('IceLake') CPUs @ 2.6 GHz and quad Nvidia A100 ('Ampere') GPUs [64GB]
  - Nvidia Mellanox HDR DragonFly++ interconnection network
- Measurements with Scalasca/2.6.1 using Score-P/8.3
- Focus of analysis (FOA): *xspecfem3D/iterate\_time*



ChEESE

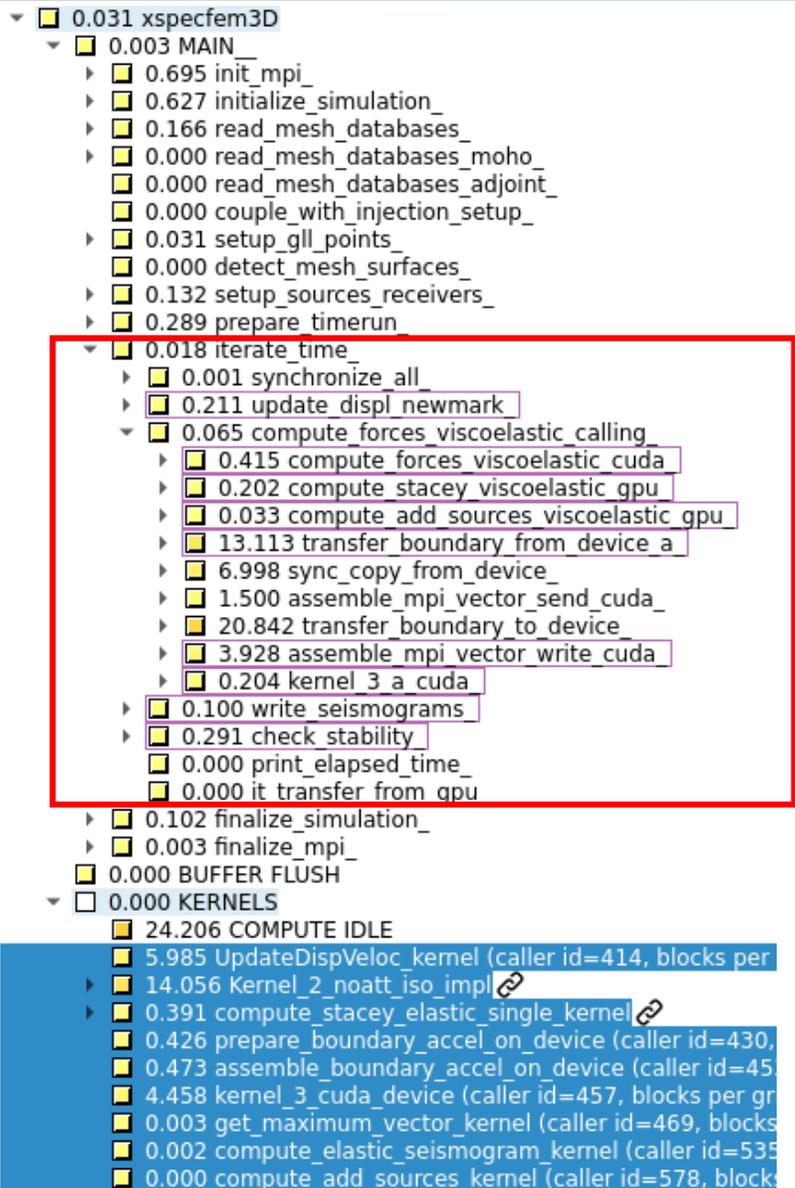
# Execution call-tree & Focus of Analysis



## • Structure

- setup/initialise (amortised in full run)
  - read (w/o MPI), MPI\_Bcast, MPI\_Reduce, etc
- solver (*iterate\_time*)
  - 48000 timesteps
    - non-blocking point-to-point communication for boundary exchange with 2D neighbours
    - data transfer to/from associated GPU device and corresponding stream synchronization
  - summary output every 500 steps
    - collective MPI\_Reduce
  - *write\_seismograms* executed only once
- Focus of Analysis selected for assessment: *iterate\_time*

# Execution call-tree & kernels



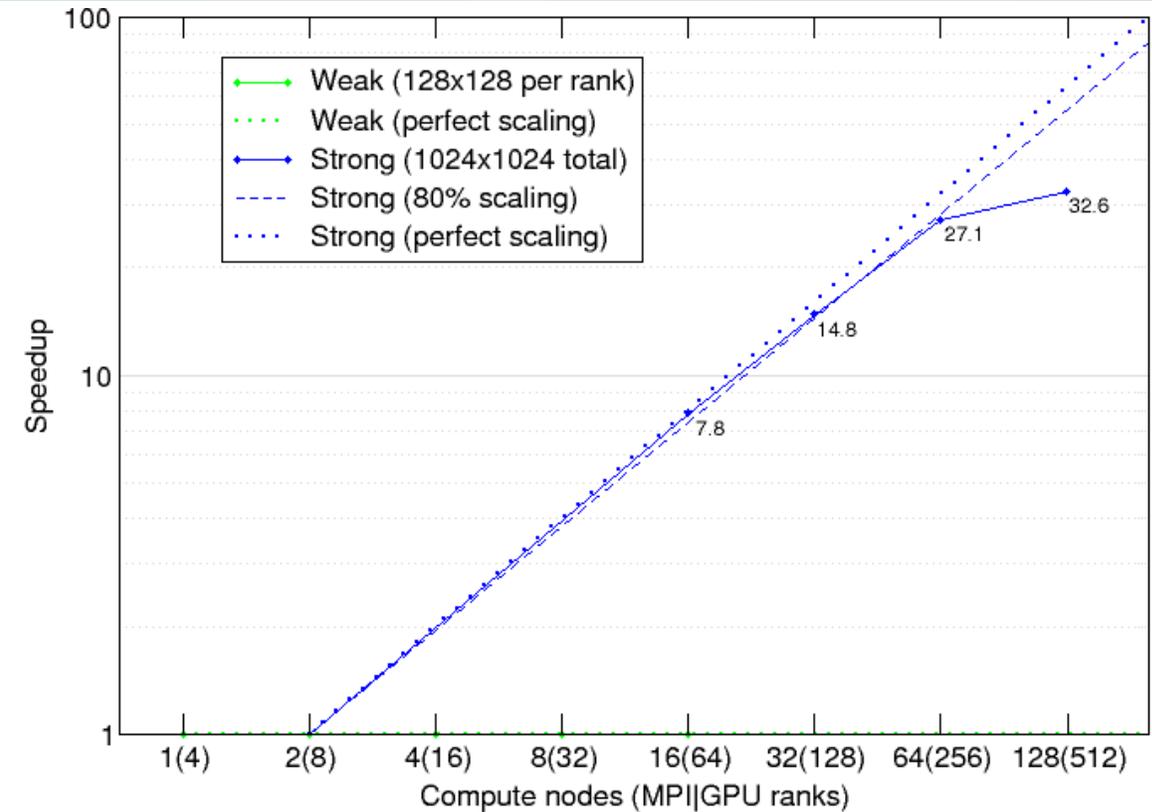
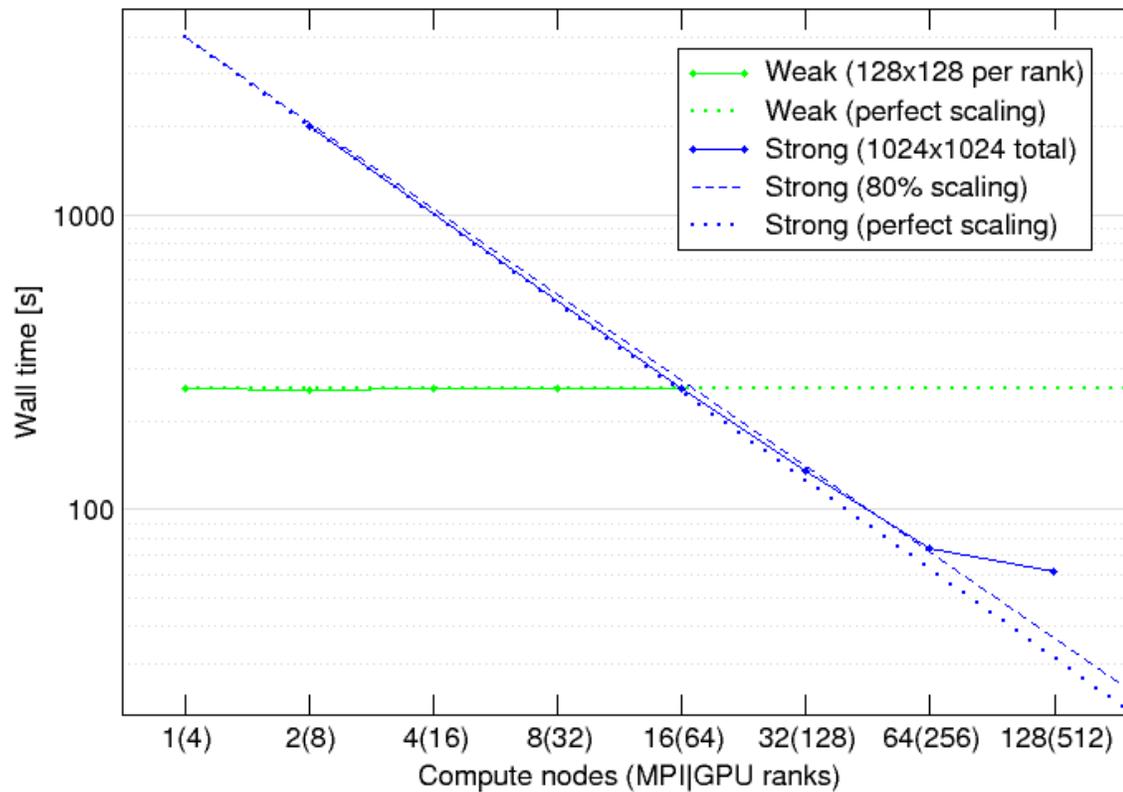
- Structure

- solver (*iterate\_time*)

- contains all of the CUDA kernel executions
      - GPU devices idle during init/finalize
    - 48000 timesteps
      - seven of nine kernels executed by all ranks
      - characteristics often determined by position in 2D grid
    - *compute\_add\_sources\_kernel* only executed by a single GPU (rank 243 of 512)
    - *compute\_elastic\_seismogram* only by 4 nearby GPUs (ranks 241, 245, 273, 277 of 512)

- Focus of Analysis selected for assessment: *iterate\_time*

# Scaling & speed-up



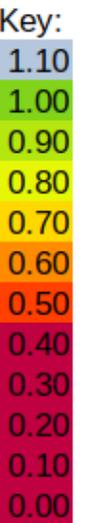
- xspecfem3d FOA *iterate\_time* on Leonardo-Booster
- Excellent weak scaling (expected to continue to higher node counts)
- Very good strong scaling (above 80% of perfect) to around 64 compute nodes (256 GPUs)

# Efficiency model (weak scaling)



Problem size		256x256	512x256	512x512	1024x512	1024x1024
MPI GPU ranks		4	8	16	32	64
Wall time [s]		255.873	255.462	256.035	255.898	255.948
<b>XPU</b>	Global scaling efficiency	0.514	0.516	0.514	0.514	0.514
	- Computation time scaling	1.000	0.992	0.977	0.972	0.963
	- Parallel efficiency	0.514	0.520	0.526	0.529	0.534
	- - Load balance efficiency	0.522	0.526	0.533	0.536	0.541
	- - Orchestration efficiency	0.986	0.988	0.988	0.988	0.988
<b>GPU</b>	Global scaling efficiency	0.986	0.987	0.985	0.986	0.986
	- Computation time scaling	1.000	1.000	1.000	1.000	1.000
	- Parallel efficiency	0.986	0.987	0.985	0.986	0.986
	- - Load balance efficiency	1.000	0.999	0.997	0.998	0.998
	- - Orchestration efficiency	0.986	0.988	0.988	0.988	0.988

- Orchestration efficiency
  - MPI communication & CUDA management
  - aka Communication efficiency



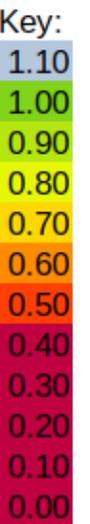
- Excellent GPU weak scaling efficiency
- Very poor CPU efficiency?
- Moderate XPU (GPU+CPU) efficiency?

# Efficiency model (strong scaling)



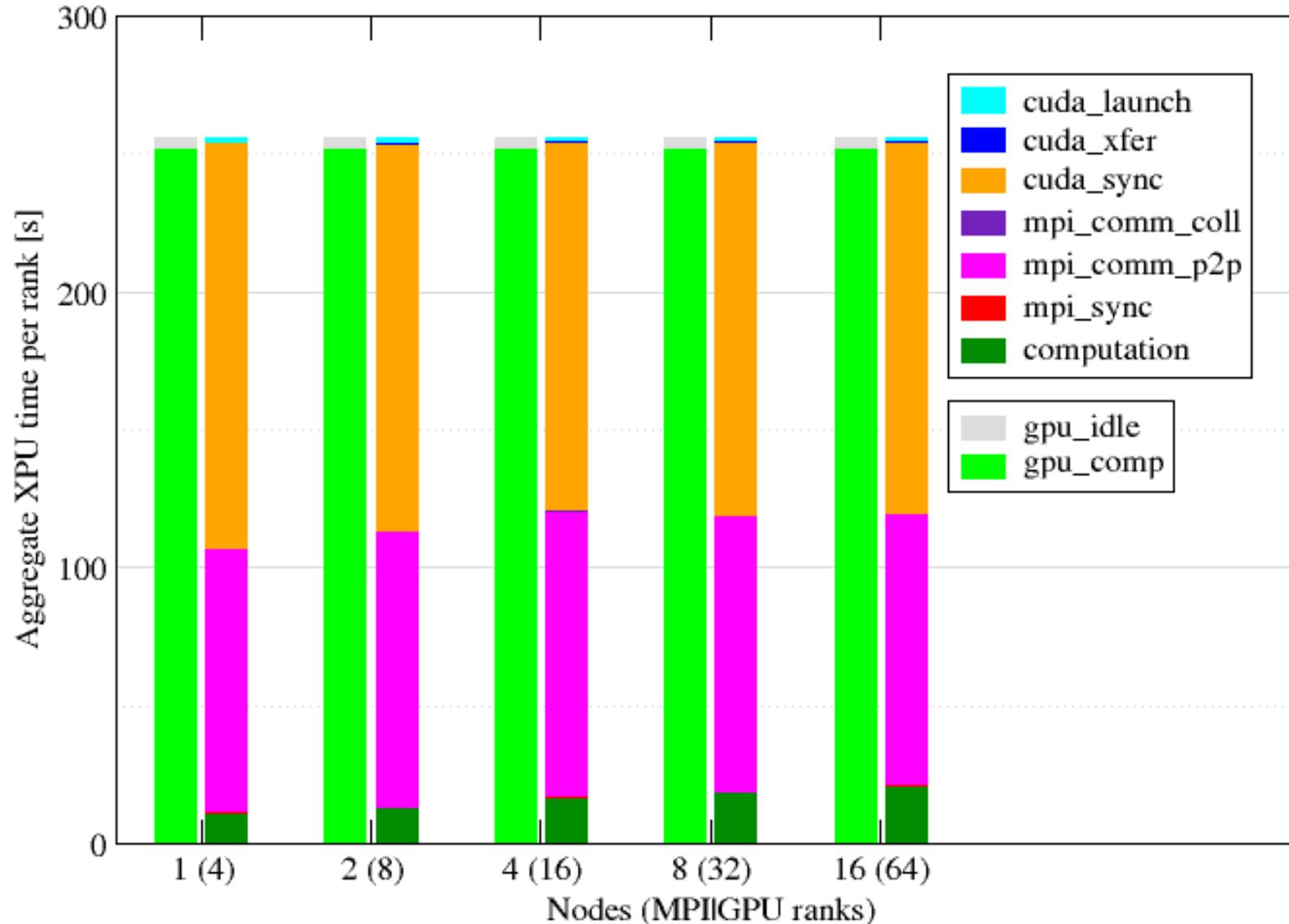
Problem size		1024x1024	1024x1024	1024x1024	1024x1024	1024x1024
MPI GPU ranks		8	64	128	256	512
Wall time [s]		2001.806	255.948	135.478	73.846	61.400
<b>XPU</b>	Global scaling efficiency	0.515	0.504	0.476	0.437	0.263
	- Computation time scaling	1.000	0.943	0.861	0.835	0.715
	- Parallel efficiency	0.515	0.534	0.553	0.523	0.367
	- - Load balance efficiency	0.518	0.541	0.583	0.586	0.671
	- - Orchestration efficiency	0.995	0.988	0.948	0.892	0.547
<b>GPU</b>	Global scaling efficiency	0.995	0.973	0.919	0.843	0.507
	- Computation time scaling	1.000	0.987	0.972	0.950	0.937
	- Parallel efficiency	0.995	0.986	0.945	0.887	0.541
	- - Load balance efficiency	1.000	0.998	0.996	0.994	0.989
	- - Orchestration efficiency	0.995	0.988	0.948	0.892	0.547

- Orchestration efficiency
  - MPI communication & CUDA management
  - aka Communication efficiency



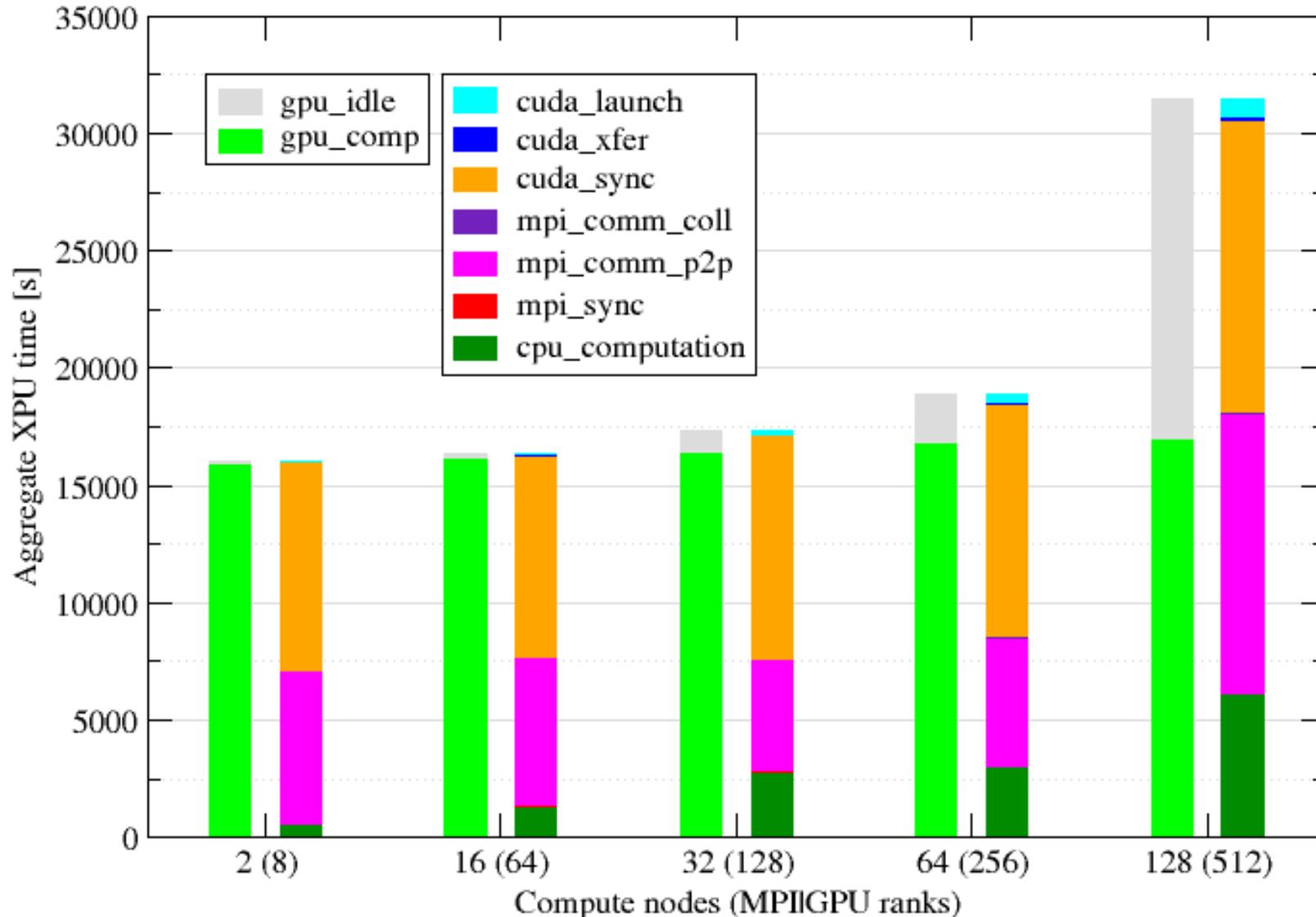
- Good GPU weak scaling efficiency to 128 GPUs (excellent load balance)
- Very poor CPU efficiency?
- Moderate XPU (GPU+CPU) efficiency?

# Weak scaling (128x128 per rank)



- Excellent weak scaling
- Little GPU idle time
- MPI communication effectively overlapped with GPU kernel computation

# Strong scaling (1024x1024 total)



- Good scaling to 256 GPUs (64 nodes)
- GPU computation time slowly grows progressively
- GPU idle time grows for 256 & particularly 512 GPUs
- CPU computation time grows substantially
  - sync\_copy\_from\_device & transfer\_boundary\_to\_device\_a
- For 512 GPUs, growing MPI communication no longer fully overlapped with GPU kernel computation

# Topological inhomogeneities



Kernel variant executed (characteristics and corresponding execution time) varies according to position in 2D grid: four corners, upper/lower & left/right edges, interior

The screenshot shows the CubeGUI-4.8.2 interface. The 'Metric tree' on the left shows a total time of 0.000 seconds, with a highlighted '32961.261 Device' metric. The 'Call tree' in the center shows a list of kernels, with 'compute\_stacey\_elastic\_single\_kernel' variants highlighted in blue. The 'Peer percent' window on the right displays a 16x32 grid of colored squares representing different kernel variants. The status bar at the bottom shows 'Ready' and various performance metrics.

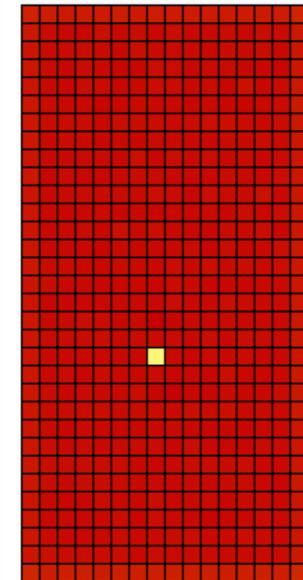
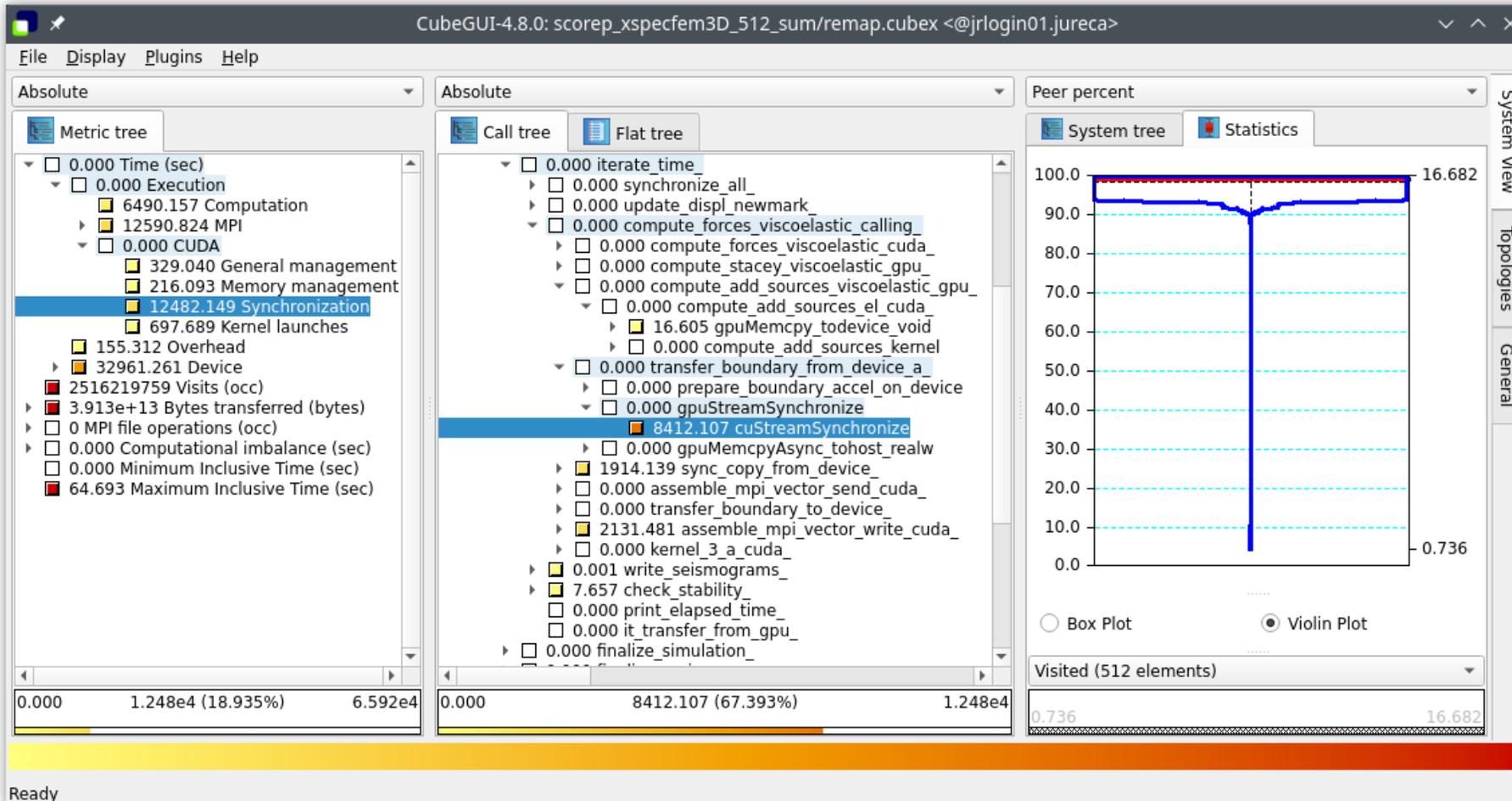
16x32 grid



# GPU computation imbalance



36% of CPU execution time is CUDA synchronization, 67% of which is 16s within *transfer\_boundary\_from\_device\_a* following *compute\_add\_sources\_kernel* that's only executed by a single GPU (source rank 243 of 512)



# Summary of observations



- *iterate\_time* (solver) chosen as focus of analysis
  - negligible time for initialization/finalization
- Excellent weak scaling up to 16 nodes (64 GPUs) and likely beyond
  - Computation very well balanced over GPUs; Excellent GPU utilization
  - MPI P2P communication time grows with scale, but effectively overlapped with GPU computation kernels
- Good strong scaling speedup up to 64 nodes (256 GPUs)
  - Computation remains very well balanced over GPUs
  - Orchestration efficiency progressively diminishes
    - *compute\_add\_sources\_kernel* execution by a single GPU seems the main origin
  - MPI P2P communication time grows significantly, becomes no longer fully overlapped with GPU computation kernels



# Performance Optimisation and Productivity

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://www.youtube.com/POPHPC)

