# POP2_AR_065: HemeLB_GPU on JUWELS performance assessment report

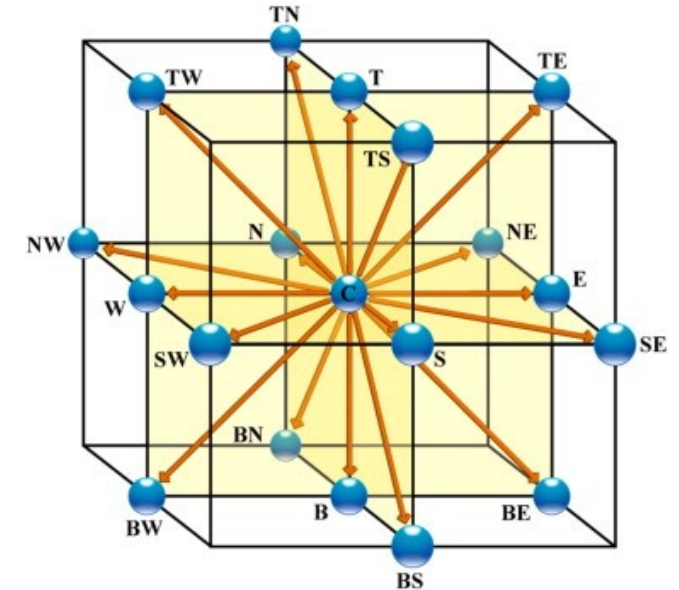## Brian Wylie, Jülich Supercomputing Centre (JSC)

## b.wylie@fz-juelich.de, July 2020

# Background

- Applicant: Ioannis Zacharoudiou, UC London, UK (developer)  [HPC CoE CompBioMed]

- Code: **HemeLB_GPU** (3D macroscopic blood flow simulation) C++ & MPI & CUDA

- Platform: **JUWELS** (@ JSC)

  - 56 dual 20-core Intel Xeon Platinum 6148 'Skylake' compute nodes each accelerated with 4 Nvidia V100 'Volta' GPUs

  - GCC/8.3.0 & CUDA/10.1.105 & ParaStationMPI/5.4

- Testcase: 1.78 GiB CBM2019_Arteries_patched geometry

  - 66,401,494 lattice sites; 1+38 iolets; simulation of 2,000 time-steps (of 100μs)

- Scale: up to 128+1 processes (32 compute nodes, each with 4 MPI ranks)

- Scalasca/Score-P summary and trace measurements

  - using selective instrumentation filter

# HemeLB lattice-Boltzmann method for haemodynamic simulation

- Fluid particles tracked at macroscopic level on a lattice grid

  - discrete set of permissible velocities

  - only nearest-neighbour interactions

- Comparable accuracy to conventional continuum CFD

  - relatively straightforward implementation of complex boundary conditions

  - exhibit superior parallel performance

    - less communication between computational subdomains

www.2020science.net/software/hemelb.html                    www.hemelb.org

# Details

- JUWELS (GPU nodes): 46 (+10 reserved for development)

  - 192 GiB compute nodes; Dual EDR-Infiniband (Connect-X4)

  - IBM Spectrum Scale (GPFS) parallel filesystem; CentOS Linux 7.8.2003

- Execution configuration:

  - gres=gpu:4, ntasks-per-node=4 (plus one extra for monitor on first node)

  - all GPUs used, but only 10% of CPU cores

- HemeLB: GPU development version (using BasicDecomposition scheme)

  - configured to offload kernels from each MPI process to associated GPU

  - MPI File writing of intermediate properties state every 1000 steps

  - 2 variants (reordered 'a' & original 'b') order of actions within each timestep
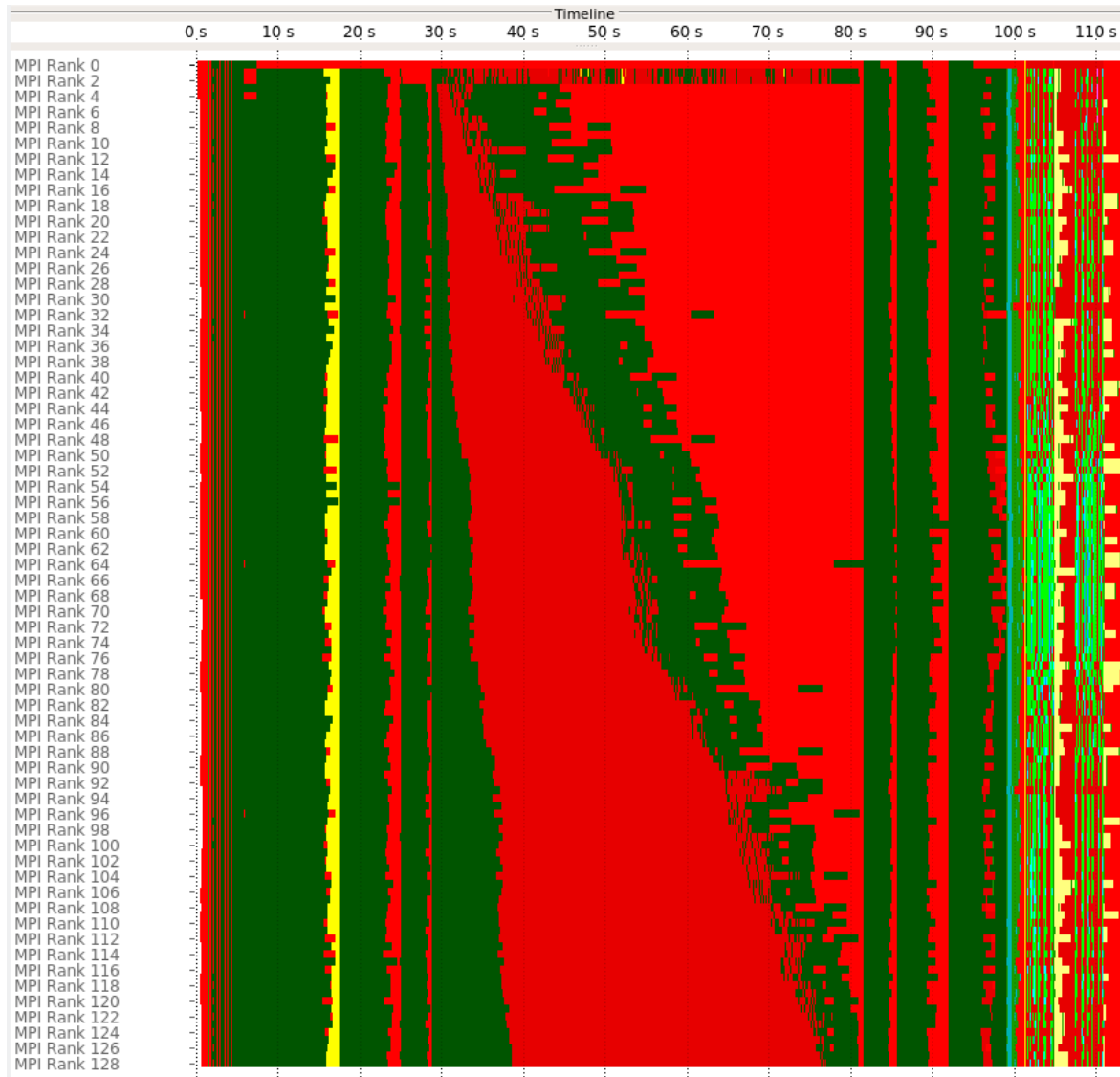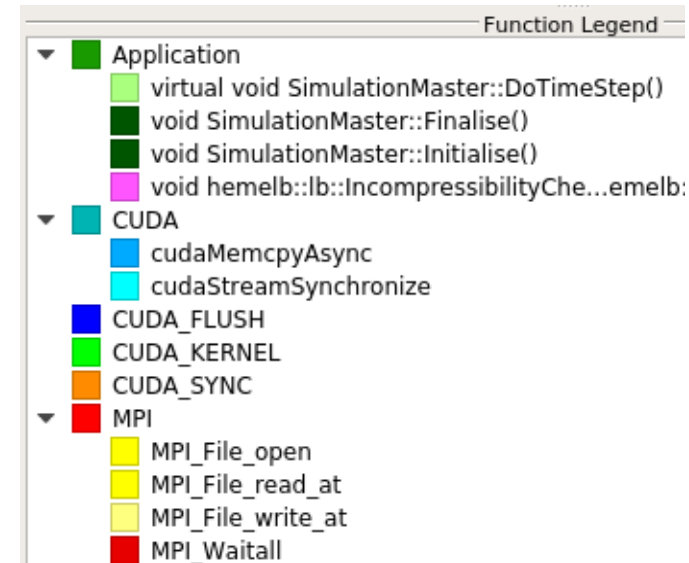
# Method

- Score-P/6.0 (GCC+ParaStationMPI) instrumenter used when building code

  - CMake: CUDACXX=scorep-nvcc (no instrumentation of CXX)

    - workaround replace "-dc" with "--relocatable-device=true -c"

  - set SCOREP_WRAPPER_INSTRUMENTER_FLAGS=
    "--mpp=mpi --thread=none --cuda --instrument-filter=hemelb.filt"

- Scalasca/2.5 runtime measurement configuration

  - SCAN_TRACE_ANALYZER=none
    SCOREP_MPI_ENABLE_GROUPS='coll','env','io','p2p','rma','topo','xnonblock' # 'cg'
    SCOREP_CUDA_ENABLE=runtime,memcpy,kernel,sync,flushatexit
    SCOREP_CUDA_BUFFER=10M
    SCOREP_TOTAL_MEMORY=64M

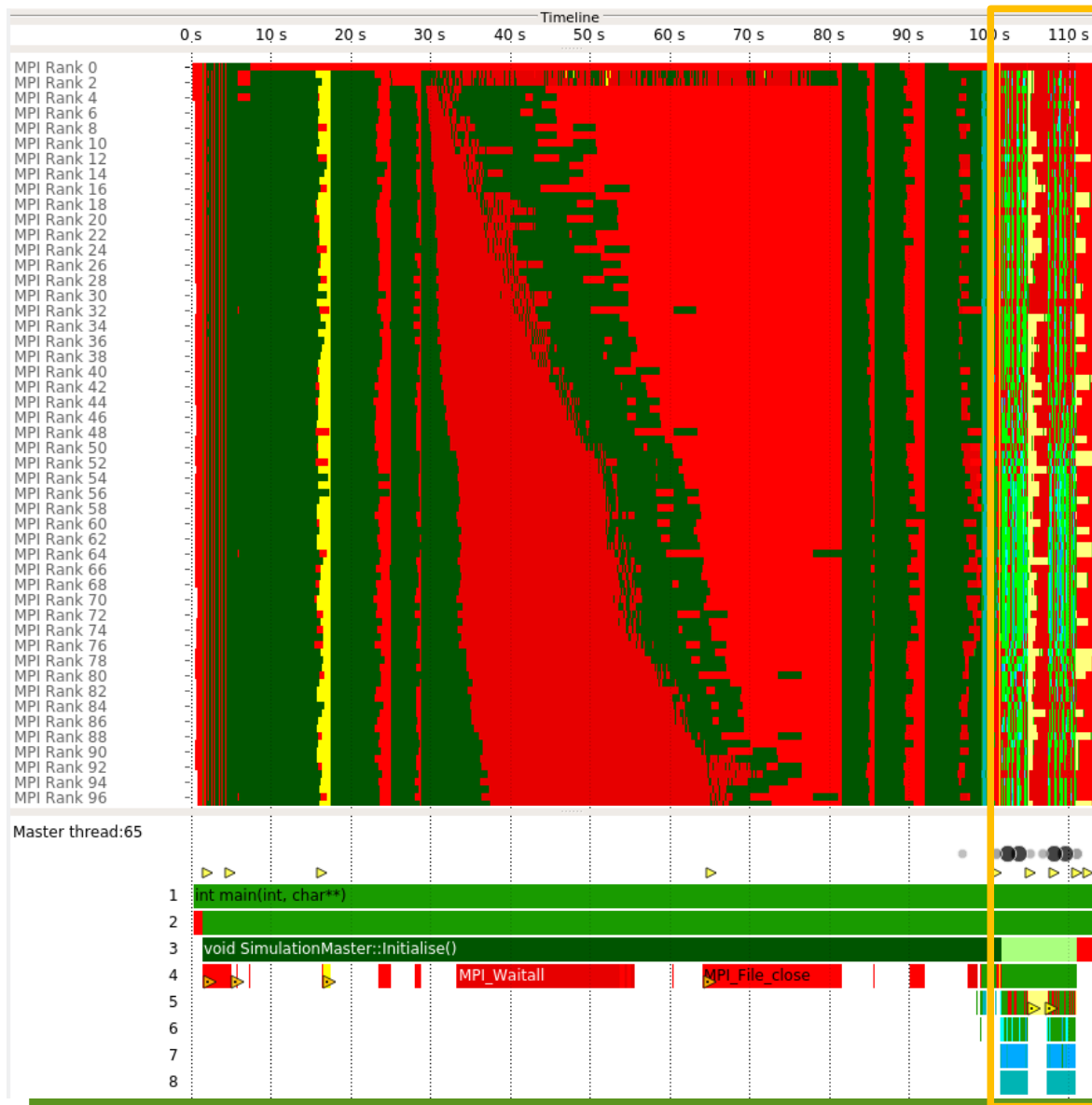- Profiles explored with CUBE/4.4.4, execution traces examined with Vampir/9.8.0

# HemeLB_GPU execution timeline (v1.20a with file writing)



- 32 nodes: 129 MPI processes driving 128 GPUs

  - monitor rank 0 not participating in simulation

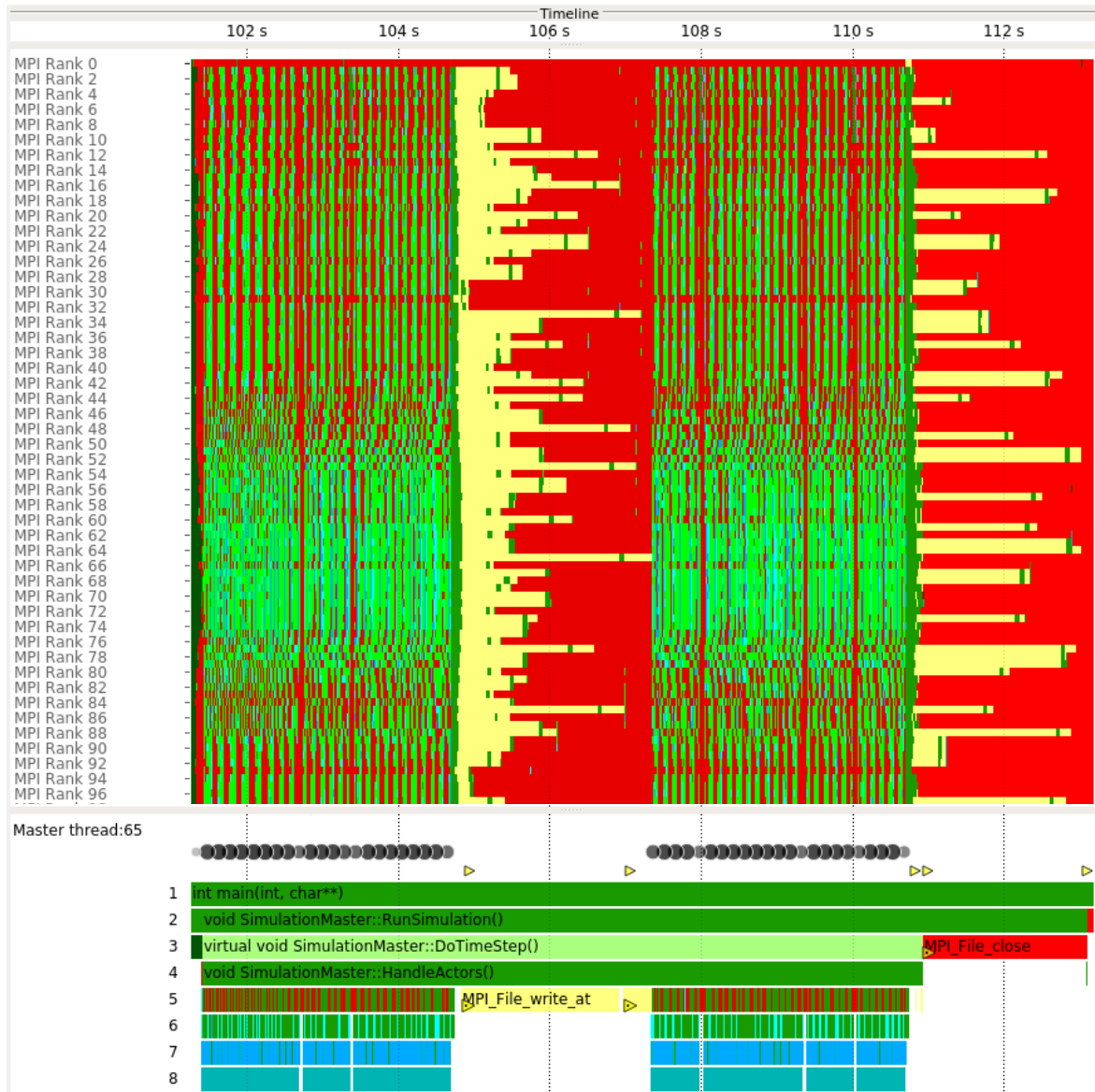  - reader ranks 1 & 2 distribute simulation geometry, then participate in simulation

**Function Legend**

- ▼ ■ Application
  - ■ virtual void SimulationMaster::DoTimeStep()
  - ■ void SimulationMaster::Finalise()
  - ■ void SimulationMaster::Initialise()
  - ■ void hemelb::lb::IncompressibilityChe…emelb:
- ▼ ■ CUDA
  - ■ cudaMemcpyAsync
  - ■ cudaStreamSynchronize
- ■ CUDA_FLUSH
- ■ CUDA_KERNEL
- ■ CUDA_SYNC
- ▼ ■ MPI
  - ■ MPI_File_open
  - ■ MPI_File_read_at
  - ■ MPI_File_write_at
  - ■ MPI_Waitall

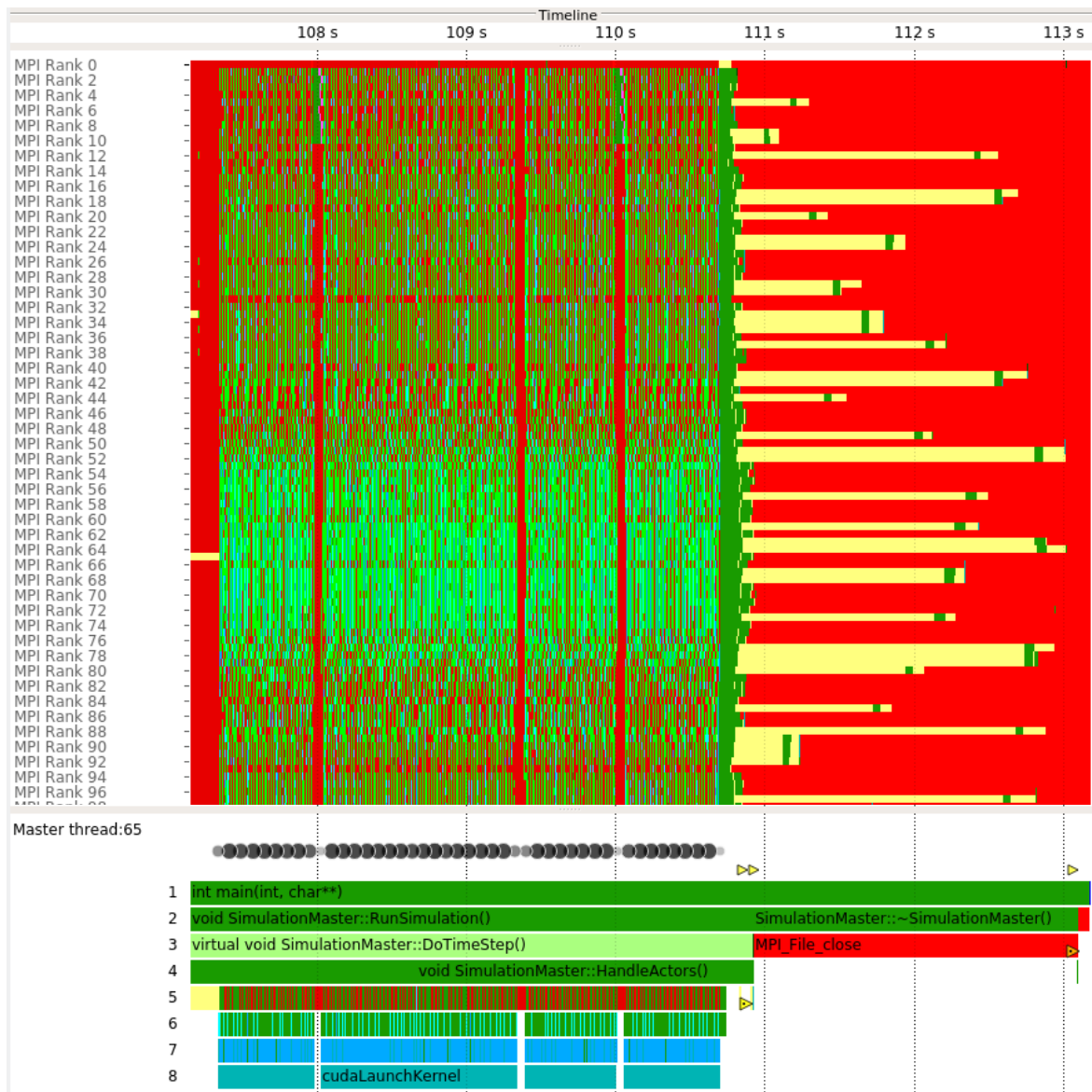# HemeLB_GPU execution timeline: Focus of Analysis (FOA) *DoTimeStep*



- 32 nodes: 129 MPI processes driving 128 GPUs

  – monitor rank 0 not participating in simulation

  – reader ranks 1 &2 distribute simulation geometry, then participate in simulation

- SimulationMaster class constructor/destructor methods initialisation and finalisation

- FOA is RunSimulation *DoTimeStep* routine

  – 2000 simulation steps: approx. 101-113 seconds

  – CUDA kernel offloads to dedicated GPU

- 32 nodes: 129 MPI processes driving 128 GPUs

  - monitor rank 0 not participating in simulation

  - reader ranks 1 &2 distribute simulation geometry, then participate in simulation

- SimulationMaster class constructor/destructor methods initialisation and finalisation

- FOA is RunSimulation *DoTimeStep* routine

  - 2000 simulation steps: approx. 101-113 seconds

  - CUDA kernel offloads to dedicated GPU

  - property file writing after each 1000 steps

    - MPI_File_write_at by each process

# HemeLB_GPU execution timeline: FOA zoom (1000 steps incl. writing)



- 32 nodes: 129 MPI processes driving 128 GPUs

  - monitor rank 0 not participating in simulation

  - reader ranks 1 &2 distribute simulation geometry, then participate in simulation

- SimulationMaster class constructor/destructor methods initialisation and finalisation

- FOA is RunSimulation *DoTimeStep* routine

  - 2000 simulation steps: approx. 101-113 seconds

  - CUDA kernel offloads to dedicated GPU

  - property file writing after each 1000 steps

    - MPI_File_write_at by each process, very imbalanced

    - amount written varies 27-84MB
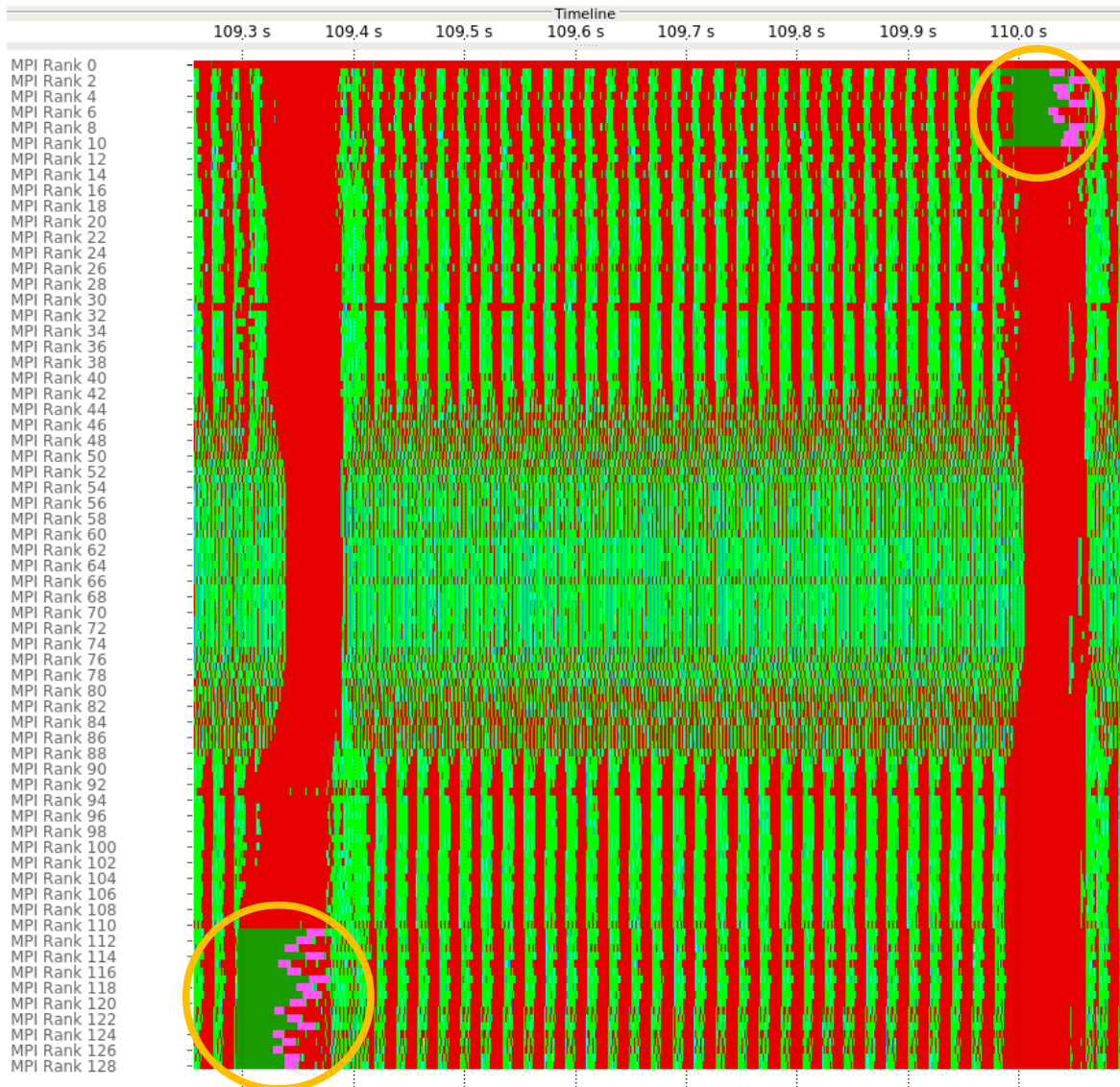
- 32 nodes: 129 MPI processes driving 128 GPUs
  - monitor rank 0 not participating in simulation
  - reader ranks 1 &2 distribute simulation geometry, then participate in simulation
- SimulationMaster class constructor/destructor methods initialisation and finalisation
- FOA is RunSimulation *DoTimeStep* routine
  - 2000 simulation steps: approx. 101-113 seconds
  - CUDA kernel offloads to dedicated GPU
  - property file writing after each 1000 steps
  - *IncompressibilityChecker* each 200 steps (mostly)
    - only lowest 10 simulation ranks
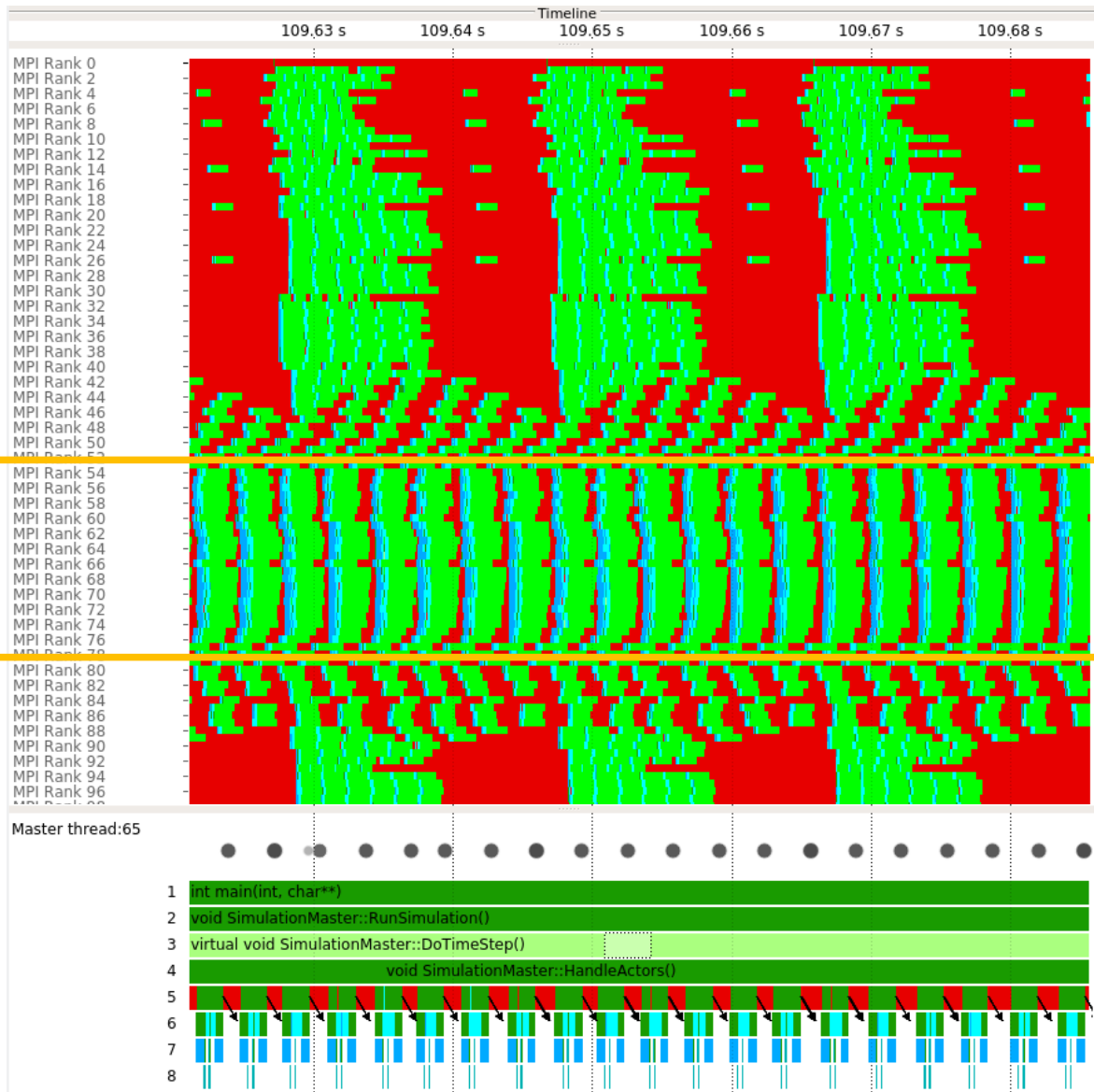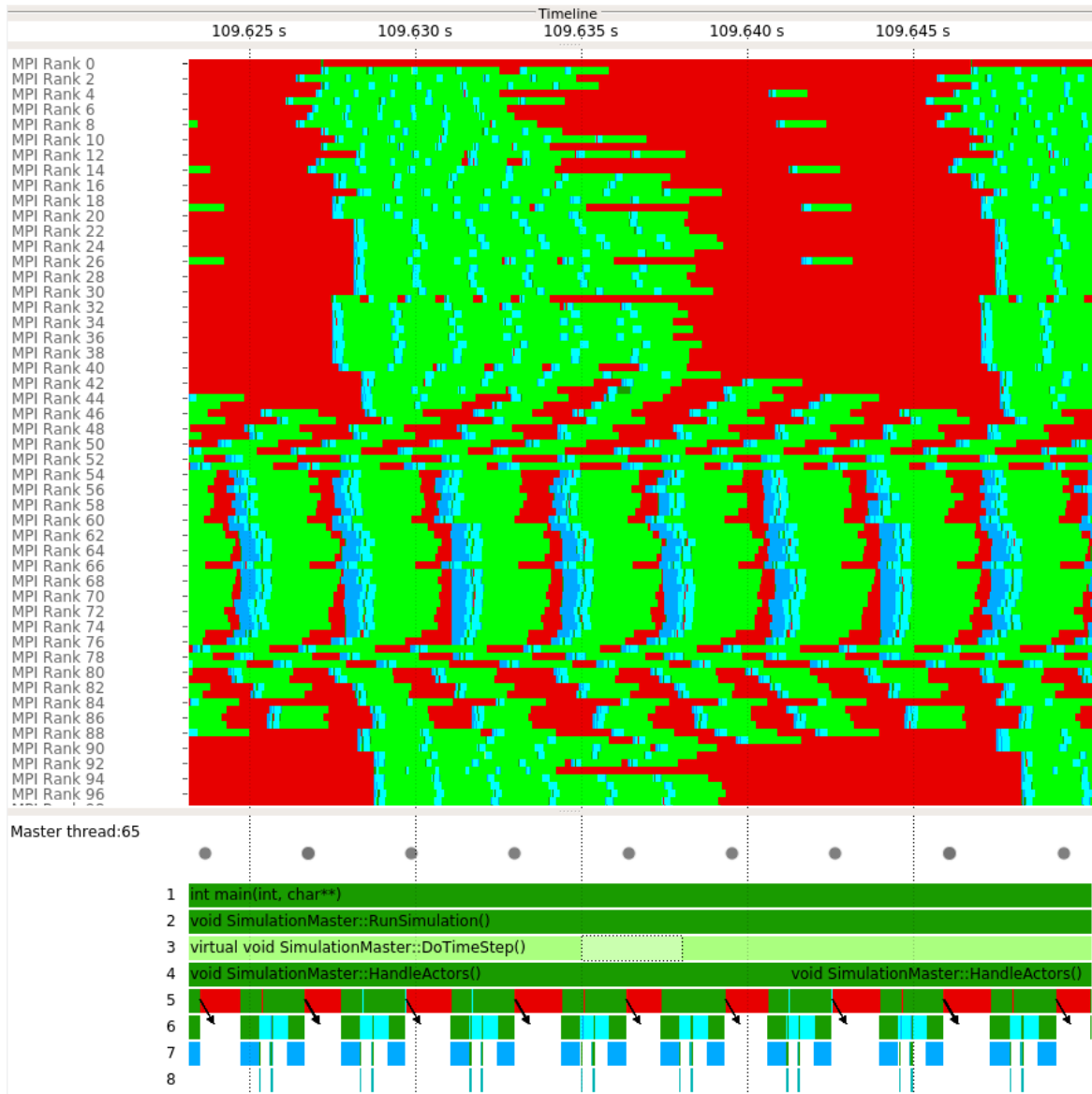    - **or** highest 14 simulation ranks (32 nodes)

# HemeLB_GPU execution timeline: FOA zoom (250 steps)



- 32 nodes: 129 MPI processes driving 128 GPUs

  - monitor rank 0 not participating in simulation

  - reader ranks 1 &2 distribute simulation geometry, then participate in simulation

- SimulationMaster class constructor/destructor methods initialisation and finalisation

- FOA is RunSimulation *DoTimeStep* routine

  - 2000 simulation steps: approx. 101-113 seconds

  - CUDA kernel offloads to dedicated GPU

  - property file writing after each 1000 steps

  - *IncompressibilityChecker* each 200 steps (mostly)

    - only lowest 10 simulation ranks

    - **or** highest 14 simulation ranks (32 nodes)

- 32 nodes: 129 MPI processes driving 128 GPUs

  - monitor rank 0 not participating in simulation

  - reader ranks 1 &2 distribute simulation geometry, then participate in simulation

- SimulationMaster class constructor/destructor methods initialisation and finalisation

- FOA is RunSimulation *DoTimeStep* routine

  - 2000 simulation steps: approx. 101-113 seconds

  - CUDA kernel offloads to dedicated GPU

  - property file writing after each 1000 steps

  - *IncompressibilityChecker* each 200 steps (mostly)

  - interior ranks (54-76) have "uniform" steps

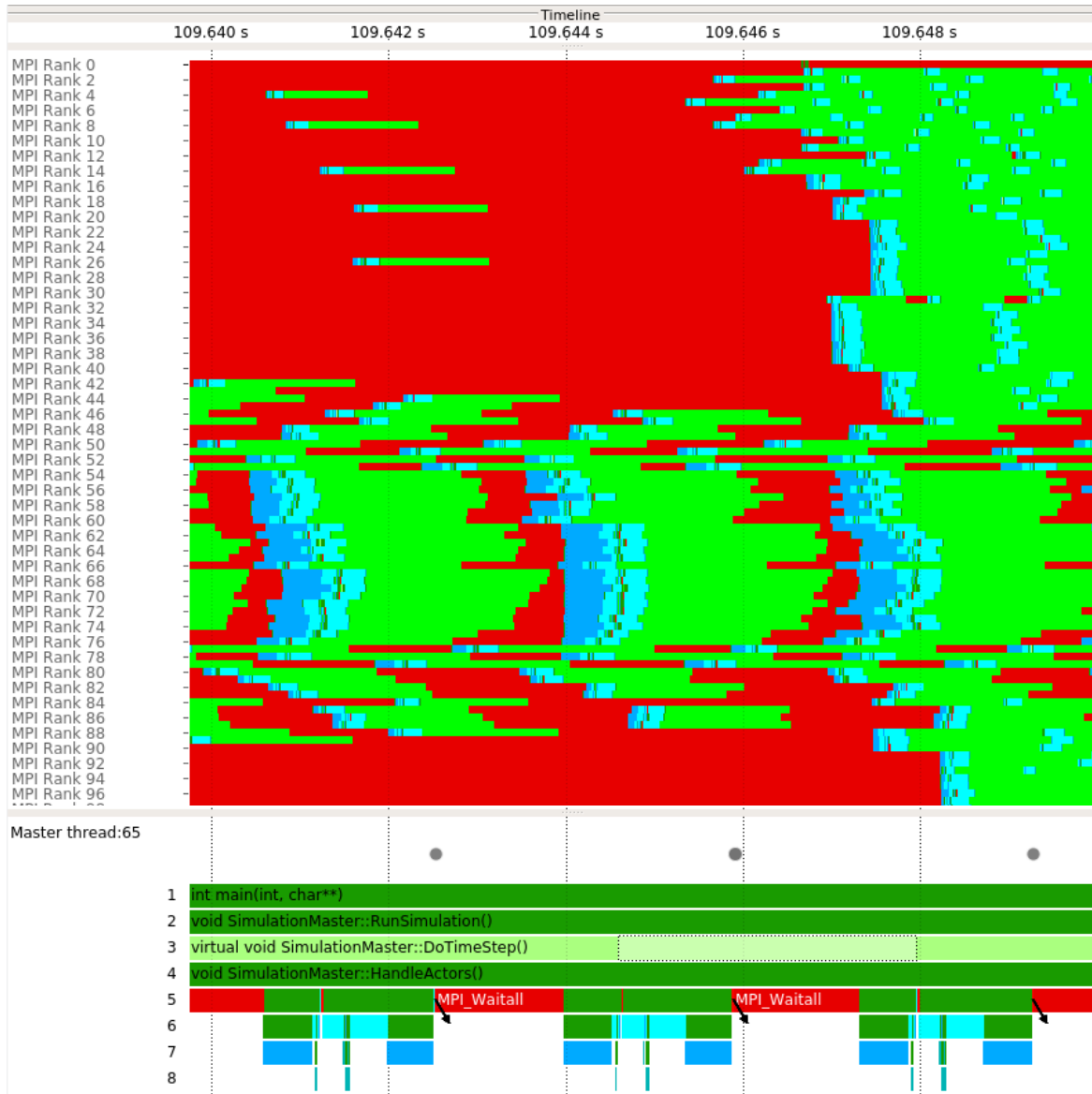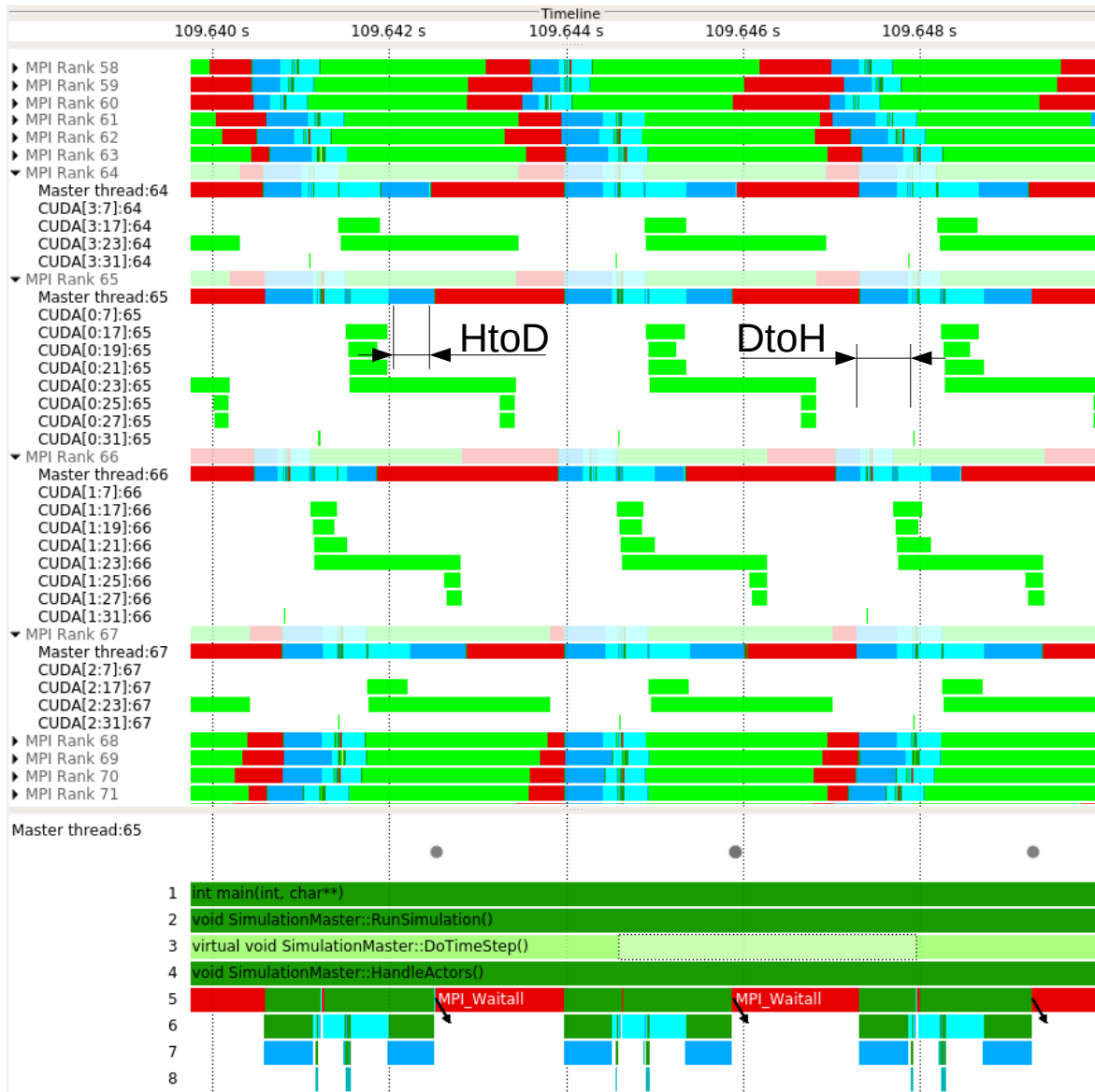    - others blocked every 6 steps, waiting for communication

# HemeLB_GPU execution timeline: FOA zoom (8 steps)



- 32 nodes: 129 MPI processes driving 128 GPUs

    - monitor rank 0 not participating in simulation

    - reader ranks 1 &2 distribute simulation geometry, then participate in simulation

- SimulationMaster class constructor/destructor methods initialisation and finalisation

- FOA is RunSimulation *DoTimeStep* routine

    - 2000 simulation steps: approx. 101-113 seconds

    - CUDA kernel offloads to dedicated GPU

    - property file writing after each 1000 steps

    - *IncompressibilityChecker* each 200 steps (mostly)

        - interior ranks (54-76) have "uniform" steps

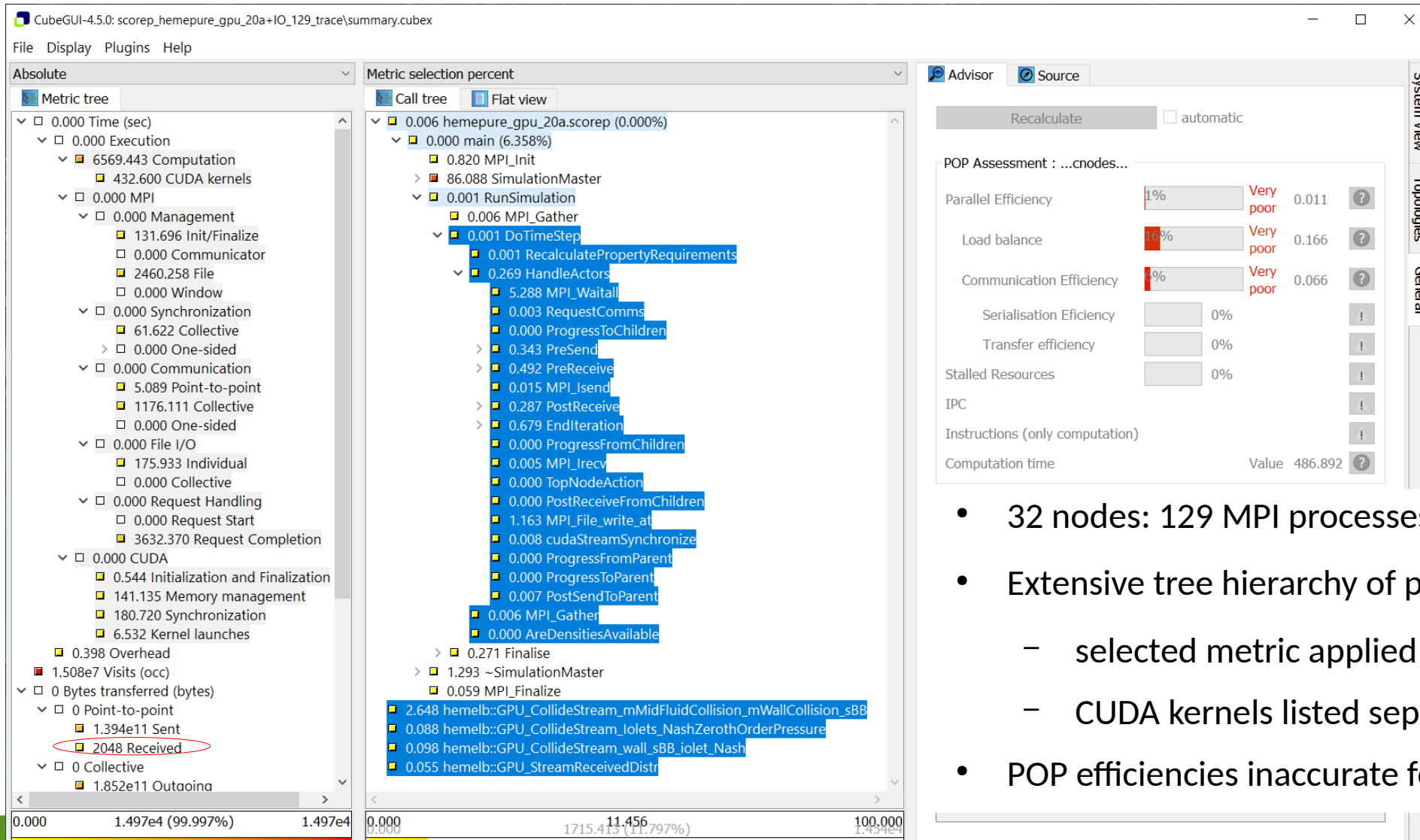        - others blocked every 6 steps, waiting for communication

- 32 nodes: 129 MPI processes driving 128 GPUs

  - monitor rank 0 not participating in simulation

  - reader ranks 1 &2 distribute simulation geometry, then participate in simulation

- SimulationMaster class constructor/destructor methods initialisation and finalisation

- FOA is RunSimulation *DoTimeStep* routine

  - 2000 simulation steps: approx. 101-113 seconds

  - CUDA kernel offloads to dedicated GPU

  - property file writing after each 1000 steps

  - *IncompressibilityChecker* each 200 steps (mostly)

    - interior ranks (54-76) have "uniform" steps

    - others blocked every 6 steps, waiting for communication

- 32 nodes: 129 MPI processes driving 128 GPUs

  - monitor rank 0 not participating in simulation

  - reader ranks 1 &2 distribute simulation geometry, then participate in simulation

- SimulationMaster class constructor/destructor methods initialisation and finalisation

- FOA is RunSimulation *DoTimeStep* routine

  - 2000 simulation steps: approx. 101-113 seconds

  - CUDA kernel offloads to dedicated GPU

    - some ranks/GPUs have additional kernels (to process iolets) executed concurrently

    - critical PostReceive *cuMemcpyDtoHAsync* in Read_DistrFunctions_CPU_to_GPU_totalSharedFs

      - no longer overlapped with kernels

# HemeLB_GPU profile



- 32 nodes: 129 MPI processes driving 128 GPUs

- Extensive tree hierarchy of performance metrics

    - selected metric applied to routines in call-tree

    - CUDA kernels listed separately at bottom

- POP efficiencies inaccurate for CPU+GPU hybrid!

# HemeLB_GPU code structure and Focus of Analysis (FOA) DoTimeStep



- 32 nodes: 129 MPI processes driving 128 GPUs

  - v1.20a with file writing, 2000 time-steps

- SimulationMaster class constructor/destructor methods initialisation and finalisation

- FOA is RunSimulation *DoTimeStep* routine

  - loops through key actions

    - PreSend, PreReceive, Send, PostReceive, EndIteration
      - (a) Send after PreReceive
      - (b) Send before PreReceive
    - plus periodic MPI file writing by all processes
    - and occasional *IncompressibilityChecker* by subset

  - launches CUDA kernels on GPUs

    - execute asynchronously on specific CUDA streams
    - 4 distinct CUDA kernels

## 20b

# HemeLB_GPU Simulation CUDA kernels/streams

- PreSend
  - cudaStreamSynchronize [31]
  - [17] *GPU_CollideStream_mMidFluidCollision_mWallCollision_sBB*
  - 2x cudaMemcpyAsync / cuMemcpy**HtoD**Async_v2 + cudaStreamSynchronize [29,30]
  - [19,20] *GPU_CollideStream_Iolets_NashZerothOrderPressure*
  - [21,22] *GPU_CollideStream_wall_sBB_iolet_Nash*

  local partition edge

- PreReceive
  - [23] *GPU_CollideStream_mMidFluidCollision_mWallCollision_sBB*
  - cudaStreamSynchronize [17,18,19,20,21,22]
  - Read_DistrFunctions_GPU_to_CPU_totalSharedFs / cuMemcpy**DtoH**Async_v2 [35]
  - cudaStreamSynchronize [35]
  - [25,26] *GPU_CollideStream_Iolets_NashZerothOrderPressure*
  - [27,28] *GPU_CollideStream_wall_sBB_iolet_Nash*

  local partition interior

- PostReceive
  - Read_DistrFunctions_CPU_to_GPU_totalSharedFs / cuMemcpy**HtoD**Async_v2 [33]
  - cudaStreamSynchronize [33]
  - [31] *GPU_StreamReceivedDistr*

- EndIteration
  - cudaStreamSynchronize [23,24,25,26,27,28]
  - cudaMemcpyAsync / cuMemcpy**DtoD**Async_v2 [31]
  - Read_Macrovariables_GPU_to_CPU / cuMemcpy**DtoH**Async_v2 + cudaStreamSynchronize [34]

# HemeLB_GPU CUDA kernels/streams

Stream   7: Memcpy HtoD [sync]
Stream 17: **PreSend** kernel1 hemelb::GPU_CollideStream_mMidFluidCollision_mWallCollision_sBB
Stream 19: **PreSend** kernel2 hemelb::GPU_CollideStream_Iolets_NashZerothOrderPressure (inlet)
Stream 20: **PreSend** kernel2 hemelb::GPU_CollideStream_Iolets_NashZerothOrderPressure (outlet)
Stream 21: **PreSend** kernel3 hemelb::GPU_CollideStream_wall_sBB_iolet_Nash (inlet)
Stream 22: **PreSend** kernel3 hemelb::GPU_CollideStream_wall_sBB_iolet_Nash (outlet)
Stream 23: **PreReceive** kernel1 hemelb::GPU_CollideStream_mMidFluidCollision_mWallCollision_sBB
Stream 25: **PreReceive** kernel2 hemelb::GPU_CollideStream_Iolets_NashZerothOrderPressure (inlet)
Stream 26: **PreReceive** kernel2 hemelb::GPU_CollideStream_Iolets_NashZerothOrderPressure (outlet)
Stream 27: **PreReceive** kernel3 hemelb::GPU_CollideStream_wall_sBB_iolet_Nash (inlet)
Stream 28: **PreReceive** kernel3 hemelb::GPU_CollideStream_wall_sBB_iolet_Nash (outlet)
Stream 29: ~~Memcpy HtoD [async] (stream_ghost_dens_inlet)~~
Stream 30: ~~Memcpy HtoD [async] (stream_ghost_dens_outlet)~~
Stream 31: ~~Memcpy DtoD [async]~~ + **PostReceive** kernel4 hemelb::GPU_StreamReceivedDistr
Stream 33: ~~Memcpy HtoD [async] (stream_memCpy_CPU_GPU_domainEdge)~~
Stream 34: ~~Memcpy DtoH [async] (stream_Read_Data_GPU_Dens)~~
Stream 35: ~~Memcpy DtoH [async] (mNet_cuda_stream)~~

# HemeLB_GPU Simulation strong scaling (with intermediate file writing)



Legend:
- (8p) 1.18 Simulation
- (4p*) 1.20b Simulation
- (4p*) 1.20b kernels.max
- (4p*) 1.20b kernels.mean
- (4p*) 1.20a Simulation
- (4p*) 1.20a kernels.max
- (4p*) 1.20a kernels.mean

Axes: Time [s] (y-axis), Nodes (4g/node) (x-axis)

- Reference execution v1.18 with 8ppn
  - multiple processes offloading GPU kernels generally unproductive
- Comparison of v1.20a & v1.20b (4ppn)
  - v1.20a generally better
- 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 1 node, 0.50 for 32 nodes)
  - similar for both versions

# HemeLB_GPU Simulation strong scaling (w/o file writing)



- File writing disabled during Simulation
  - previously every 1000 steps
- Scaling greatly improved
  - particularly beyond 8 nodes
- CUDA kernels on GPUs
  - essentially unaffected

**Chart — Relative speed-up vs Nodes (4g/node)**

Legend:
- Perfect
- 80%
- (8p) 1.18 Simulation
- (4p*) 1.20b Simulation
- (4p*) 1.20b Simulation-IO
- (4p*) 1.20a Simulation
- (4p*) 1.20a Simulation-IO

Data labels: 1.67, 1.65, 2.83, 3.06, 5.45, 6.56, 8.68, 12.7, 15.2, 20.7

- File writing disabled during Simulation
  - previously every 1000 steps
- Speed-up improved
  - 21x for v1.20a on 32 nodes (65% scaling efficiency)
  - 80% scaling efficiency retained to 16 nodes
  - most of the loss of scaling occurs moving from shared-memory MPI within a single node to off-node communication via interconnect

# HemeLB_GPU Simulation strong scaling efficiency (v1.20a with file writing)

| | 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 |

Key:
1.1
1.0
0.9
0.8
0.7
0.6
0.5
0.4
0.3
0.2
0.1
0.0

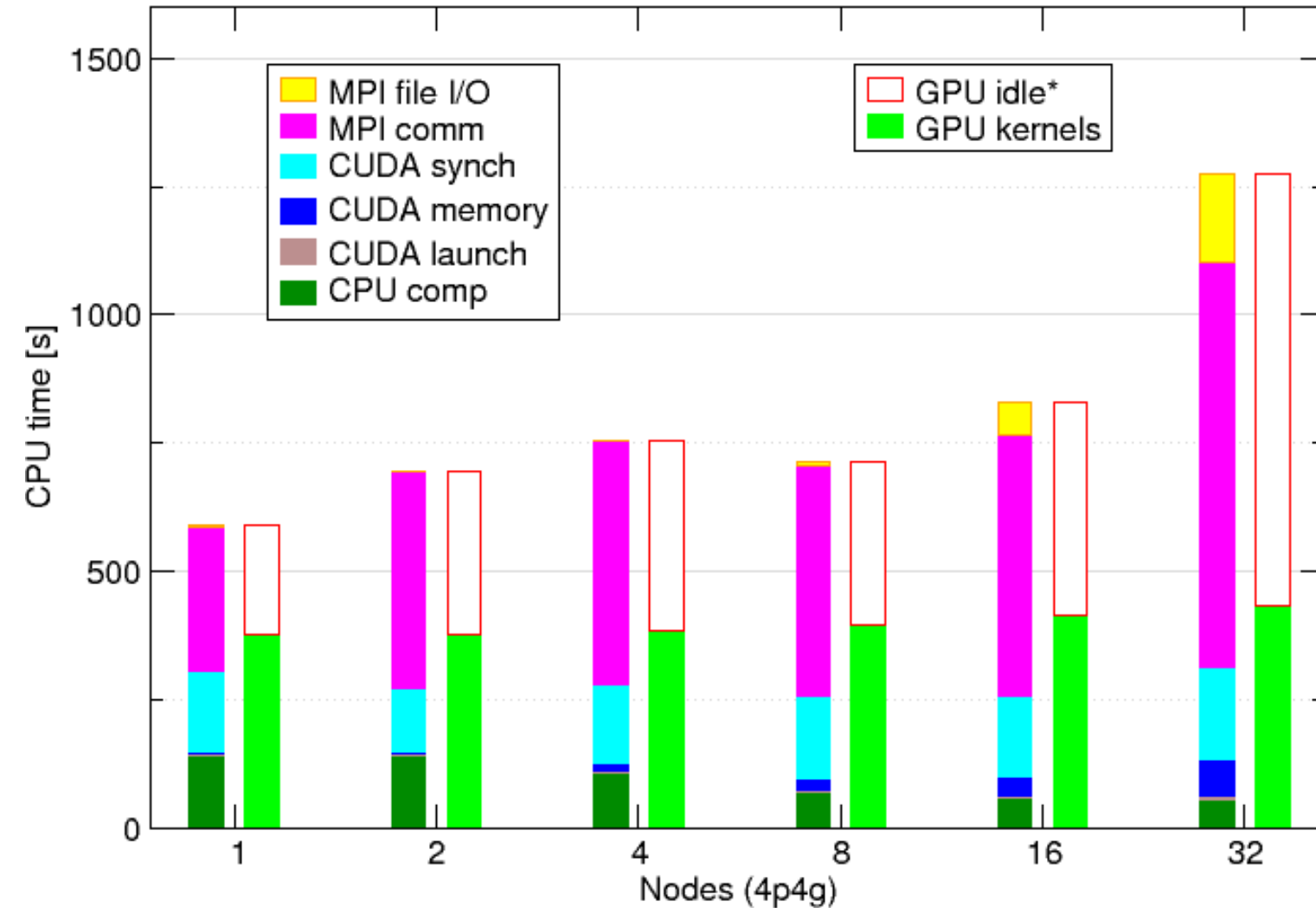- Only considering GPUs (ignoring all CPU cores, 90% of which are completely unused)

- Parallel efficiency determined by load balance and communication (including file I/O)

- Single (quad-GPU) node already suffers significant communication inefficiency

  – but doesn't degrade much as additional nodes are included

- Load balance of GPUs deteriorates progressively

- GPU computation scaling remains reasonably good

# HemeLB_GPU Simulation time breakdown (v1.20a)



- 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 or more nodes

# HemeLB_GPU Simulation time breakdown (v1.20a w/o file writing)



- Disabled Simulation file writing
  - no impact on GPU kernels, nor CUDA operations on CPU
  - reduced CPU computation & MPI communication (and no MPI file I/O)

- CUDA kernels on GPUs
  - reduced GPU idle time, but still significant (roughly half)

- MPI processes on CPUs
  - computation time decreases
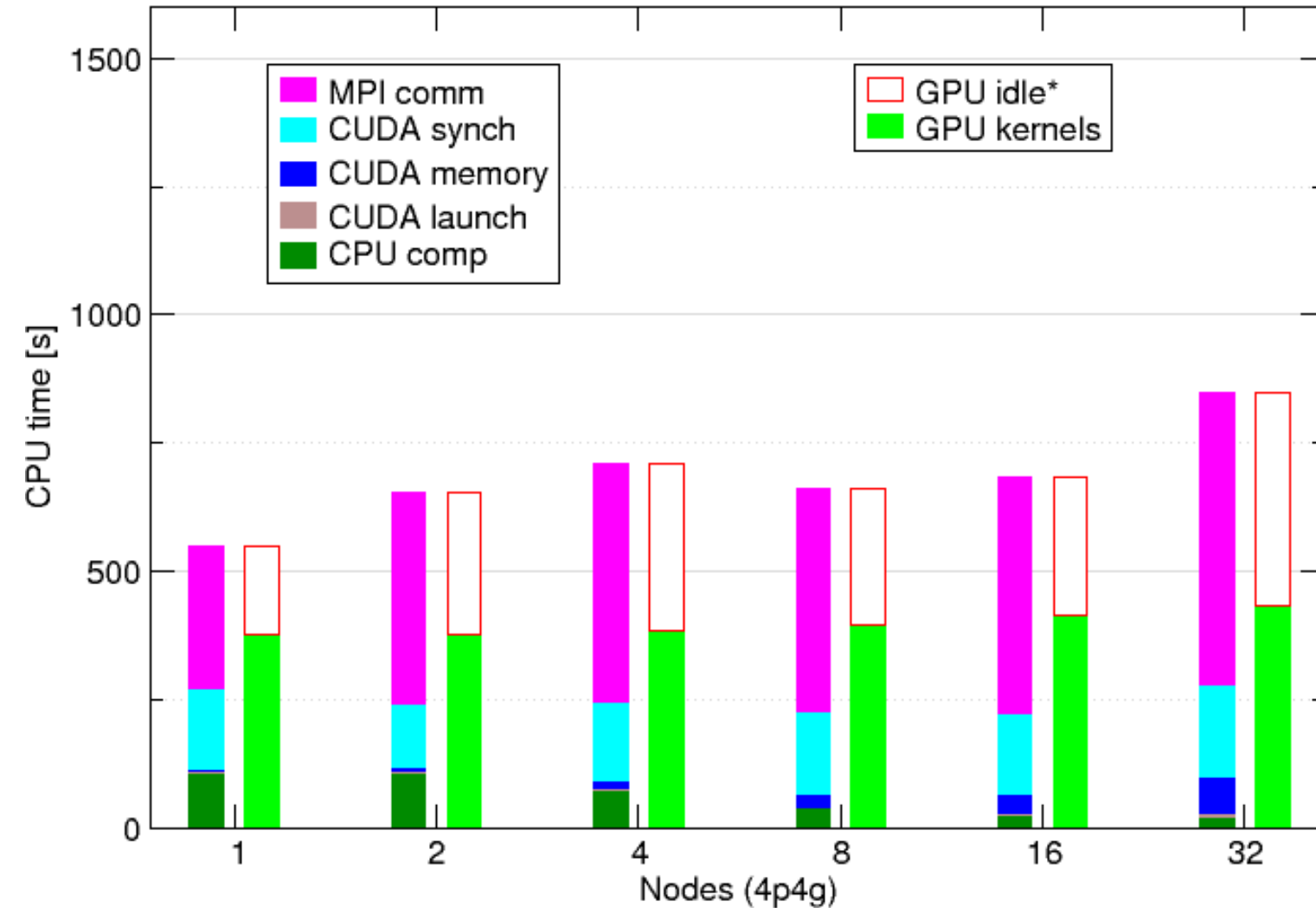  - MPI communication time dominates, at all scales but growing with scale

# HemeLB_GPU Simulation time breakdown (v1.20b)



- CUDA kernels on GPUs

  - less than half of Simulation time (therefore GPUs mostly idle)

  - total kernel time scales very well (0.94 scaling efficiency)

- MPI processes on CPUs

  - computation time roughly constant

  - CUDA synchronization time fairly constant, but time for memory management increases somewhat

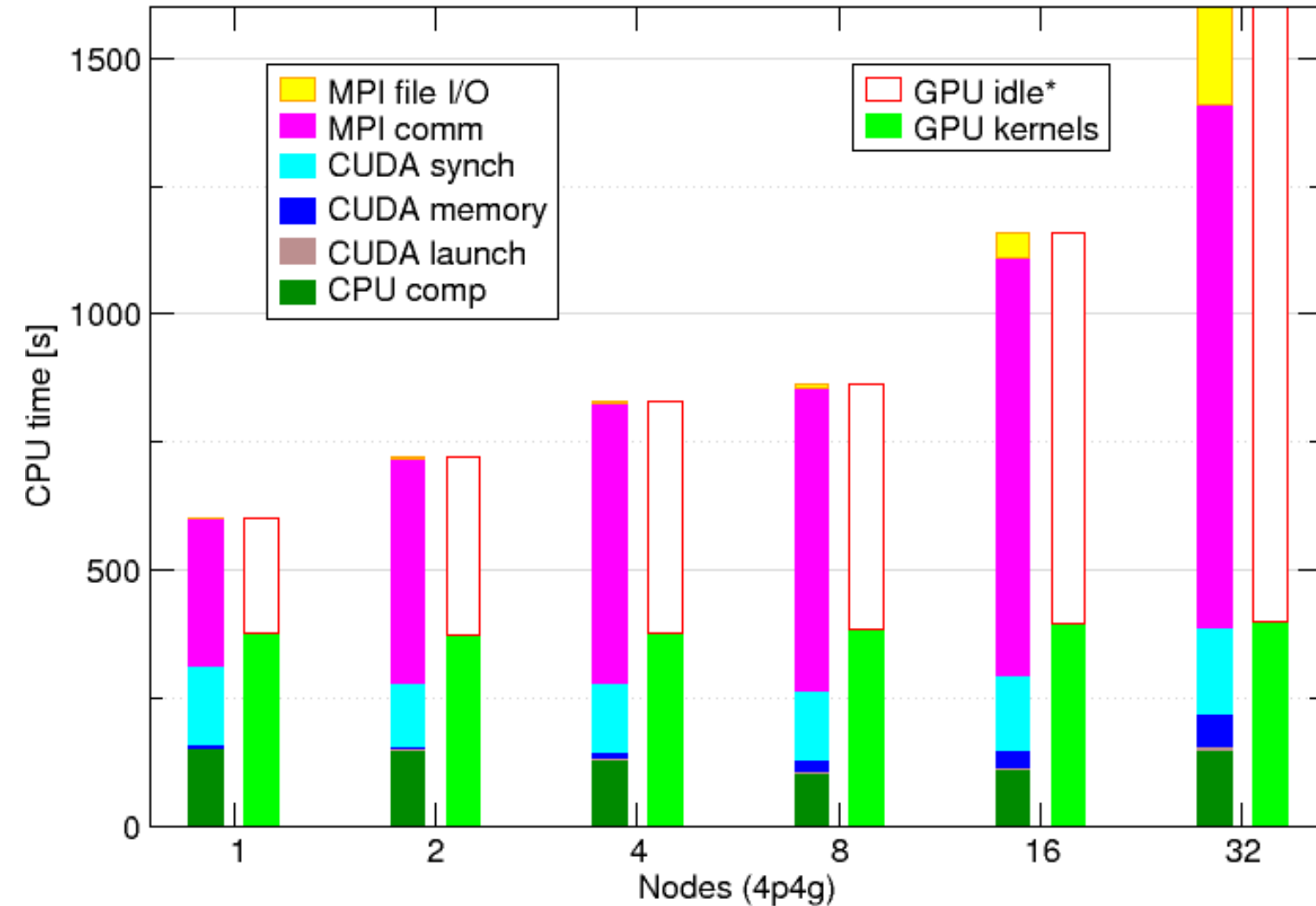  - MPI communication time dominates, with much more time for file writing with 16 or more nodes

# HemeLB_GPU Simulation time breakdown (v1.20b w/o file writing)



- Disabled Simulation file writing

  - no impact on GPU kernels,
    nor CUDA operations on CPU

  - reduced CPU computation & MPI
    communication (and no MPI file I/O)

- CUDA kernels on GPUs

  - reduced GPU idle time, but still
    significant (more than half)

- MPI processes on CPUs

  - computation time roughly constant

  - MPI communication time dominates,
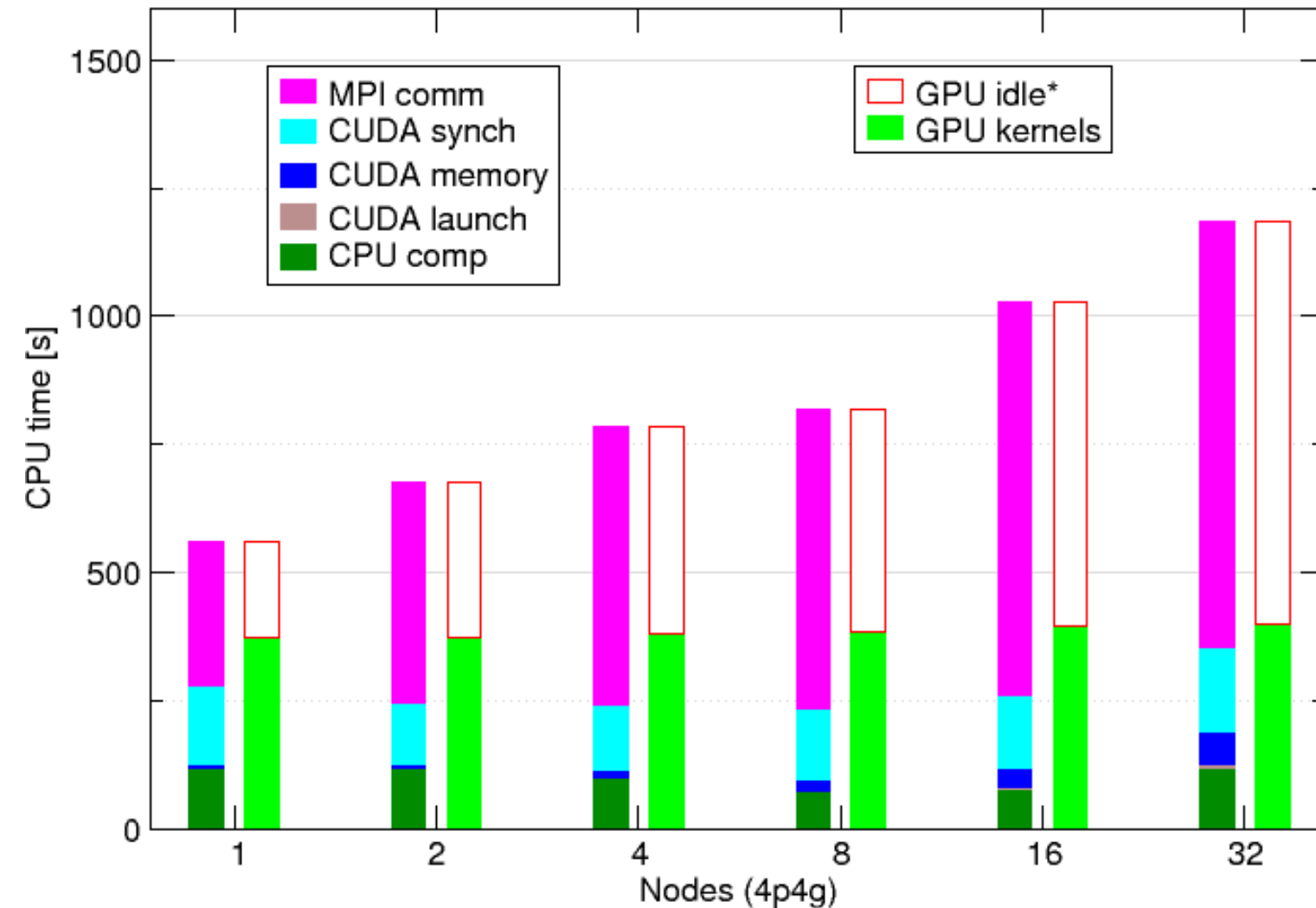    at all scales but growing with scale

# HemeLB_GPU Simulation GPU time balance (v1.20a w/o file writing)



- 32 nodes: 128 MPI processes/GPUs
  - disabled Simulation file writing
- Considerable variations by rank/GPU
  - mid-range ranks 49-80 take longer, partially due to PreSend 17_k1
  - PreReceive 23_k1 dominates, but PreSend 17_k1 also significant
  - rank 65 most heavily overloaded, along with ranks 66 & 60
    - apparently due to processing for iolets at partition edges (kernels k2 & k3)

# HemeLB_GPU Simulation GPU time balance (v1.20a w/o file writing)



32 nodes: 128 MPI processes/GPUs

- disabled Simulation file writing

Considerable variations by rank/GPU

- mid-range ranks 52-76 take longer, partially due to PreSend 17_k1

- PreReceive 23_k1 dominates, but PreSend 17_k1 also significant

- rank 65 most heavily overloaded, along with ranks 66 & 60

  • apparently due to processing for iolets at partition edges (kernels k2 & k3)
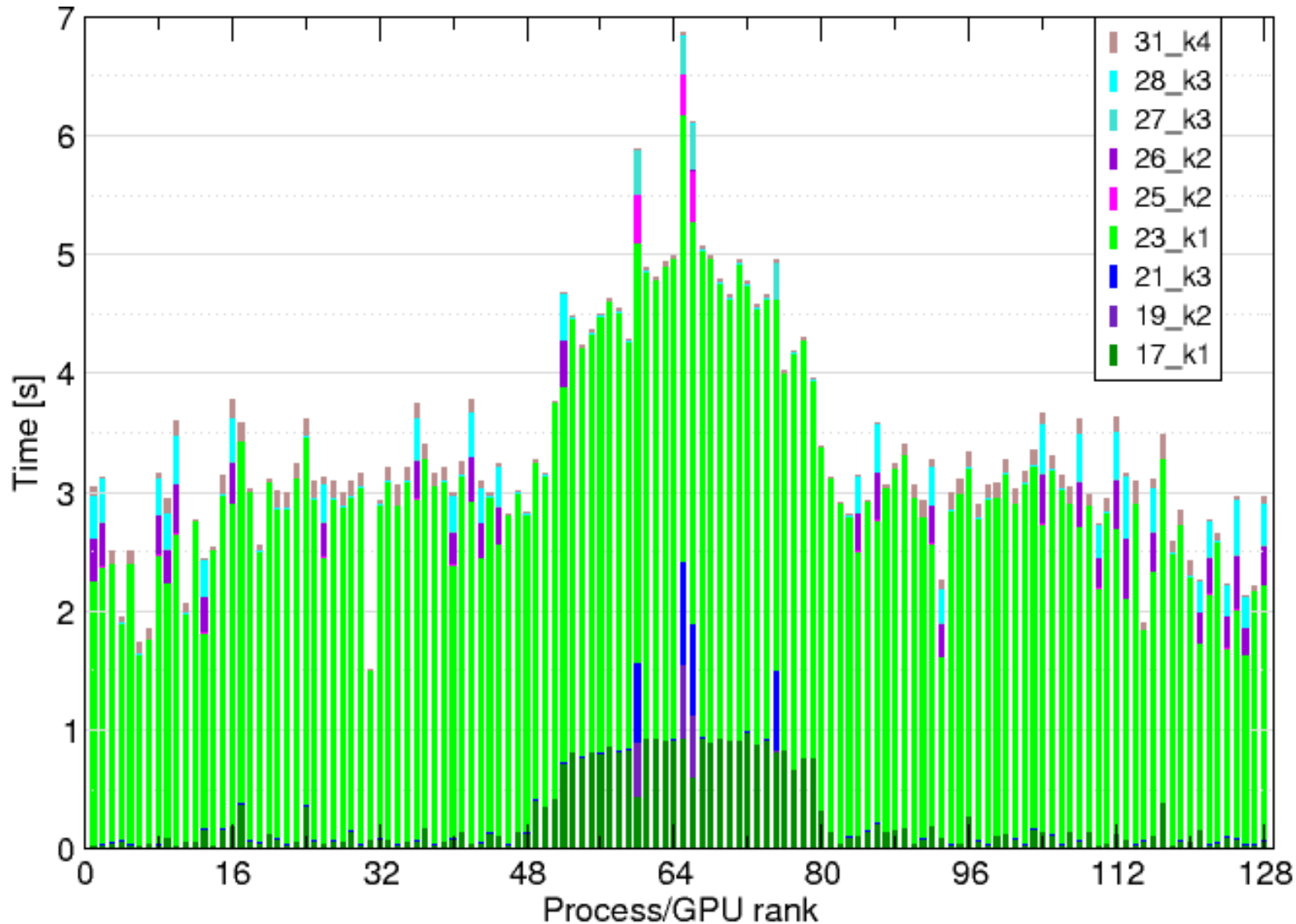
- most GPUs idle half of the time!

# HemeLB_GPU Simulation GPU time balance (v1.20a w/o file writing)



32 nodes: 128 MPI processes/GPUs

- disabled Simulation file writing

Considerable variations by rank

- GPU global memory correlates with number of fluid sites

- CUDA kernel processing time correlates with number of blocks?

- what about iolets?

# HemeLB_GPU Simulation CPU file writing balance (v1.20a)



32 nodes: 128 MPI processes/GPUs

- 4 writes per rank: MPI_File_write_at

- 6.93 GiB written in total (to 2 files)

- 174s total writing time (0.01-3.75)

  - 4 procs: 3.55s (0.89-1.05)

  - 8 procs: 3.50s (0.44-0.71)

  - 16 procs: 3.76s (0.24-0.47)

Considerable variations by rank

- erratic writing time

  - varies from run to run

- time for writing not particularly correlated with amount written

Also results in additional MPI waiting before starting next simulation time step

# HemeLB_GPU Simulation CPU time balance (v1.20a w/o file writing)



32 nodes: 128 MPI processes/GPUs

- disabled Simulation file writing

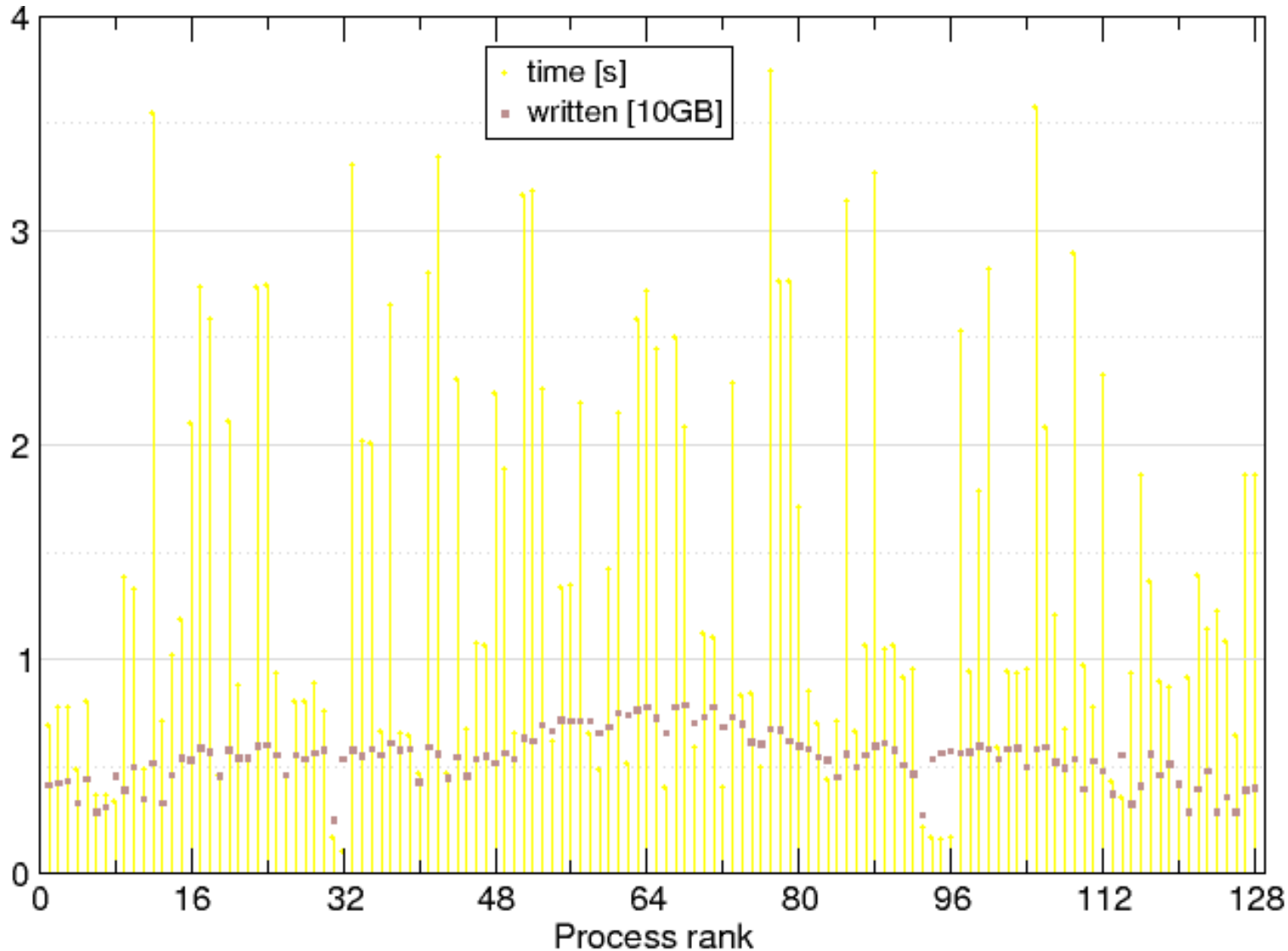Considerable variations by rank

- computation time roughly constant, but noticably higher for highest and lowest ranks

- CUDA memory management much higher for processes in the middle

- CUDA synchronization time varies significantly

- MPI communication time dominates
  - however, almost all waiting time while GPUs compute kernels

# HemeLB_GPU Simulation CPU cudaMemcpyAsync



32 nodes: 128 MPI processes/GPUs

Considerable variations by rank

- CUDA memory management time much higher for processes in the middle

Some copies require much longer than others to initiate

- 51% PreReceive DtoH (Read_DistrFunctions_GPU_to_CPU_total SharedFs)

- 41% PostReceive HtoD (Read_DistrFunctions_CPU_to_GPU_total SharedFs)

- 4% EndIteration DtoD

# HemeLB_GPU Simulation MPI communication imbalance



32 nodes: 128 MPI processes/GPUs

Considerable variations by rank

- ranks 63 & 75 have 12 partners, whereas several have only one

- most mid-range ranks 57-76 have at least 6 partners, and also exchange twice as much data as the others

- however, these mid-range ranks spend less time (waiting) in MPI communication

# HemeLB_GPU Simulation MPI communication matrix


Aggregated Message Volume — MPI communication matrix (Sender vs Receiver)

- 32 nodes: 128 MPI processes/GPUs
  - 4000 MPI_Waitall calls by each rank
  - 2334-24334 MPI_Isend/Irecv calls (messages)
    - max for ranks 63 & 75
  - 0.03 - 5.19 GB sent/received
    - max for rank 65
- Considerable variations by rank
  - heaviest communication for interior processes

# Timeline detail of 3 time-steps

- 2 nodes: 9 MPI processes driving 8 GPUs

  - monitor rank 0 not executing simulation

- One MPI process (#5) mostly doing CUDA synchronization in EndIteration & next PreSend actions

  - others mostly in MPI_Waitall

  - neighbours waiting for return messages from rank 5 (and their neighbours waiting for them)

- Additional CUDA kernels executed by that process/GPU

  - iolet processing both for domain edge and interior (both PreSend and PreReceive actions)

37

# Notes

- Time for CUDA kernels from each GPU aggregated

  - ignores concurrent execution on GPU!

  - more v1.20a concurrency results in more kernel execution time dilation?

- GPUs assumed to idle when no kernels executing

  - don't see asynchronous memory transfers on GPU!

- *IncompressibilityChecker* only executed by (diminishing) subset of processes on CPUs?

  - and not every 200 steps?

- MPI File I/O writing time & bytes by CPUs is imbalanced (particularly at growing scale)

  - associated with additional CPU computation and MPI communication time

  - but otherwise no apparent impact on CUDA kernels (since GPUs idle)

# Comparison

- v1.20a Simulation is faster, particularly with more nodes

    - but v1.20b CUDA kernels execute faster and scale better

    - however, in both cases they execute less than half of the total Simulation time, and suffer from progressively worsening load balance (0.5 efficiency with 32 nodes)

- **_GPU_CollideStream_mMidFluidCollision_mWallCollision_sBB_** kernel constitutes most GPU (non-idle) time: 99% on 1 node decreasing progressively to 92% on 32 nodes (load balance efficiency of 0.7)

    - further load imbalance originates from (less time consuming) iolet kernels

- Simulation is dominated by MPI communication (waiting) time, with MPI file writing becoming most significant for 16 or more nodes

    - v1.20a significantly reduces MPI communication time (particularly without file writing), despite slightly increasing CUDA overheads

# Summary (GPU)

- Good GPU computation scaling, however, ...

- GPU kernel load imbalance is significant

  - use alternative decomposition scheme? (such as Zoltan+ParMETIS or ALL)

    - likely to require more host memory and longer initialisation time

  - revise weighting used by BasicDecomposition?

    - currently default is unweighted, but weighted GMY+ should be used instead

- cudaMemcpyAsync becomes expensive when not overlapped with kernels

  - particularly Read_DistrFunctions_CPU_to_GPU_totalSharedFs cuMemcpyDtoHAsync

  - CUDA-aware MPI should avoid need for these copies between CPU and GPU, and generally be much more efficient for data transfers between GPUs

# Summary (CPU)

- Major inefficiencies arise from file I/O and non-GPU computation

- GPUs are idle much of the time with CPU-only activities during Simulation *DoTimeStep*

    - periodic property file writing with MPI File I/O (by all ranks)

        - possible to do asynchronously (with dedicated processes/threads)

        - or non-blocking MPI_File_iwrite_at (+ MPI_Wait)

    - periodic incompressibility check (only by some ranks): may be better done on GPUs?

- Additional 'monitor' process is extra complication

    - necessary for future coupled codes

- Initialisation of simulation takes a significant time and warrants detailed investigation

    - time grows with size of geometry, but not specifically file reading

        - investigate additional reader ranks (possibly all simulation ranks)?

# pop

## Performance Optimisation and Productivity
### A Centre of Excellence in Computing Applications & HPC

Contact:
https://www.pop-coe.eu
mailto:pop@bsc.es

# HemeLB_GPU Simulation strong scaling efficiency (v1.20a w/o file writing)

| | 1n 5p | 2n 9p | 4n 17p | 8n 33p | 16n 65p | 32n 129p |
|---|---|---|---|---|---|---|
| Simulation time [s] | 137.73 | 82.33 | 44.98 | 21.00 | 10.87 | 6.66 |
| Global scaling efficiency | 0.68 | 0.57 | 0.52 | 0.56 | 0.54 | 0.44 |
| − Parallel efficiency | 0.68 | 0.57 | 0.53 | 0.59 | 0.59 | 0.51 |
| − − Load balance efficiency (GPU) | 0.95 | 0.78 | 0.73 | 0.73 | 0.65 | 0.49 |
| − − Communication efficiency (GPU) | 0.72 | 0.73 | 0.73 | 0.80 | 0.92 | 1.03 |
| − Computation scaling (GPU) | *1.00* | 0.99 | 0.98 | 0.95 | 0.91 | 0.87 |

Key:
1.1
1.0
0.9
0.8
0.7
0.6
0.5
0.4
0.3
0.2
0.1
0.0

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

- Parallel efficiency determined by load balance

- Single (quad-GPU) node already suffers significant communication inefficiency

  − but improves as additional nodes are included?

- Load balance of GPUs deteriorates progressively

- GPU computation scaling remains reasonably good

# HemeLB_GPU Simulation strong scaling efficiency (v1.20a) hybrid CPU+GPU

| | 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 | 0.35 | 0.33 | 0.32 | 0.35 | 0.29 | 0.30 | |
| – Parallel efficiency | 0.35 | 0.33 | 0.30 | 0.31 | 0.27 | 0.28 | |
| – – Load balance efficiency | 0.77 | 0.69 | 0.65 | 0.75 | 0.66 | 0.52 | CPU+GPU |
| – – Communication efficiency | 0.46 | 0.47 | 0.46 | 0.41 | 0.40 | 0.54 | |
| – Computation scaling | 1.00 | 1.00 | 1.06 | 1.12 | 1.11 | 1.07 | |
| Global scaling efficiency (CPU) | 0.19 | 0.18 | 0.17 | 0.19 | 0.16 | 0.09 | |
| – Parallel efficiency (CPU) | 0.19 | 0.18 | 0.13 | 0.10 | 0.07 | 0.04 | |
| – – Load balance efficiency (CPU) | 0.78 | 0.68 | 0.52 | 0.45 | 0.38 | 0.64 | CPU only |
| – – Communication efficiency (CPU) | 0.25 | 0.26 | 0.26 | 0.21 | 0.17 | 0.06 | |
| – Computation scaling (CPU) | 1.00 | 1.00 | 1.32 | 2.01 | 2.47 | 2.63 | |
| File writing scaling | 1.00 | 1.02 | 0.94 | 0.35 | 0.05 | 0.02 | |
| – Bytes written scaling | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | |
| Global scaling efficiency (GPU) | 0.64 | 0.53 | 0.49 | 0.52 | 0.43 | 0.25 | |
| – Parallel efficiency (GPU) | 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 | GPU only |
| – – 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 | |

Key:
1.1
1.0
0.9
0.8
0.7
0.6
0.5
0.4
0.3
0.2
0.1
0.0