

# OPENQCD ON GPU

**Antonio Rago**

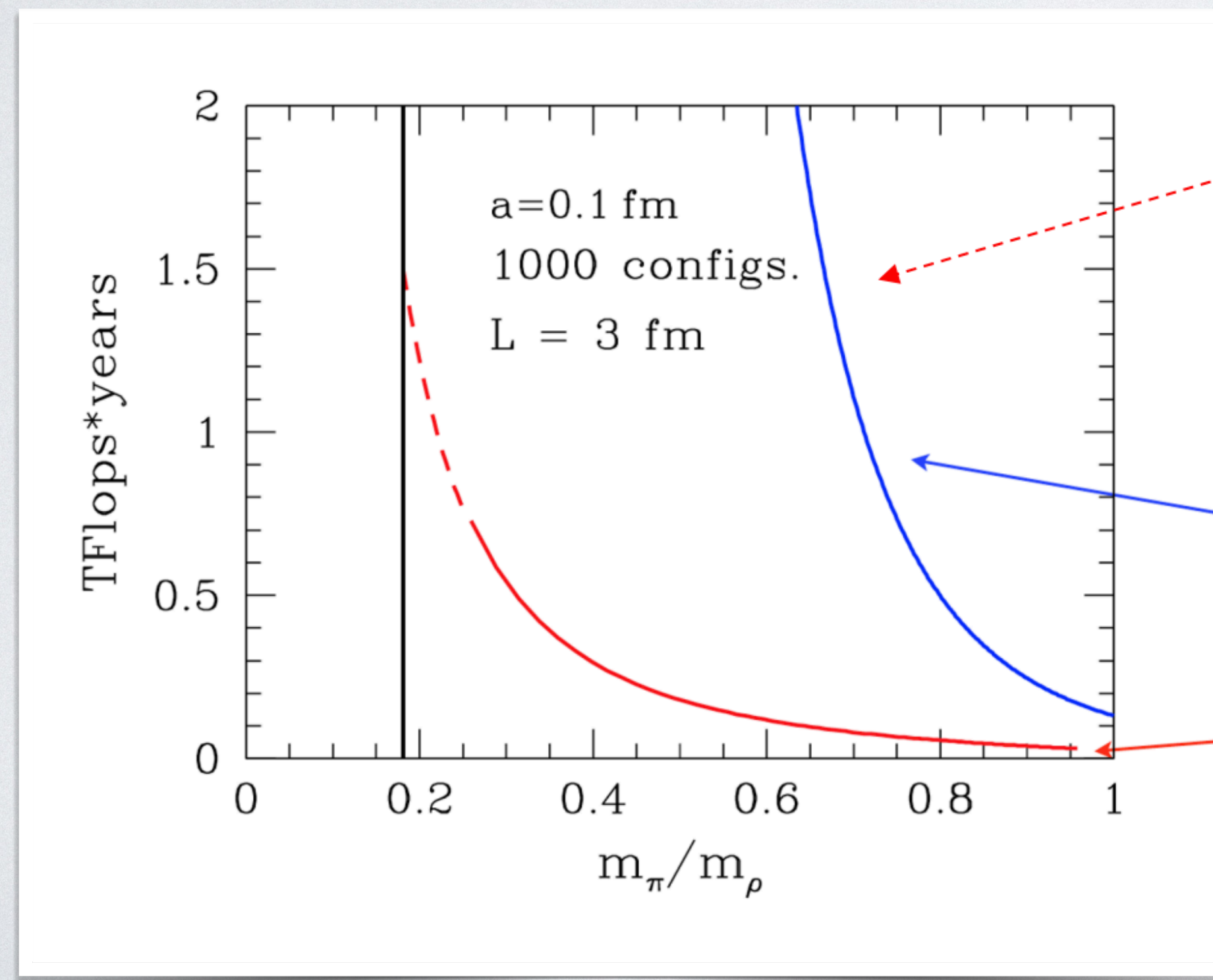
Quantum Theory Center, University of Southern Denmark



# Why openQCD?

openQCD is a simulation suite for lattice QCD, featuring an efficient implementation of the  $\mathcal{O}(a)$  improved Wilson-Dirac fermion operator.

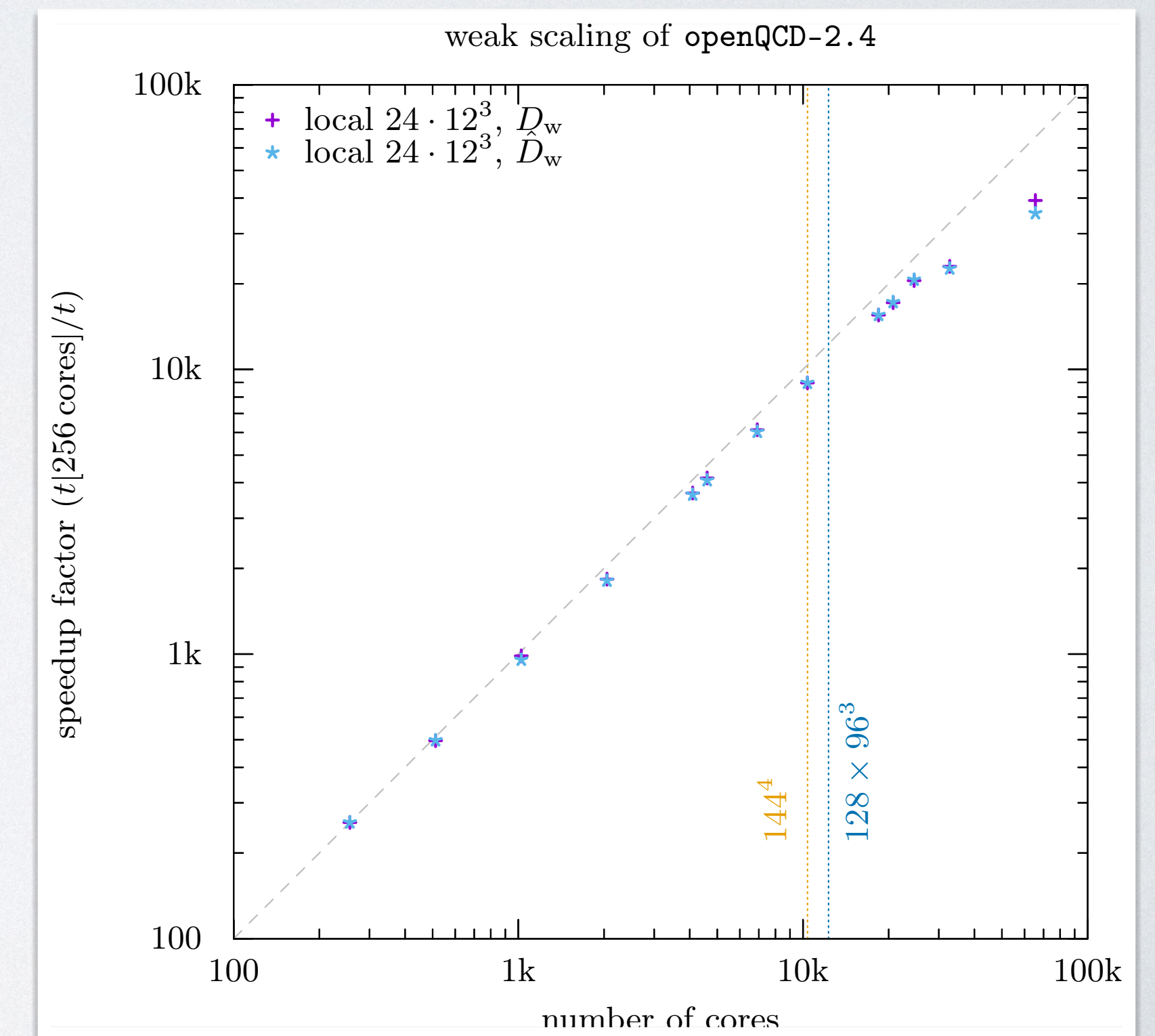
Pivotal to the performance properties of the code is the locally deflated solver.



**The Berlin Wall (2001)**

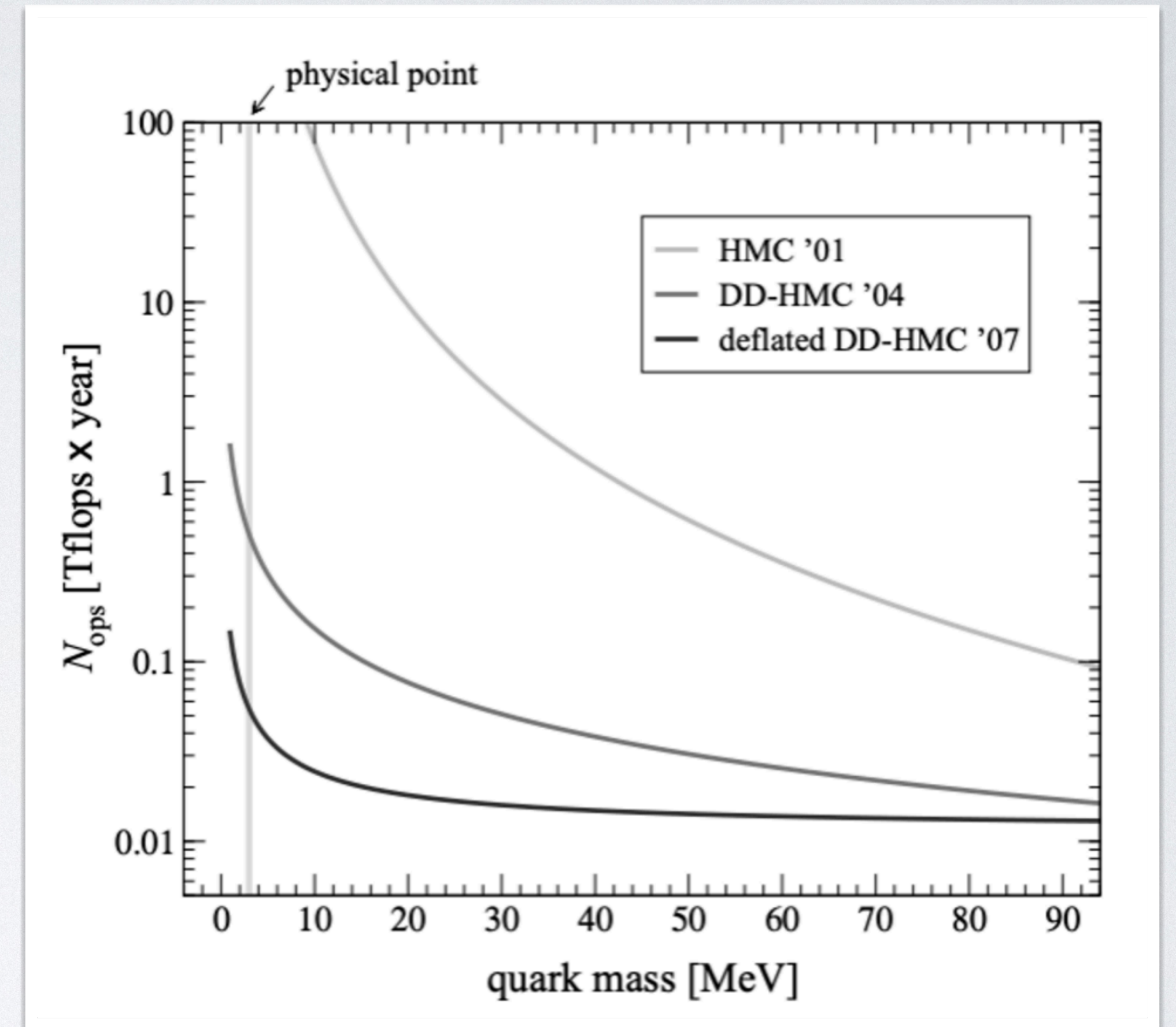
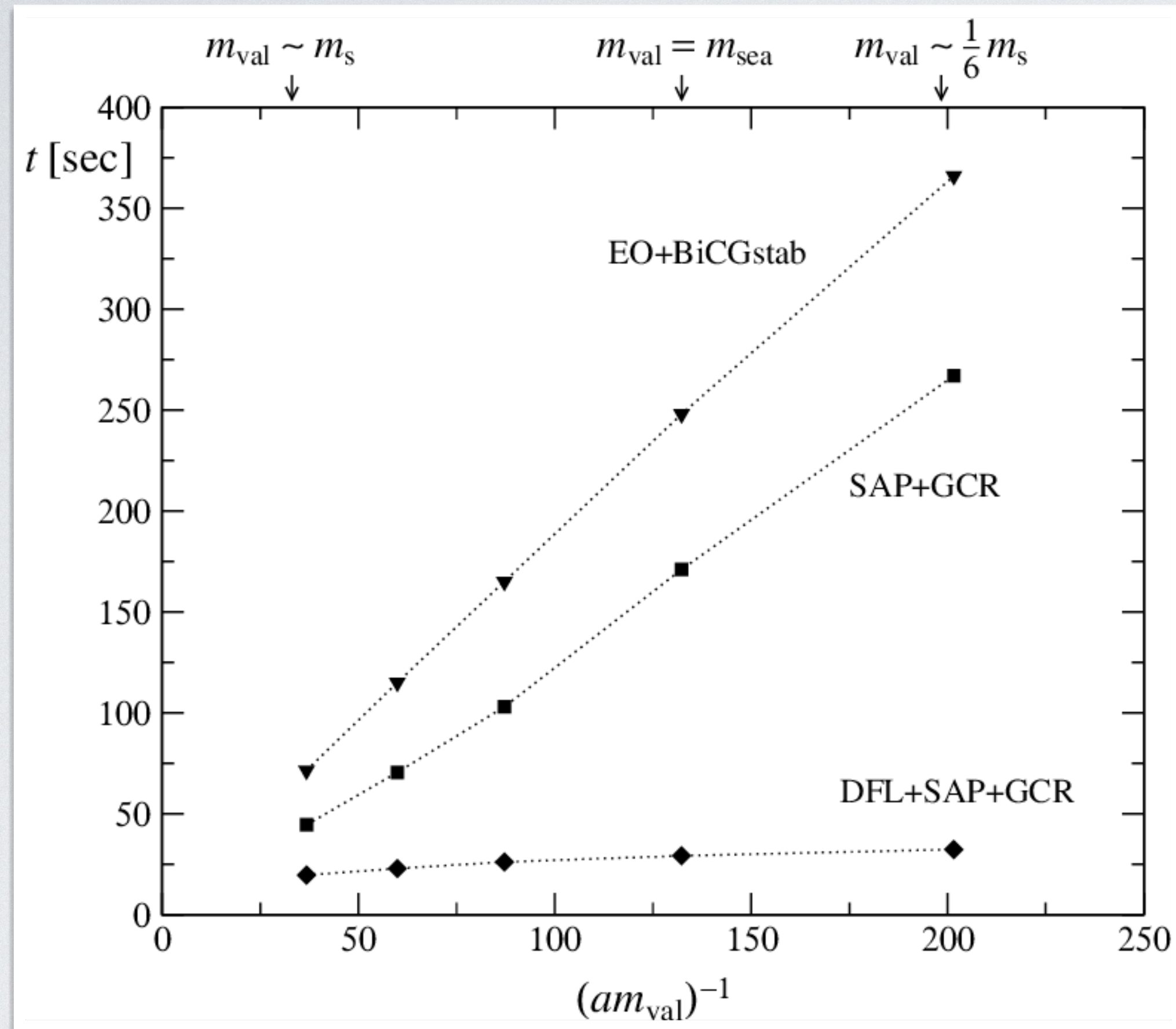
**Ukawa (2001)**

**Luscher (2004)**





# Why openQCD?



[JHEP0712:011,2007]

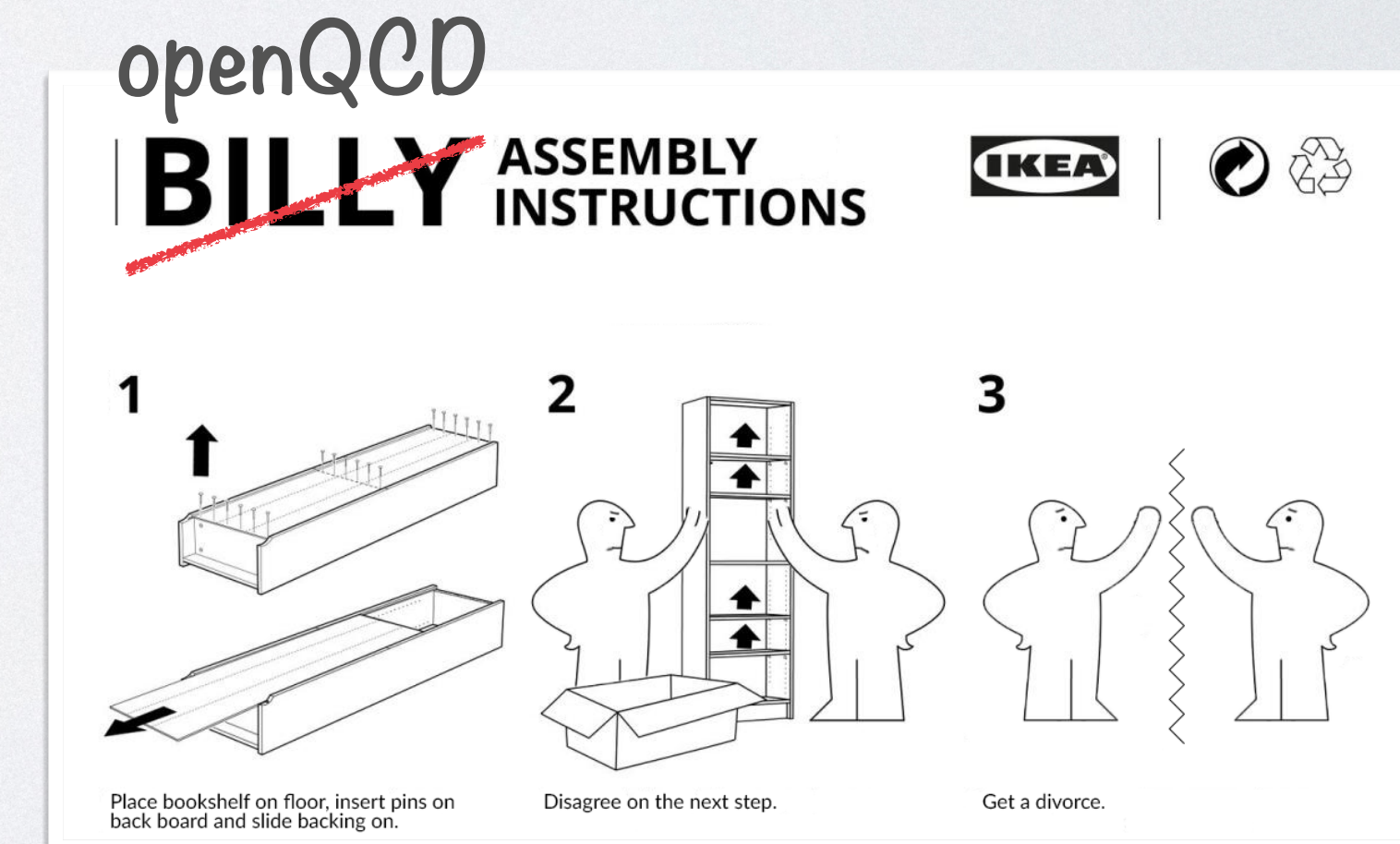
Domain Decomposition in HMC (2003-2004)

Deflated DD HMC in 2007



# Develop C++/CUDA modules to replicate openQCD functions on the GPU

- Module Development:
  - Create a comprehensive set of C++/CUDA modules mapping the CPU (relevant) functions to GPU. Ensure these modules accurately reproduce the functionality of openQCD on the GPU.
- Data Structure Translation:
  - Maintain consistent and coherent translation between GPU and CPU data structures. Ensure seamless data flow and compatibility (of fields).
- Function Naming Conventions:
  - Adopt a clear and consistent naming convention. Example: `dfl_subspace` in CPU becomes `dfl_subspace_gpu` in GPU.
- GPU only:
  - Implement all the computational part of inversion cycle on the GPU. Future Goals: Extend this to include HMC and SMD processes.

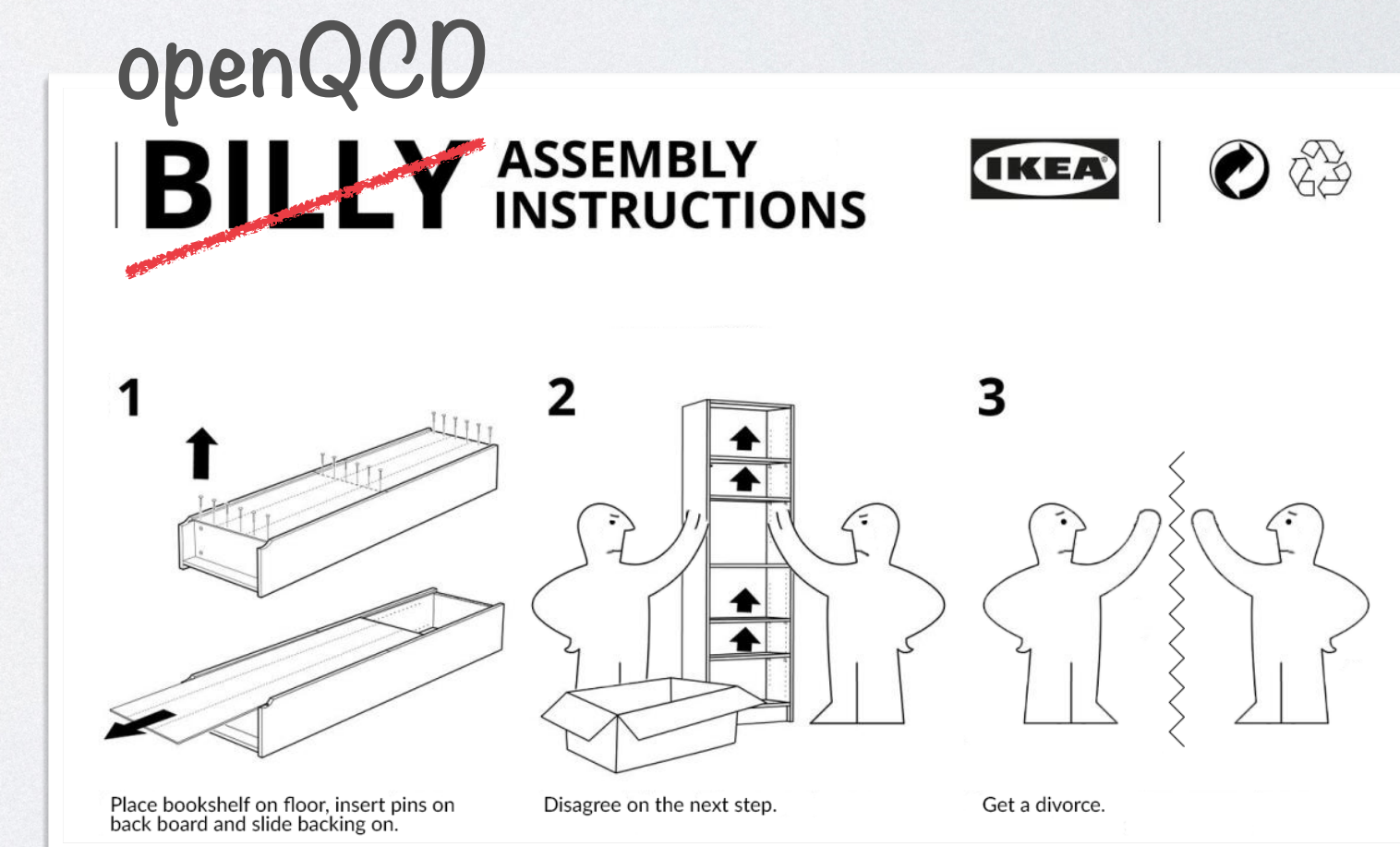




# A pragmatic approach

- Redesigned Data Structures:  
Allow full memory coalesced access in GPU data structures and memory access.
- Code Readability:  
Prioritised code readability and usability.  
Avoided optimising non-critical routines (e.g., not all data structures have full memory coalescence).
- CUDA Grid Mapping:  
Mapped CUDA grid of blocks to geometrical blocks for efficient computation.
  - Block Reductions:  
Implemented block reductions using CUB library for optimised performance.
- Database Duplication:  
Duplicated the openQCD database for GPU use.  
field status is propagated between CPU and GPU databases.

- Quadruple Precision Reduction:  
Implemented quadruple precision reduction to ensure computational accuracy.
- Random Number Generation:  
Integrated directly on the GPU robust random number generation for simulations.
- Validation:  
Guaranteed at field level for each openQCD test





# The test ground



## Hardware Configuration of the JUWELS Booster Module

- 936 compute nodes
  - 2x AMD EPYC Rome 7402 CPU, 2x 24 cores, 2.8 GHz
  - Simultaneous Multithreading
  - 512 GB DDR4, 3200 MHz
  - 4x NVIDIA A100 GPU, 40 GB HBM2e
  - 4x InfiniBand HDR (Connect-X6)
  - Diskless

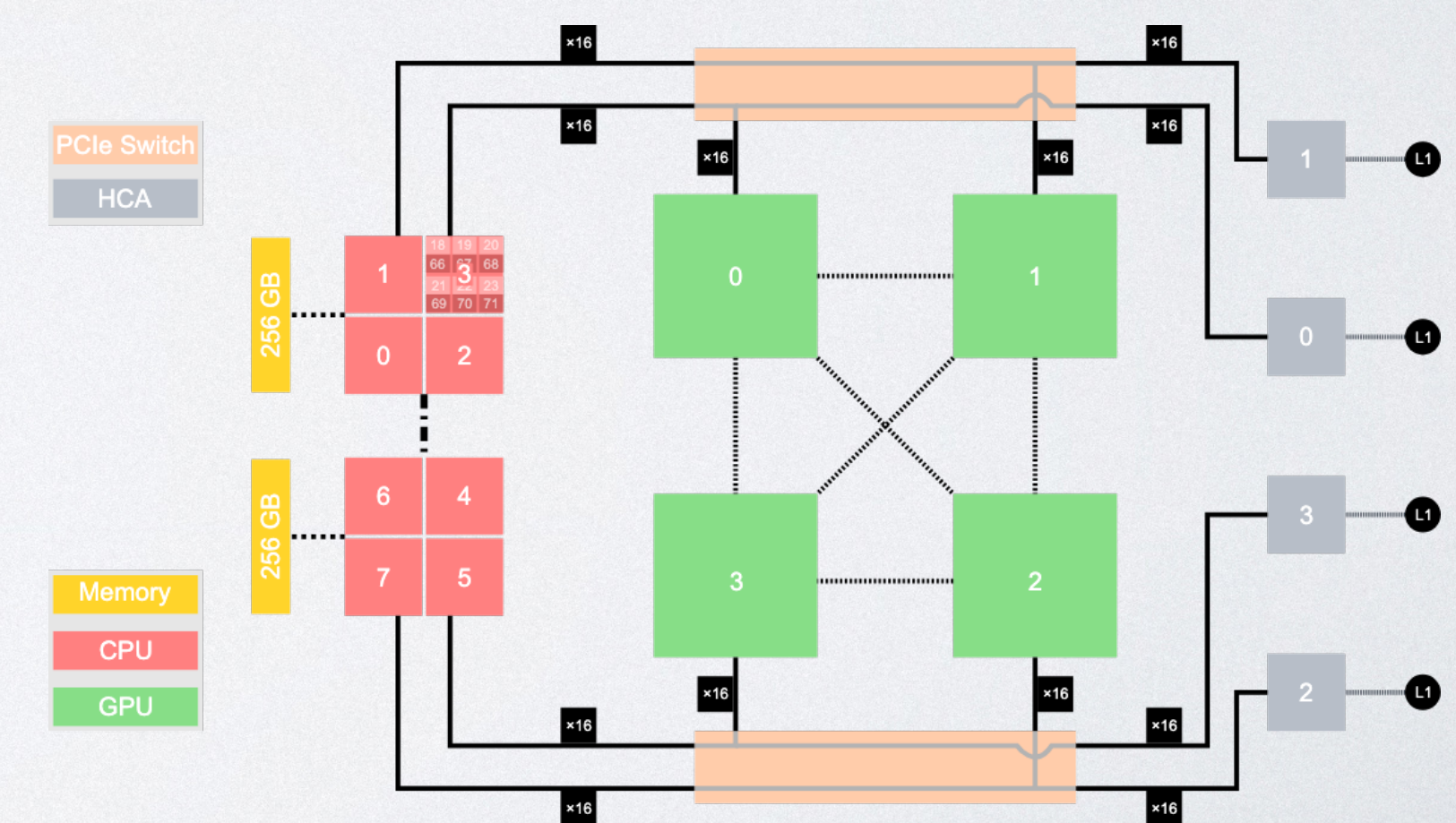
	GPU0	GPU1	GPU2	GPU3	NIC0	NIC1	NIC2	NIC3	CPU Affinity	NUMA Affinity	GPU NUMA ID
GPU0	X	NV4	NV4	NV4	PIX	SYS	SYS	SYS	18-23,66-71	3	N/A
GPU1	NV4	X	NV4	NV4	SYS	PIX	SYS	SYS	6-11,54-59	1	N/A
GPU2	NV4	NV4	X	NV4	SYS	SYS	PIX	SYS	42-47,90-95	7	N/A
GPU3	NV4	NV4	NV4	X	SYS	SYS	SYS	PIX	30-35,78-83	5	N/A
NIC0	PIX	SYS	SYS	SYS	X	SYS	SYS	SYS			
NIC1	SYS	PIX	SYS	SYS	SYS	X	SYS	SYS			
NIC2	SYS	SYS	PIX	SYS	SYS	SYS	X	SYS			
NIC3	SYS	SYS	SYS	PIX	SYS	SYS	SYS	X			

### Legend:

X = Self  
 SYS = Connection traversing PCIe as well as the SMP interconnect between NUMA nodes (e.g., QPI/UPI)  
 NODE = Connection traversing PCIe as well as the interconnect between PCIe Host Bridges within a NUMA node  
 PHB = Connection traversing PCIe as well as a PCIe Host Bridge (typically the CPU)  
 PXB = Connection traversing multiple PCIe bridges (without traversing the PCIe Host Bridge)  
 PIX = Connection traversing at most a single PCIe bridge  
 NV# = Connection traversing a bonded set of # NVLinks

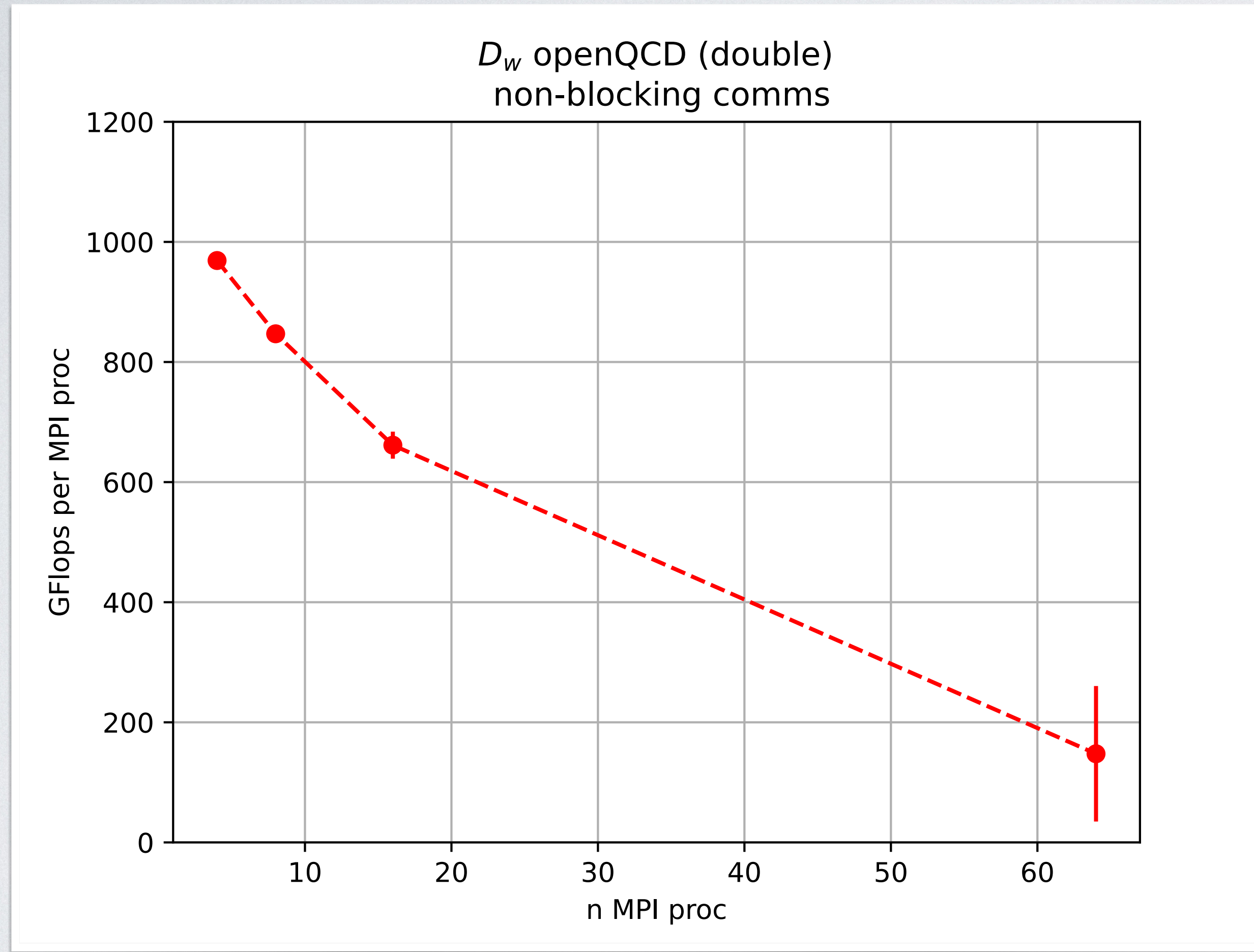
### NIC Legend:

NIC0: mlx5\_0  
 NIC1: mlx5\_1  
 NIC2: mlx5\_2  
 NIC3: mlx5\_3





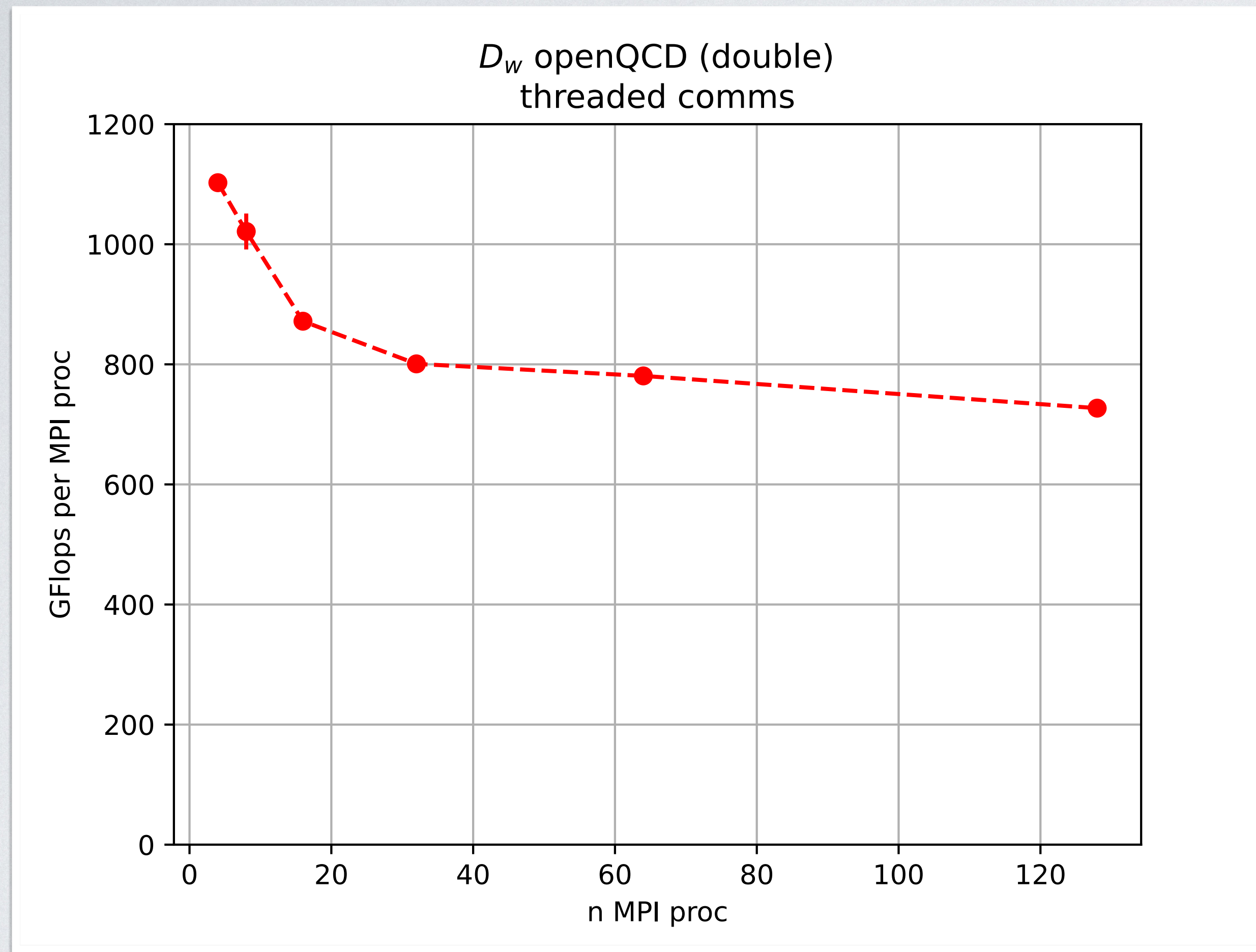
# A first test: weak scaling $D_W$ for a local $24^4$



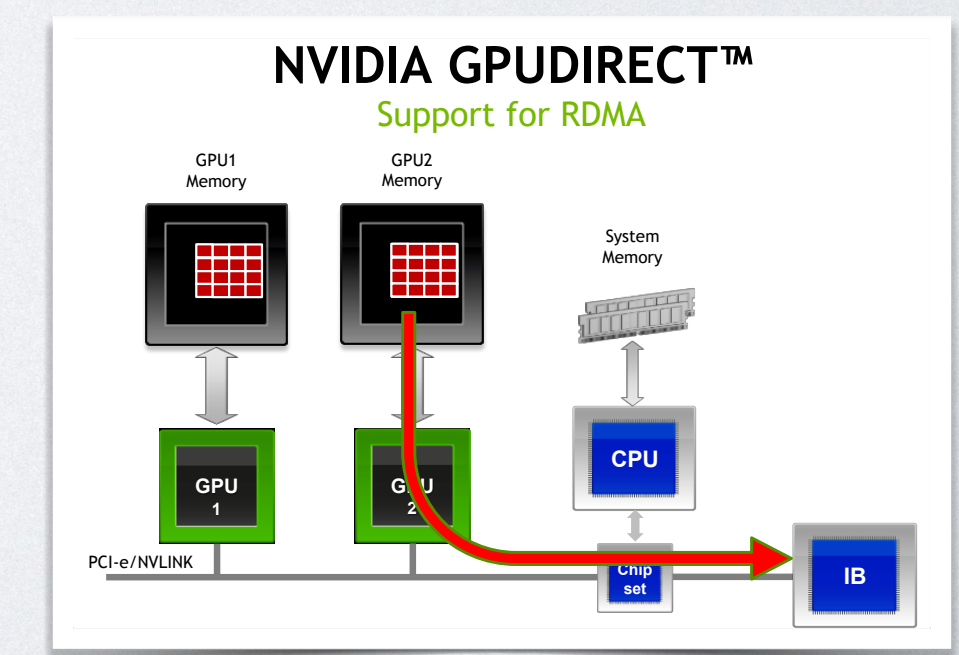
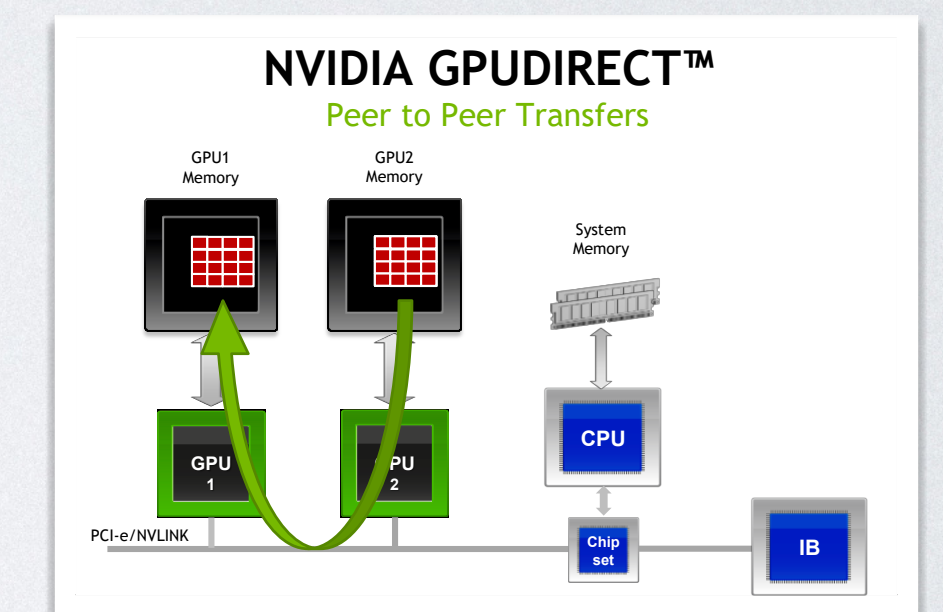
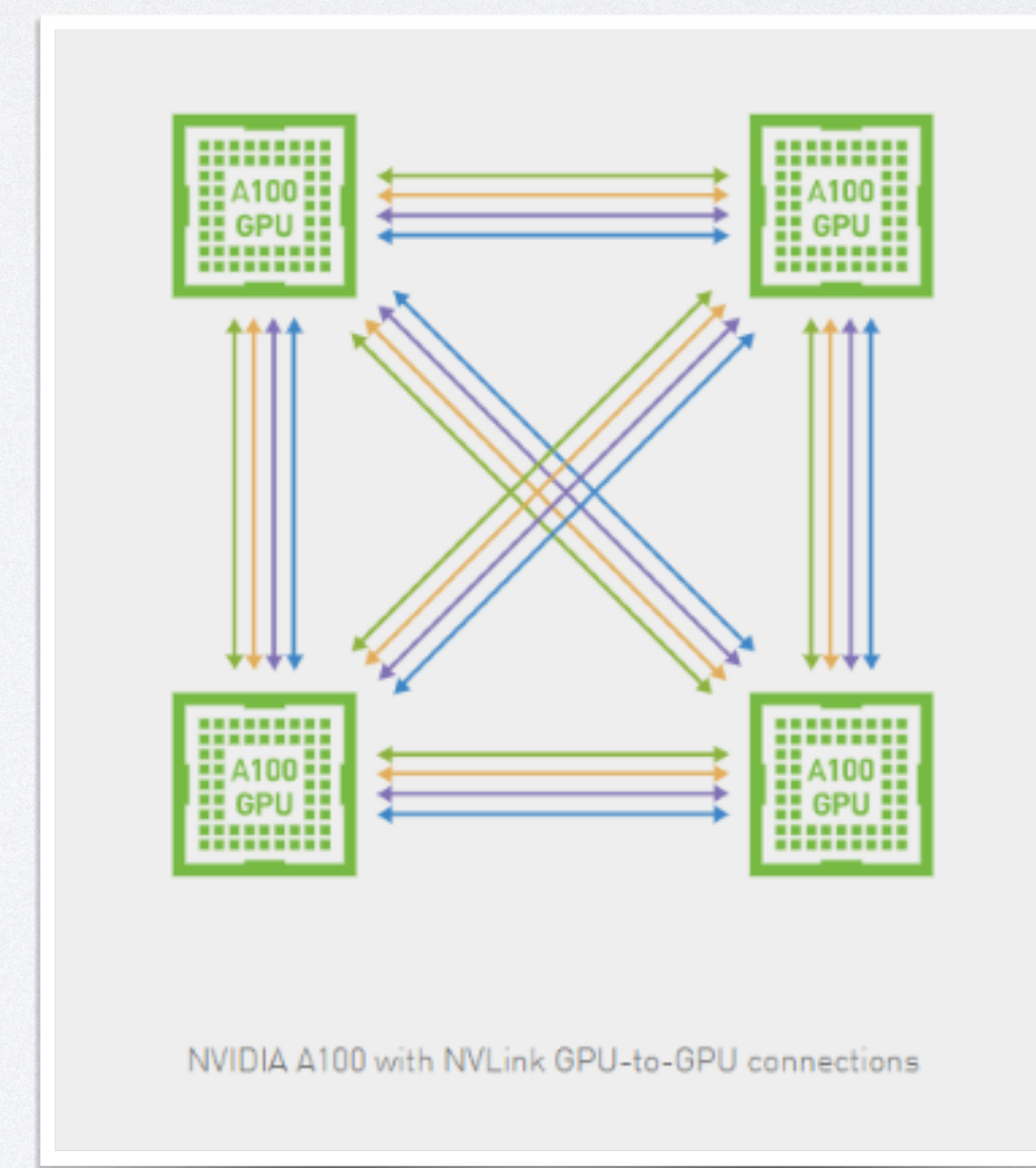
- No un-coalesced access
- EO memory arrangement for every field
- We opted for a index computation “in flight”, no lookup tables.  
Block sizes defined a compile time
- Can achieve 1 Tflops per A100 card



# A first test: weak scaling $D_W$ for a local $24^4$

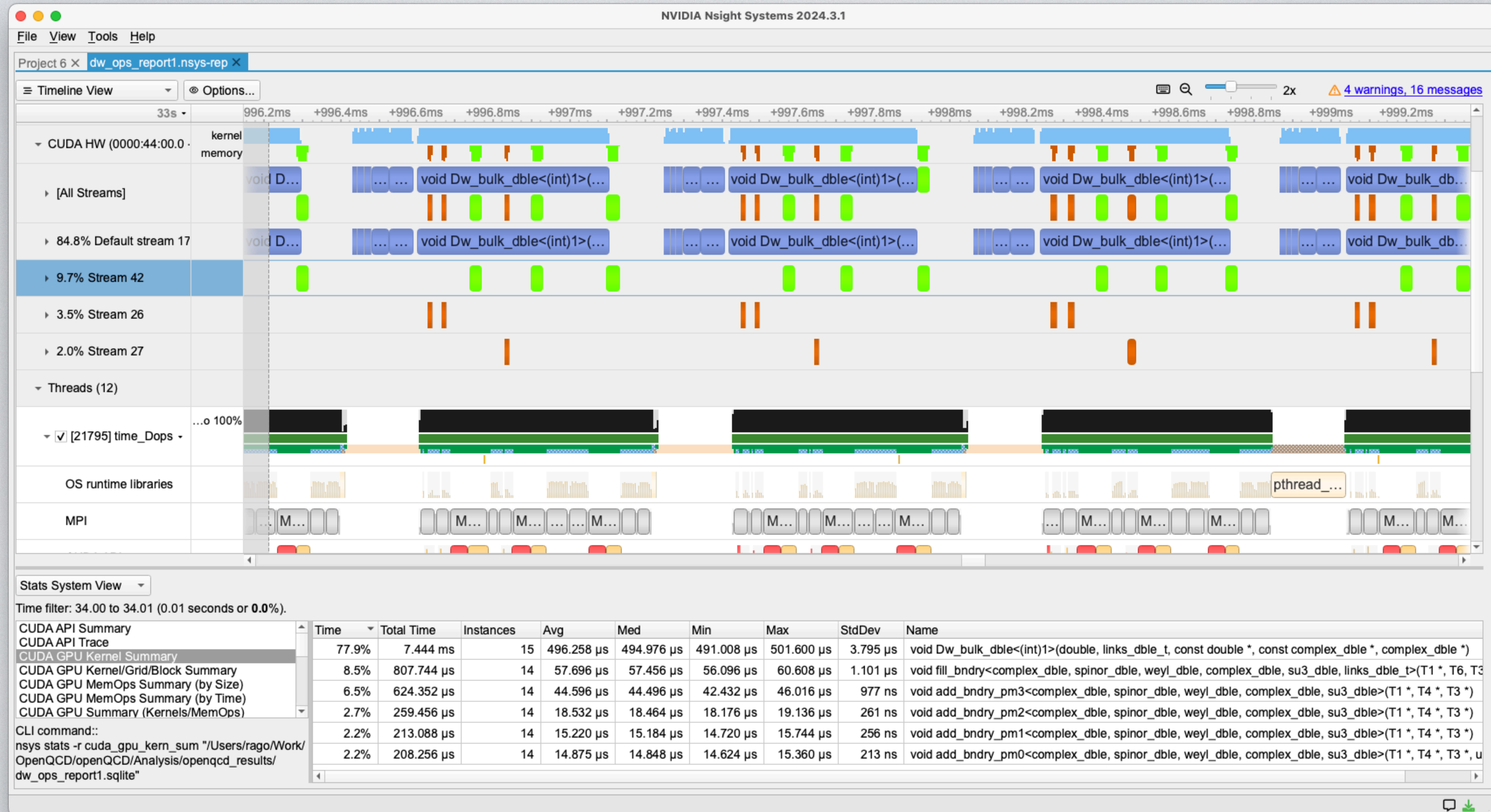


- Centralised the MPI communications and support locking and non locking strategy at compile time.
- MPI is often erratic: non-locking comms over p-threads
- Topology awareness guarantees a sensible speedup thanks to NVlink



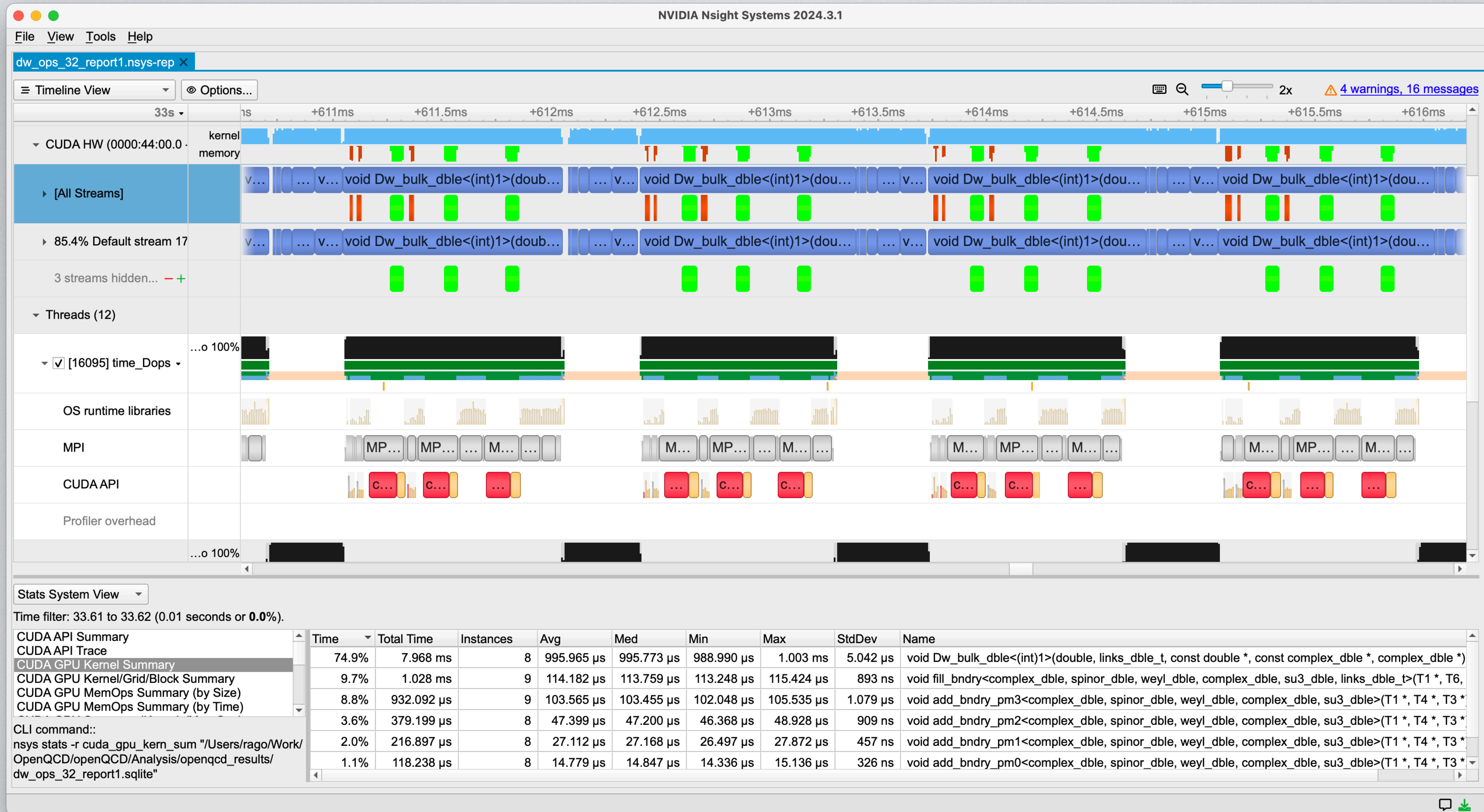


# A first test: weak scaling $D_W$ for a local $24^4$





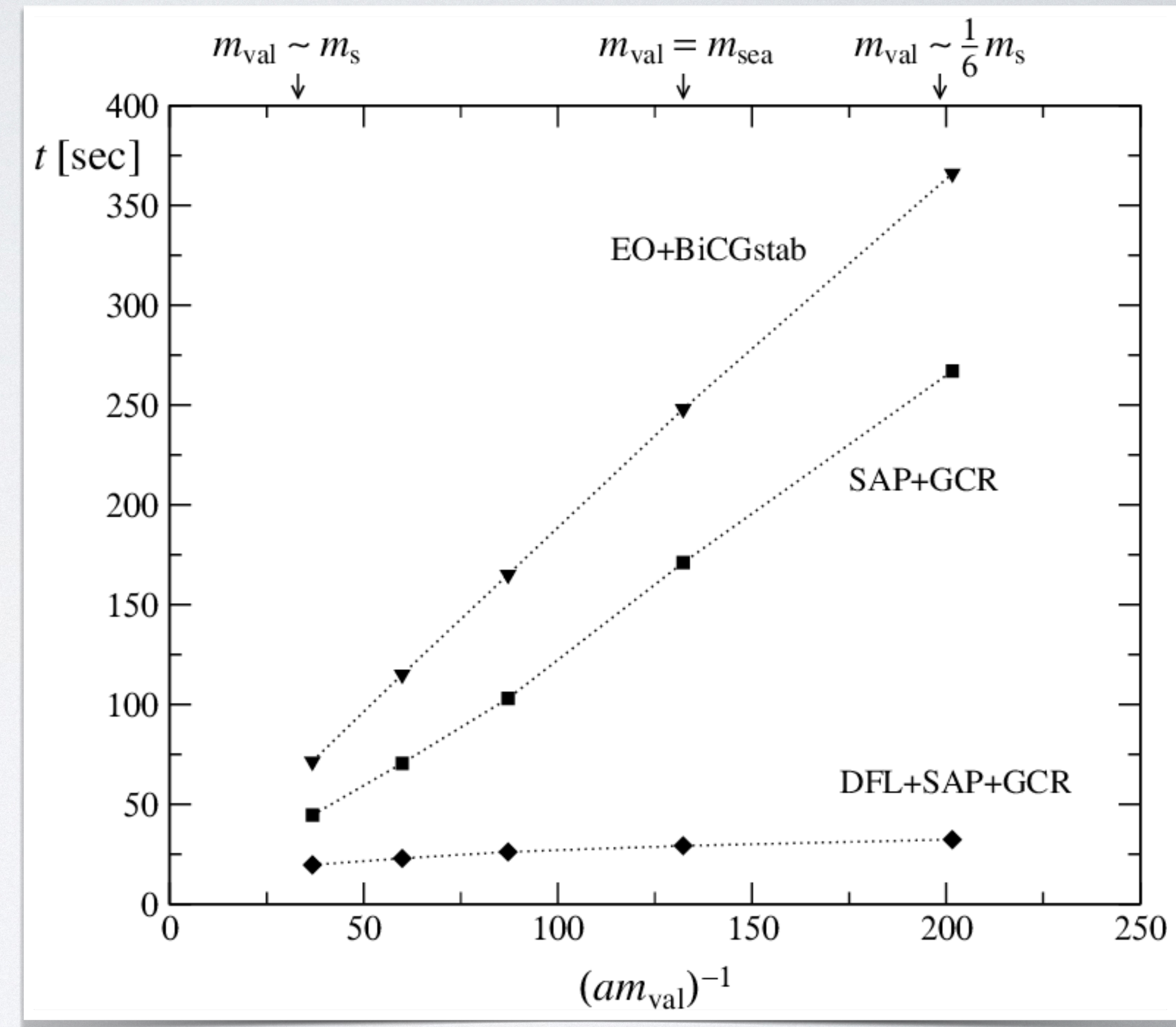
# A first test: weak scaling $D_W$ for a local $48 \times 24^3$





# A more challenging case: SAP and DFL

- SAP and Deflation are two inverters that take advantage of locality and mode coherence.
- Our strategy has been to focus on an efficient porting of the two routines.
  - block geometry allowing seamless use of different block sizes
  - Shared memory within the sap blocks
  - Block reduction through cubs thanks to the binding of the block indices to the CUDA grid

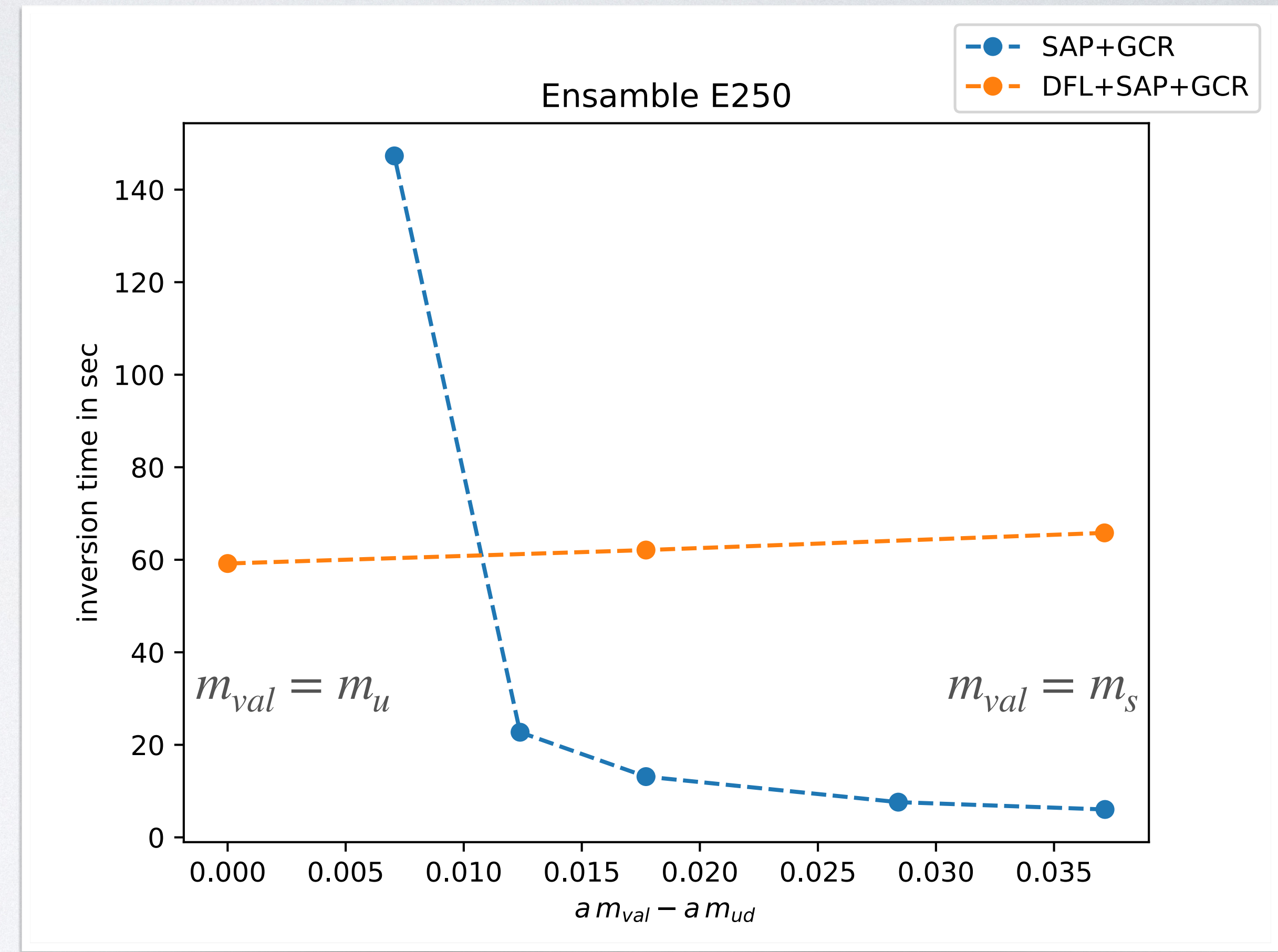




# A more challenging case: SAP and DFL

	$T \times L^3$	$M_\pi$	$M_K$	$a$
E250	$192 \times 96^3$	130 MeV	500 MeV	0.065 fm

- ✓ SAP seems to be a pretty well tuned routine, with small communication overhead.
- ✓ On the lightest point the computational cost is equally shared between the SAP and the little-GCR.
- ☐ Kernel fusion could improve the little-GCR but the overall performance is not going to change drastically.



64 Juwels booster nodes



# Modern coding standards

- git & Continuous integration
- Code coverage

The screenshot shows the GitHub repository page for 'openQCD on GPU'. The repository is private and has 16 branches and 0 tags. The latest commit is by 'latticestefan' with the message 'some fixes for merge' and 232 commits, pushed 2 days ago. The repository contains several files and folders, including '.github', 'gpu\_devel/CUDA', 'openQCD-2.4.2', '.clang-format', '.gitignore', 'LICENSE.md', and 'README.md'. The README indicates that the repository contains the openQCD port to GPU simulation code and shows two CI jobs, 'ci-gpu' and 'check-format', both in a 'passing' state. The repository is licensed under GPL-3.0 and has 0 stars, 0 watchers, and 0 forks. The language distribution is: C 85.5%, Cuda 12.2%, Makefile 2.1%, and Other 0.2%.

```
238 [check_Dw] ... OK
239 [check_Dw_bulk] ... OK
240 [check_dfl_subspace] ... OK
241 [check_set_Awhop_gpu] ... OK
242 [check_vfields] ... OK
243 [check_set_Awhop_gpu] ... OK
244 [check_set_Awblk_gpu] ... OK
245 [check_sapgcr] ... OK
246 [check_update_Awblk_gpu] ... OK
247 [check_set_Awhat_gpu] ... OK
248 [check_Aw_gpu] ... OK
249 [check_dfl_modes] ... OK
250 [check_dflsapgcr] ... OK
```

>  Complete job



Already implemented:

- ✓ SAP / DFL / GCR
- ✓ Only PBC (Periodic Boundary Conditions):
- ✓ Random Number Generator
- ✓ Quadruple Precision in Reductions
- ✓ AMD/HIP porting works and passes initial checks.

What is missing

- Open Boundary Conditions
- Exponential  $c_{sw}$
- Independent/automatic single kernel tuning
- Second stage hmc/smd (although gauge and  $c_{sw}$  forces gauge already implemented)



# Future plans and code release

## Testing and Release Strategy:

- Conduct initial testing phase with a selected group of colleagues and collaborators
- Identify and fix any preliminary bugs
- Gather feedback on usability and performances
- Ensure stability and reliability before broader release

## Second Stage: Public release

- Release the code to the public with an appropriate license upon initial publication
- Provide detailed documentation and user guides
- Include comprehensive results from testing and performance evaluations