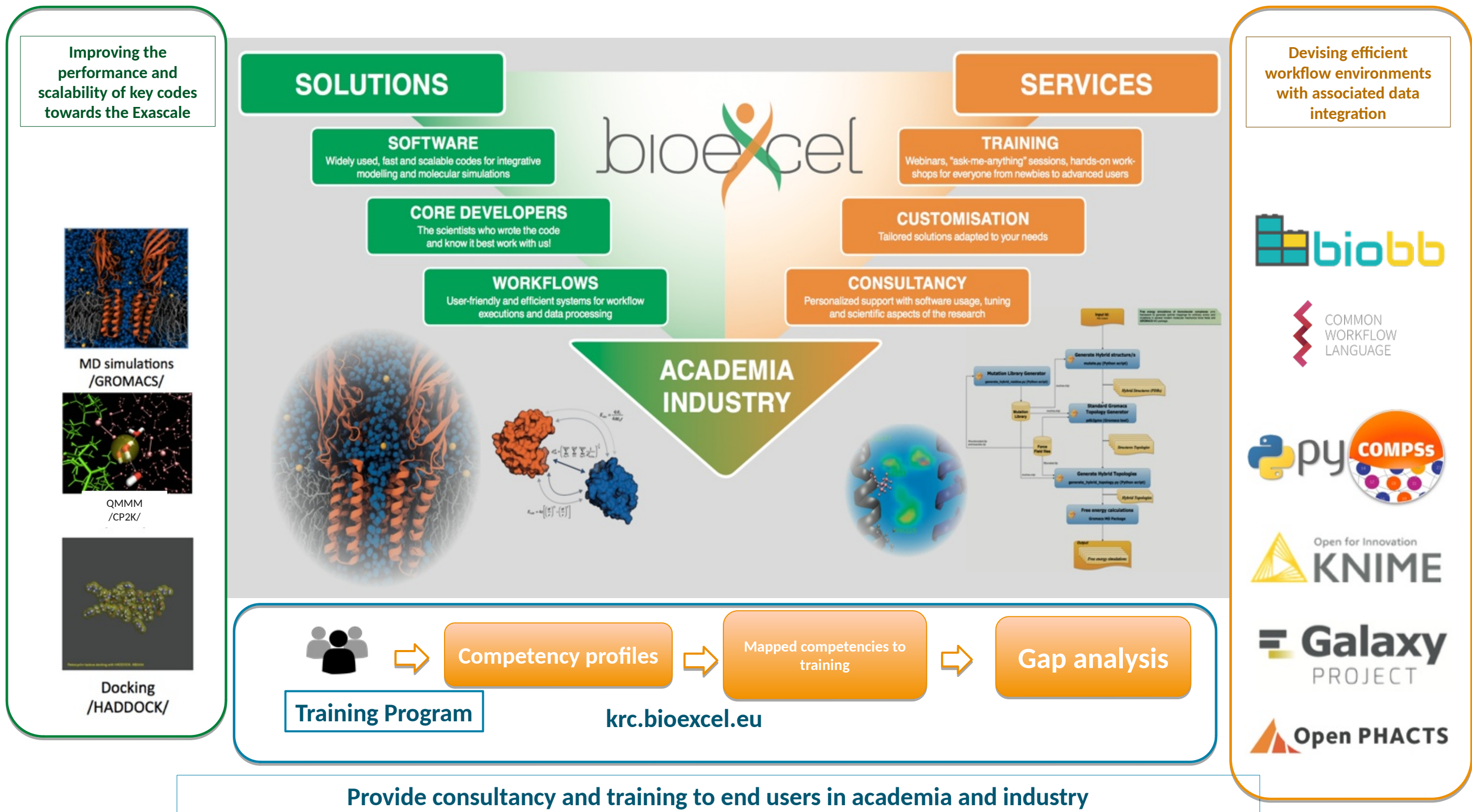# HPC codesign in GROMACS

Szilárd Páll

pszilard@kth.se
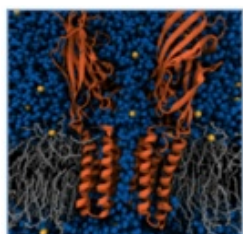
Workshop on Software Co-Design Actions in European Flagship HPC Codes
ISC 2022
June 2, 2022
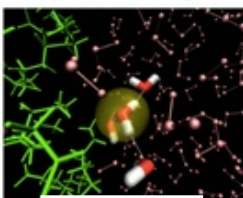
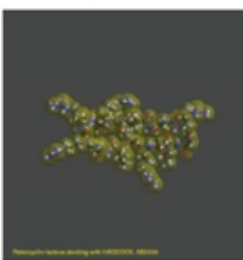# BioExcel Center of Excellence



**Improving the performance and scalability of key codes towards the Exascale**

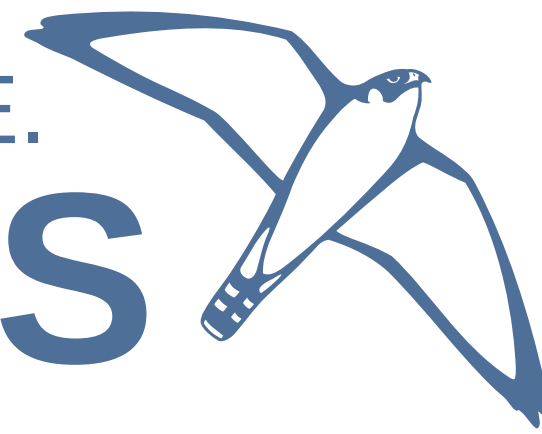MD simulations /GROMACS/

QMMM /CP2K/

Docking /HADDOCK/

**Devising efficient workflow environments with associated data integration**

**SOLUTIONS**

**SOFTWARE**
Widely used, fast and scalable codes for integrative modelling and molecular simulations

**CORE DEVELOPERS**
The scientists who wrote the code and know it best work with us!

**WORKFLOWS**
User-friendly and efficient systems for workflow executions and data processing

**SERVICES**

**TRAINING**
Webinars, "ask-me-anything" sessions, hands-on work-shops for everyone from newbies to advanced users

**CUSTOMISATION**
Tailored solutions adapted to your needs

**CONSULTANCY**
Personalized support with software usage, tuning and scientific aspects of the research

**ACADEMIA INDUSTRY**

Training Program → Competency profiles → Mapped competencies to training → Gap analysis

krc.bioexcel.eu

**Provide consultancy and training to end users in academia and industry**

Goals:
- Develop key applications (incl. GROMACS) for exascale;
- Develop workflow solutions
- Training/support to academia and industry
- Establish a long-term organizational structure

FAST. FLEXIBLE. FREE.
# GROMACS

- **Classical MD** code
  - supports all major force-fields
  - broad algorithm support

- **Development**:
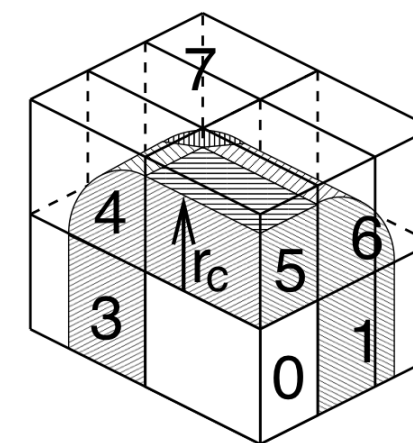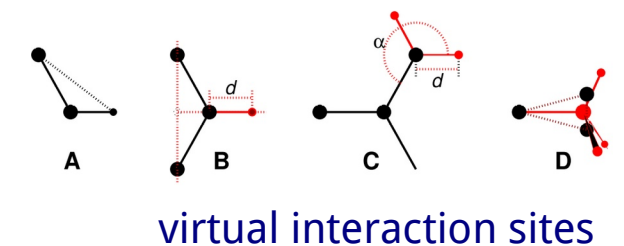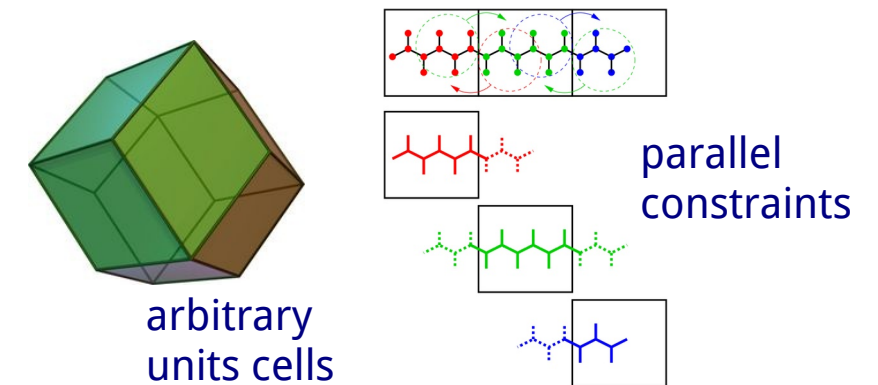
  Stockholm Sweden & partners worldwide

- **Large user base**:
  - One of the top HPC codes worldwide

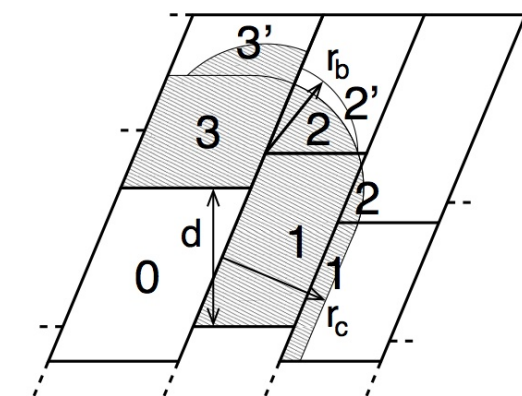    deployed on most clusters
  - 10k's academic & industry users

- **Open source**: LGPLv2

- **Open development:**
  - code review & bug-tracker: https://gitlab.com/gromacs

parallel constraints
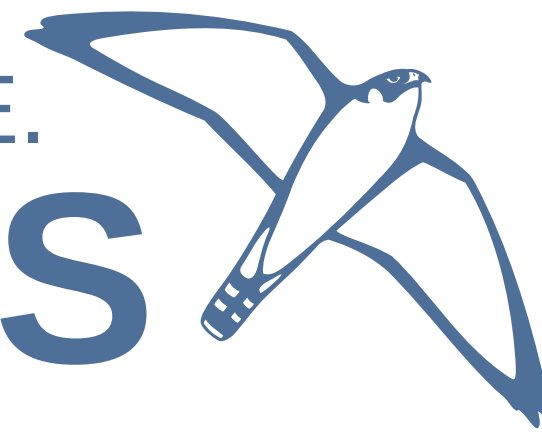
arbitrary units cells

virtual interaction sites

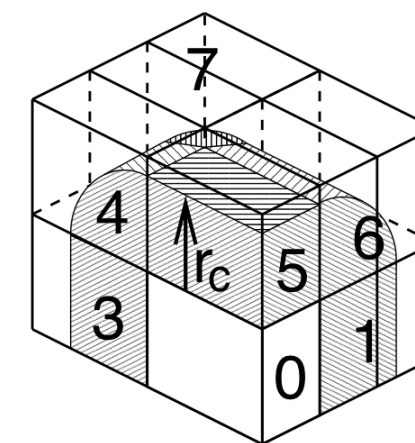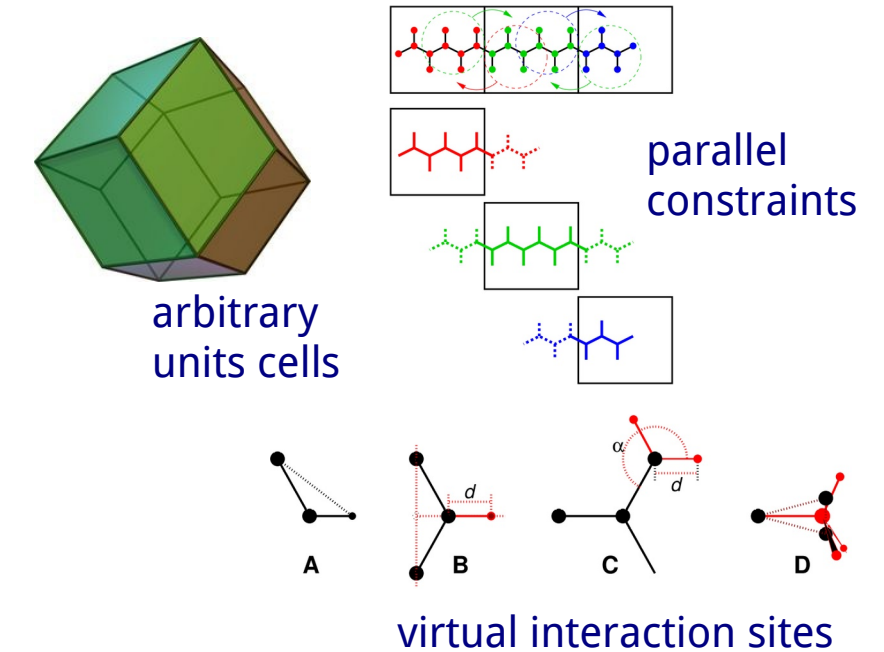Eighth shell domain decomposition

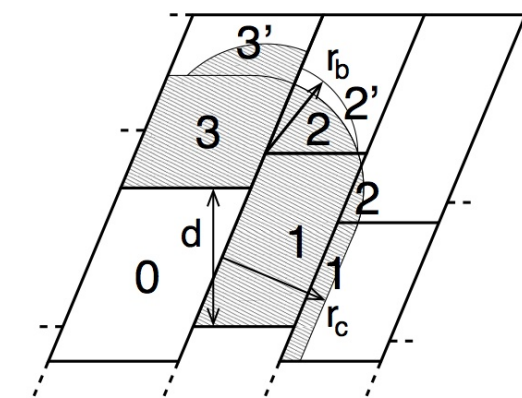Triclinic unit cell with load balancing and staggered cell boundaries

FAST. FLEXIBLE. FREE.
# GROMACS

- Focus on **high performance:**

  efficient algorithms & highly-tuned parallel code

- **Bottom-up performance oriented design:**

  – absolute performance over "just scaling"

- Focus on **portability**

  – Linux distro integration and CI

  – regular testing on all HPC arch

  – SIMD portability library, GPU abstraction layer

  – open standards-based languages/APIs

- **Modern development workflow**

  – mandatory open code review for >10 years

  – tiered CI testing / verification

parallel constraints
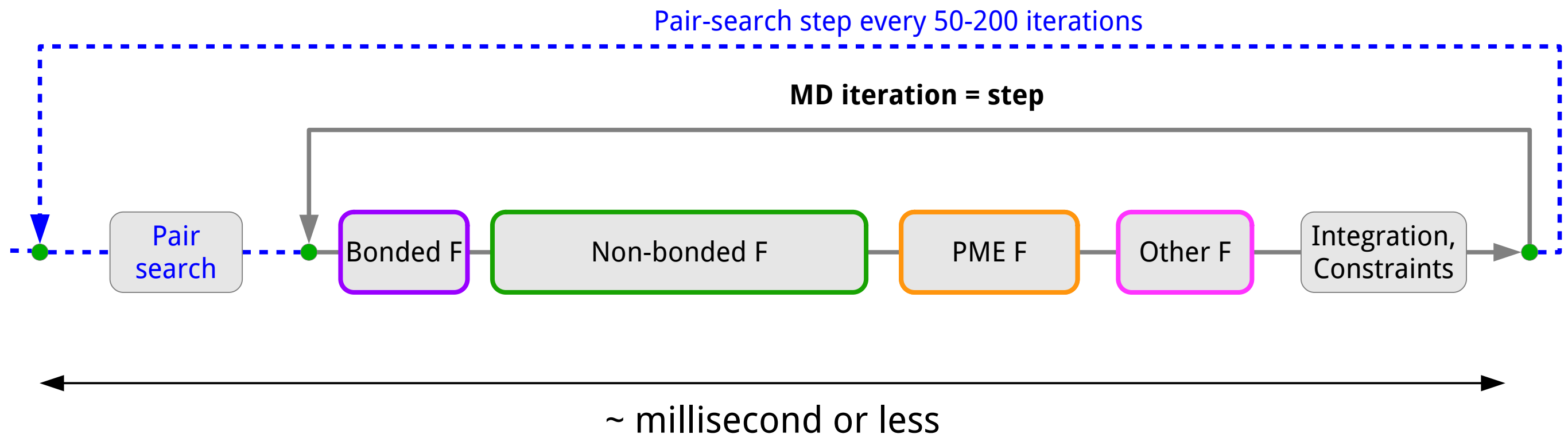
arbitrary units cells

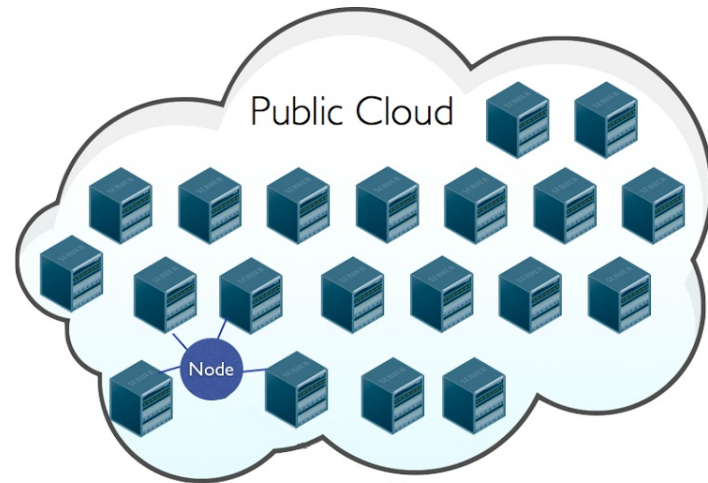virtual interaction sites

Eighth shell domain decomposition

Triclinic unit cell with load balancing and staggered cell boundaries

# MD: computational challenge

Pair-search step every 50-200 iterations

**MD iteration = step**

| Pair search | | Bonded F | Non-bonded F | PME F | Other F | Integration, Constraints |

~ millisecond or less

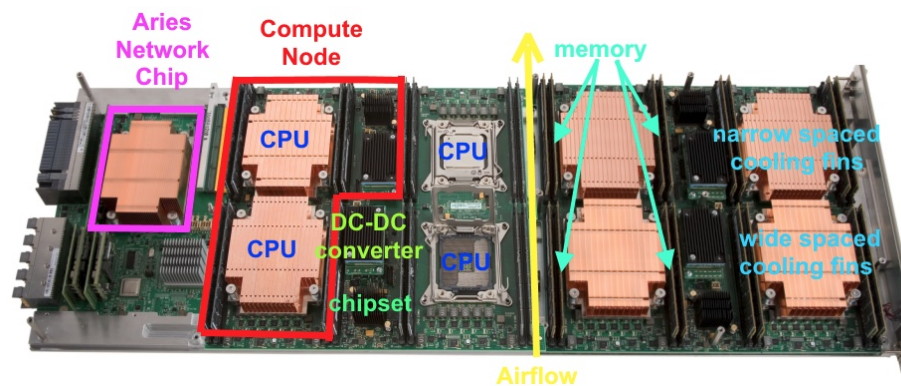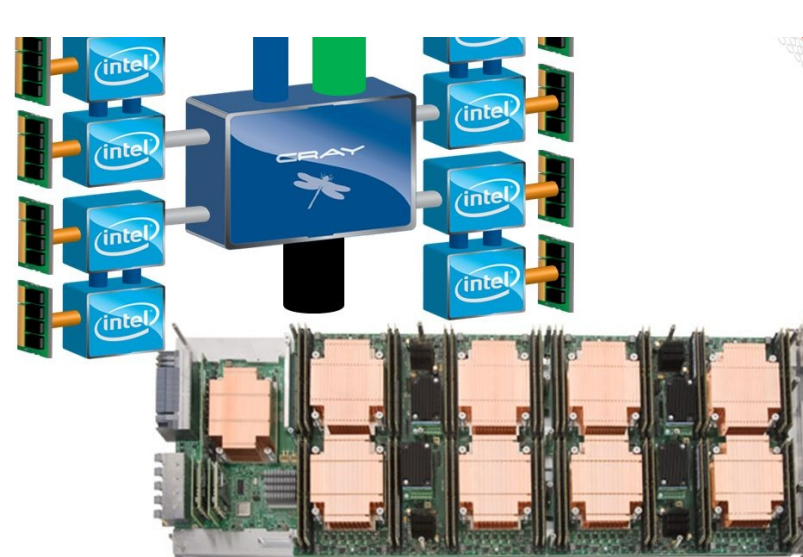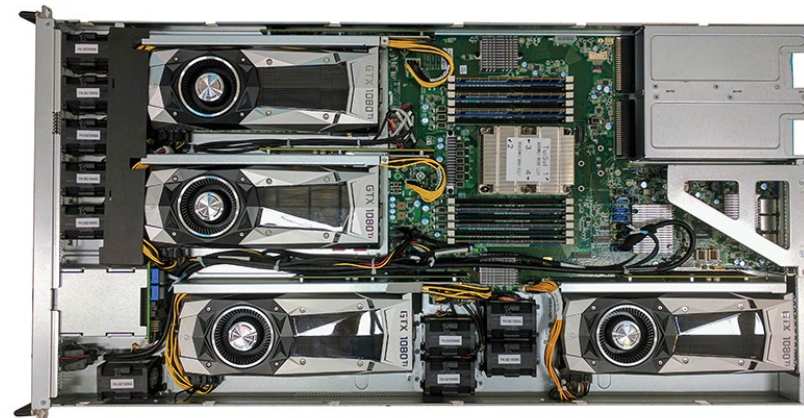- Simulation vs real-world **time-scale gap**
  - Every simulation: $10^8$ – $10^{15}$ steps
  - Every step: $10^6$ – $10^9$ FLOPs
- Main goal of parallelization:
  - study molecular systems: tackle the time- or length-scale challenge
  - typically requires: **strong scaling**, increasingly **ensemble**
- MD codes at peak: ~**100 µs / step** (on commodity hardware)
  - <100 atoms/core at peak
  - <10000 atoms / GPU

# Multiple levels of hardware parallelism



Public Cloud

Node

up to 512-bit vector units/core
=>
up to 16 single precision ops/clock

AMD EPYC

Aries Network Chip

Compute Node

memory

CPU

CPU

CPU

CPU

DC-DC converter

chipset

narrow spaced cooling fins

wide spaced cooling fins
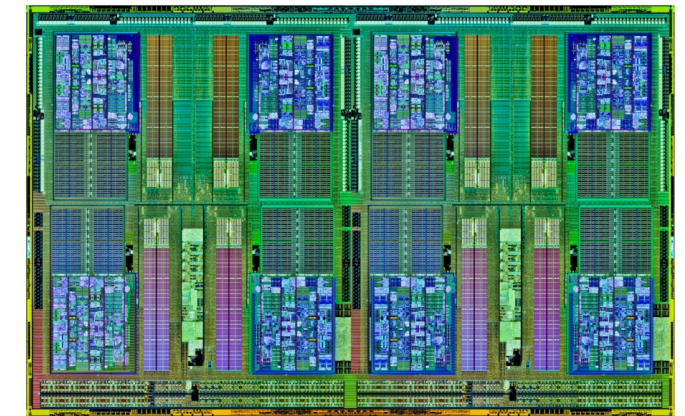
Airflow

nvidia TESLA

**Compute cluster or cloud**
Networked computers:
topology, bandwidth, latency

**Compute node / workstation**

NUMA topology, PCIe
Shared under CC BY-SA 4.0. doi.org/10.5281/zenodo.6620848

**Multicore CPU & manycore GPU**
caches, interconnects

# Multiple levels of **hardware parallelism**
# Multiple levels of **parallelization**



- Mapping the problem to the hardware:

   **expose parallelism** (algorithms) & **express parallelism** (implementation)

- Need to choose the right:

   **granularity** & **abstraction** (problem & hardware-specific)

# HPC nodes today/soon



CSC Puhti :2
CPU + 4 GPU+
NVlink, 2 NIC

JUWELS-Booster: 2 CPU +
4 GPU w NVlink + 4 NIC

DGX A100

AMD Exascale
architecture: LUMI,
Frontier, Dardel

Intel Exascale
architecture: Aurora

ORNL Summit:
NVLink
CPU/GPU-GPU

# Multiple levels of **hardware parallelism**
# Multiple levels of **parallelization**

## Exascale challenge:

- Increasing **parallelism**

  → need to express more concurrency

- Increasing **complexity** (interconnects, memories, NUMA)

  → tackle using runtimes or in application?

- Increasing **diversity**

  → zoo of programming models

  → algorithms, portability/testing, performance portability

- **Heterogeneity** is here to stay

  • ignore or embrace?

  • Wait for integration or tune for many generations?

**Compute cluster or cloud**
Networked computers:
topology, bandwidth, latency

**Compute node / workstation**

NUMA topology, PCIe
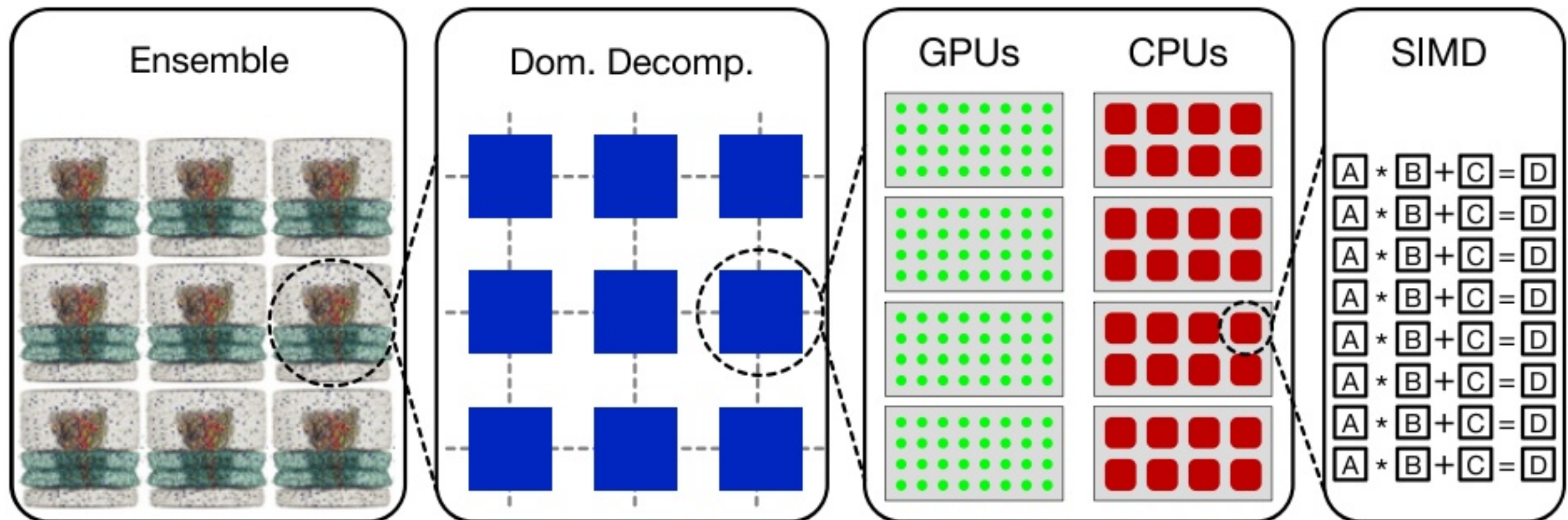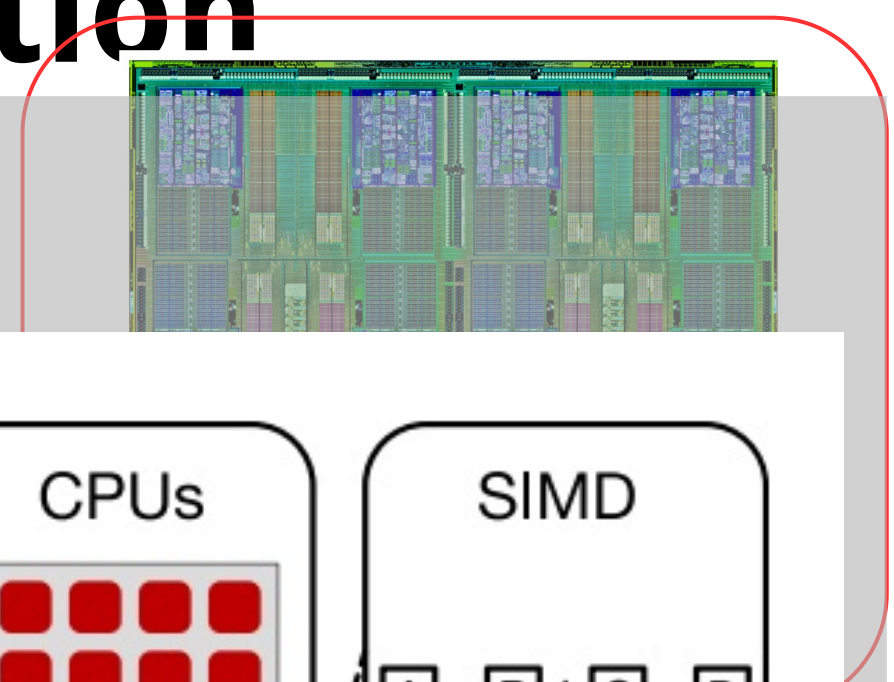Shared under CC BY-SA 4.0. doi.org/10.5281/zenodo.6620848

**Multicore CPU +
manycore GPU**
caches, interconnects

# GROMACS parallelization

- **Multi-level hierarchical** parallelization**:** target each level of hardware parallelism individually

  - Intra-node:

    - OpenMP multi-threading
      - static loop schedule, cache optimized work decomposition layout, sparse reductions
    - SIMD C++ library abstraction:
      - 14 flavors supported
    - GPU abstraction layer
      - CUDA, OpenCL, SYCL
    - thread-MPI: pthreads-based MPI for ease of use

  - Inter-node:

    - MPI: SPMD / MPMD
    - Dynamic load balancing, task balancing

# Why codesign?

- interdisciplinarity challenge

    → many hard problems need cross-disciplinary solutions

- MD: need for performance

    – GROMACS: design focus

- portability

    – GROMACS design focus

- hardware evolution...

# GROMACS & codesign

- Petascale → Exascale

  – required algorithm & parallelization redesign

  – Codesign has been & remains core component

- Physics / math + algorithms + HW

  – mainly intra-team/community

  – innovate (reformulate algorithms, accuracy-based algorithms)

  – enable (domain experts method dev, CS experts micro-bench / port)

- Algorithms + HW + vendors / CS-experts

  – mainly inter-team collaboration

  – **align goals** for collaboration so **benefits both ways**!

  – Long-term: many steps forward and several major successes

# Algorithm redesign for modern architectures

## Cluster pair-interaction algorithm for SIMD/SIMT



4x4 setup on SIMD-16



## Accuracy-based automated list buffer improves SIMD algorithm parallel efficiency



## Multi-level heterogeneous data and task load-balancing: intra-GPU, intra-node, inter-node



Regularized lists: balanced execution



Short-range cut-off 0.9 nm

Coulomb

Lennard–Jones

Increase cut-off → increase grid spacing

Short-range: non-bonded

PME to PP LJ cut-off fixed

Long-range: PME

LJ-PME

## Dual pair list with dynamic pruning



Pair-search step every 50-400 iterations

List pruning every 5-15 steps

MD iteration = step

DD/Pair search — Dyn. Prune — Bonded F — Non-bonded F — PME F — Other F — Integration, Constraints

# Embracing heterogeneity

- Heterogeneous design at the core:
  - "somewhat" complex schedule.
    → "But there is also always some reason in madness."
    - Heterogeneity for performance &
      **flexibility:** think of the (sometimes) silent codesign partners, method devs

# Dual pair list

- Trading costly data regularization for force computation not ideal!

- Instead: keep regularized particle data longer, **shift the cost trade-off**

- Use two buffers and lists:

  <span style="color:red">outer</span> / <span style="color:green">inner</span>

- Periodically re-prune

  <span style="color:red">outer</span> → <span style="color:green">inner</span>

- List lifetime / search frequency:
  - <span style="color:red">outer</span> list less frequently (costly)
  - <span style="color:green">inner</span> list more frequently (cheap)



$r_{list\text{-}outer}$

$r_{cut}$

$r_{list\text{-}inner}$

# Accuracy-based balancing: **dual pair list** reducing decomposition & search cost



Pair-search step every 20-100 iterations

MD iteration = step

DD / Pair search → Bonded F → Non-bonded F → PME F → Other F → Integration, Constraints

**Trade** search/DD cost → non-bonded pair interaction cost

Pair-search step every 50-400 iterations

List pruning every 5-15 steps

MD iteration = step

DD/Pair search → Dyn. Prune → Bonded F → Non-bonded F → PME F → Other F → Integration, Constraints

**Instead:**
Trade search/DD cost → dynamic re-pruning
Keeping the non-bonded pair interaction cost (near) constant

# Vendor-collaboration codesign: long-term practice

**Change in CUDA runtime API overhead**
between driver v295 and v331

**CUDA runtime API overhead vs input size**

CUDA RT overhead (us/step)

GTX 680
K20

System size (1000-s of atoms)

ns/step

measurement/driver version

Legend (power chart): CPU 2TPC, GPU, PKG rest

power consumed (W)

4C + GPU
2C + GPU
3C
1C

CPU 2TPC, Rest 2TPC

per-step nonbonded/total time (ms)

4C + GPU
2C + GPU
3C
1C

Nvidia Tesla P100, cuFFT 9.0
Nvidia Tesla V100, cuFFT 9.0
Intel Core i7-5960X, FFTW 3.3.6-pl2

3D FFT time per step (ms)

**Estimated scaling improvement with CUDA priorities**

Reaction-field, 192k water system, cut=0.9

18.5%

Performance (ns/day)

standard - two streams
single stream w/o non-local synch

#GPUs

Iris 655 sync / 1T CPU 2.74 GHz
i7-8559U CPU 2TPC@2.74 GHz
Tegra X1 @994.5 MHz
Tesla P4 @1113 MHz

91.47
144.92
179.75
15.34

peak kernel performance (ns/atom)

| CUDA RT call | single-stre (us/step | |
|---|---|---|
| H2D pair list | 0.3 | |
| H2D x+q | 7.2 | |
| NB kernel | 9.8 | |
| D2H forces | 5.2 | 9.2 |
| D2H E+shift F | 1.4 | 1.4 |
| cudaStreamSync | 3.4 | 2.3 |
| Clear kernel | 9.2 | 8.8 |
| **Total** | **37.2** | **49.5** |

# Direct GPU communication

- Alan Gray & Gaurav Garg (NVIDIA)

- Goal:
  - avoid CPU staging, accelerate critical path
  - target intra-node interconnects, e.g. NVLInk

- Two flavors:
  - thread-MPI: single-node (since 2021)
    - P2P copies (put/get), exchange CUDA events allows remote sync
    - Single process + multiple GPUs: bottlenecks required CUDA driver threading optimizations
  - CUDA-aware MPI: multi-node (since 2022)
    - requires host sync before issuing MPI call

**DGX 1V**



**DGX A100**



Shared under CC

# Multi-GPU/rank force offload scheme

# Multi-GPU/rank GPU-resident scheme

# Multi-GPU/rank GPU-resident scheme

# Multi-node GPU resident & direct GPU communication



remote MPI rank

remote MPI rank

Pair-search & domain-decomposition step

MPI comm: receive NLoc x

MPI comm: send NLoc F

**Transfer directly between GPU—GPU (or CPU—GPU) memory spaces**

CPU

DD comm

DD

Loc PS

Nloc PS

H2D
Loc pair list

H2D
NLoc pair list

GPU

Local stream

Non-local stream (high priority)

List pruning

un-pack

Bonded F

Non-Local non-bonded F

pack

preempted by non-local

Local non-bonded F

Reduce F

Rolling prune

Clear buffers

D2H x

Spread

3D-FFT fwd

Solve

3D-FFT back

Gather

Integration Constraints

# Multi-GPU resident step:
# single-node P2P direct GPU comm

- The entire inner loop including communication can be enqueued ahead of time

  - if there is no CPU task (Other F)

  - enables more efficient scheduling

  - overlap launch cost with work

  - CUDA graphs

- Challenges:

  - integrating CPU tasks

  - load balancing



Enqueue all work for 100-400 iterations

remote MPI rank

MPI comm: receive NLoc x

MD step

Other F

H2D f

List tuning

Local non-bonded F
preempted by non-local

un-pack

Bonded F

Non-Local non-bonded F

pack

Spread | 3D-FFT fwd | Solve | 3D-FFT back | Gather

Reduce F

Rolling prune | Clear buffers

D2H x

Integration Constraints

# Multi-node GPU resident step & GPU-aware MPI comm



MPI doesn't integrate into the async tasking model!

# Direct GPU communication performance

- Major benefit on fast interconnects with GPU-resident steps

- Modest improvements on low-end interconnects

Reaction-field electrostatics:
only halo-exchange

PME electrostatics:
halo-exchange & PME MPMD communication

- Staged GPU comm.
- Direct GPU comm.
- Direct GPU comm. & GPU Update

**Scaling limited by 1 PME GPU**

System: STMV 1M atoms

Hardware: DGX-1V

Unused to mimic pre-Exascale machine architectures

# DD halo exchange peak strong scaling

- JUWELS-booster:
  - 2x24-core AMD EPYC Rome
  - 4xA100

- ~50% parallel efficiency up to 12 nodes
  - only ~20000 atoms/GPU

- Peak at 48-64 nodes:
  - >500 ns/day

**DD scaling of a large biomolecule (1M atom STMV)**

# DD halo exchange peak strong scaling

- JUWELS-booster:
  - 2x24-core AMD EPYC Rome
  - 4xA100

**Scaling large homogeneous systems to ~400 GPUs**

**DD scaling of a large biomolecule (1M atom STMV)**



Input size (million atoms)
- 1.44M
- 5.7M
- 23M
- 46M

RF JUWELS-Booster
RF JUWELS-Booster optimized mapping

#nodes (4xA100)

#nodes (4GPUs/node)

# PME decomposition

- **GROMACS team + Gaurav Garg (NVIDIA)**

- remove the limitation of single dedicated PME GPU

- 3D FFTs strong-scaling challenge:

  typical size $32^3$-$256^3$, hardly scale

- Released in 2022:

  **Hybrid mode**: FFT on CPU

- In development (upstreamed):
  - major algorithmic and parallelization optimizations
  - **HeFFT** and **cuFFTmp** for GPU-resident mode



From PP

To PP

**Single PME Rank**

| Recv X | PME Spline + Spread | 3D R2C FFT | PME Solve | 3D C2R FFT | PME Gather | Send PME F |

Decompose work to use multiple ranks / GPUs

**PME Rank 0**

| Recv X | PME Spline + Spread | Halo exchang | Distributed R2C FFT | PME Solve | Distributed C2R FFT | Halo exchang | PME Gather | Send PME F |

**PME Rank 1**

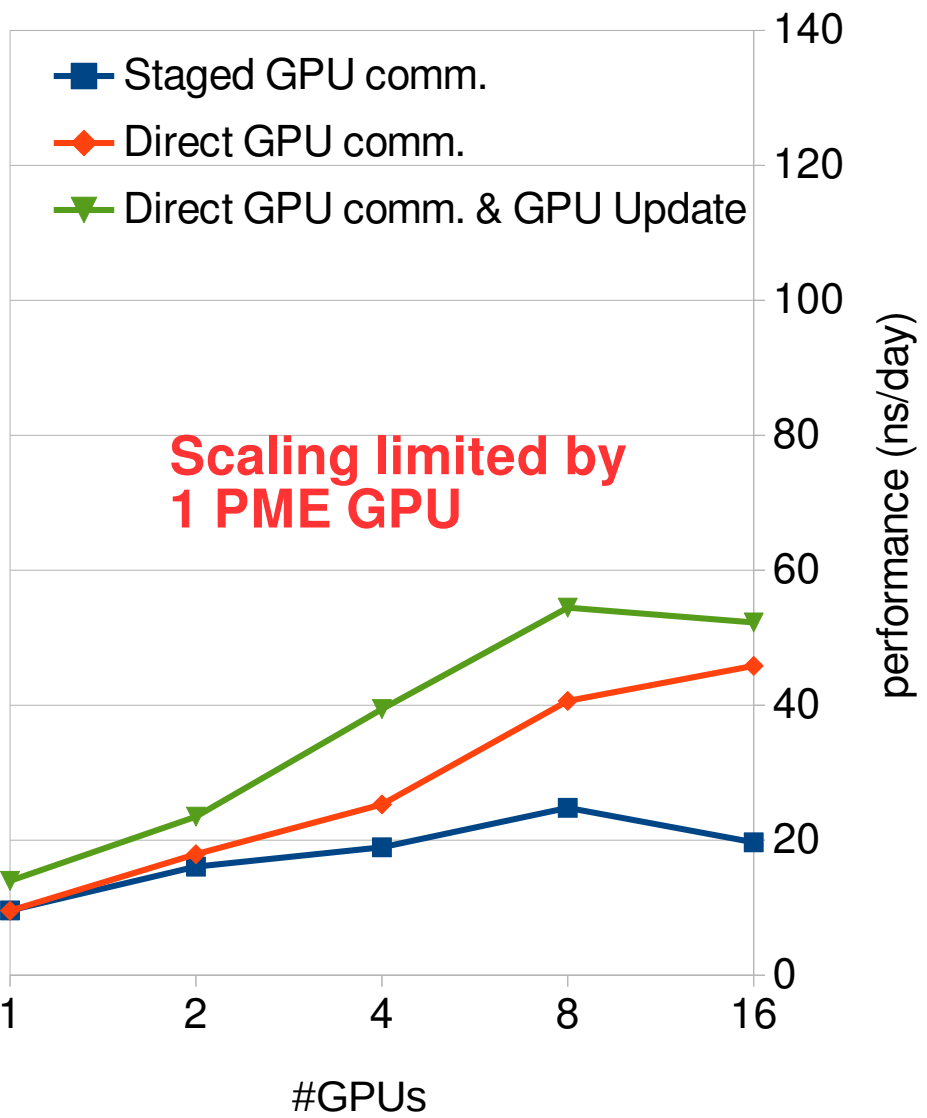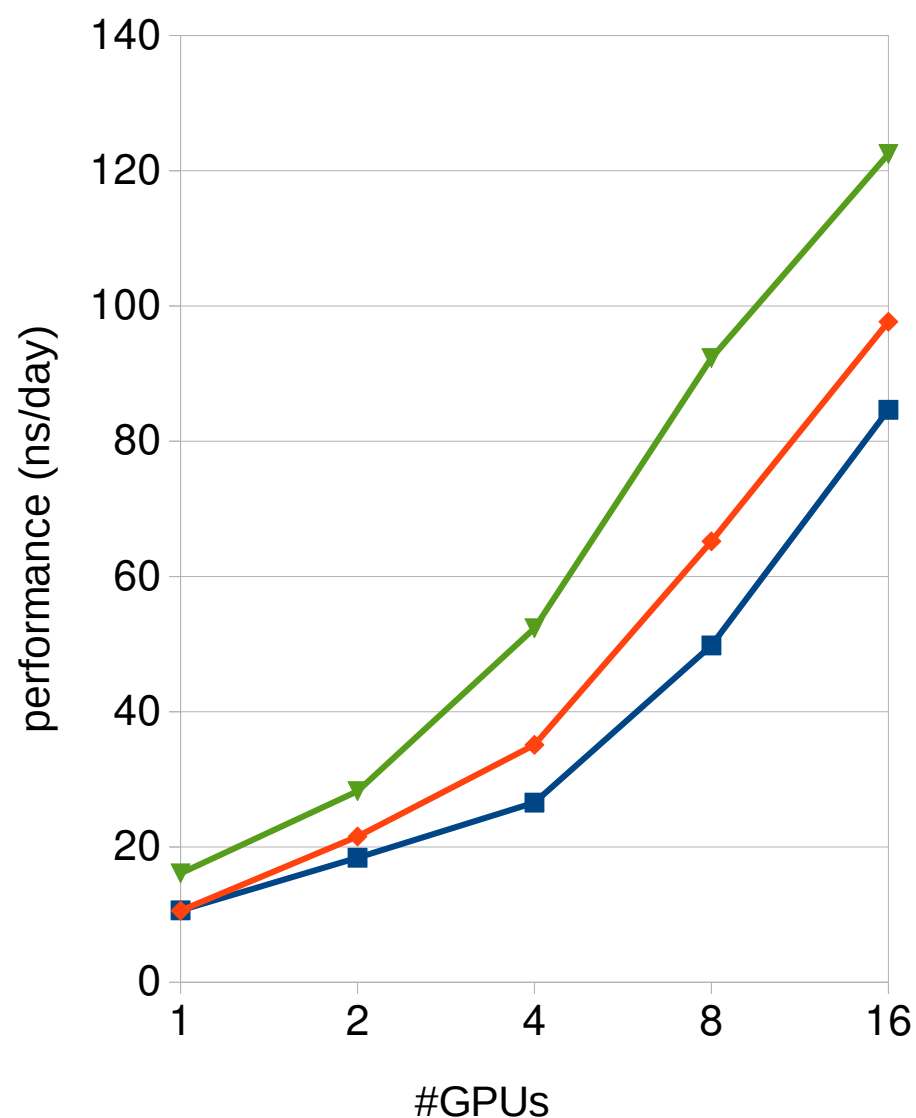| Recv X | PME Spline + Spread | Halo exchang | Distributed R2C FFT | PME Solve | Distributed C2R FFT | Halo exchang | PME Gather | Send PME F |

# Direct GPU communication with PME decomposition

- Major benefit on fast interconnects with GPU-resident steps

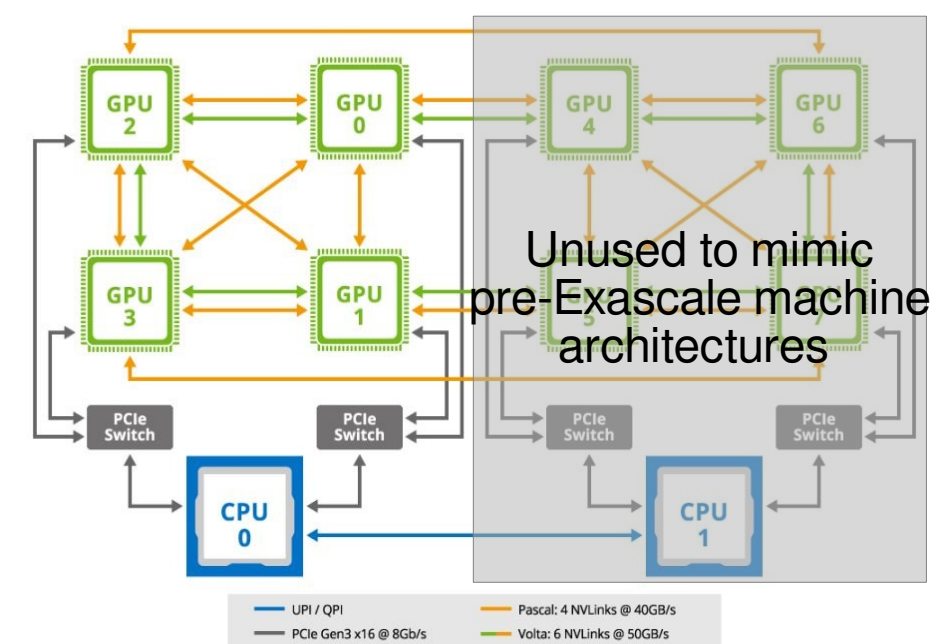- Modest improvements on low-end interconnects



PME electrostatics:
halo-exchange & PME MPMD communication

Legend:
- Staged GPU comm.
- Direct GPU comm.
- Direct GPU comm. & GPU Update
- & PME decomp

System: STMV 1M atoms

Hardware: DGX-1V

Unused to mimic pre-Exascale machine architectures

UPI / QPI
PCIe Gen3 x16 @ 8Gb/s
Pascal: 4 NVLinks @ 40GB/s
Volta: 6 NVLinks @ 50GB/s

# PME scaling improvements with cuFFTmp (in development)

STMV on NVIDIA Selene



- Strong scales reasonably well to 16-24 nodes:
  - only 10-15k (!) atoms per GPU
  - further improvements planned

- Peak can still be lower than CPU-only machines
  - algorithm improvements needed
  - next-gen hardware expected to help

# Asynchronous scheduling: CUDA graphs

Most work done by

**Alan Gray (NVIDIA)**

**Capture work in inner loop to build CUDA graph**

Decom-position

Conv. form x

Non-bonded F

Conv. Red. f

Bonded F

PME

Integration, Constraints

3990.1 ms    3990.15 ms    3990.2 ms    3990.25 ms    3990.4 ms

**Challenge:**
lacking priority support
PME kernels not prioritized

graph_376

0 _Z28pme_spline_and_spread_kernelILi4ELb1ELb1ELb1ELb1ELi1ELb0EL14ThreadsPerAtom0EEv22PmeGpuCudaKernelParams

4 _Z15regular_fft_r2cILj100ELj10ELj16EL9padding_t1EL9twiddle_t0EL20loadstore_modifier_t2EL8layout_t0EjfEv18kernel_arguments_tIT6_E

5 _Z11regular_fftILj100ELj10ELj16EL9padding_t1EL9twiddle_t0EL20loadstore_modifier_t2EL8layout_t1EjfEv18kernel_arguments_tIT6_E

6 _Z11regular_fftILj100ELj10ELj16EL9padding_t1EL9twiddle_t0EL20loadstore_modifier_t2EL8layout_t1EjfEv18kernel_arguments_tIT6_E
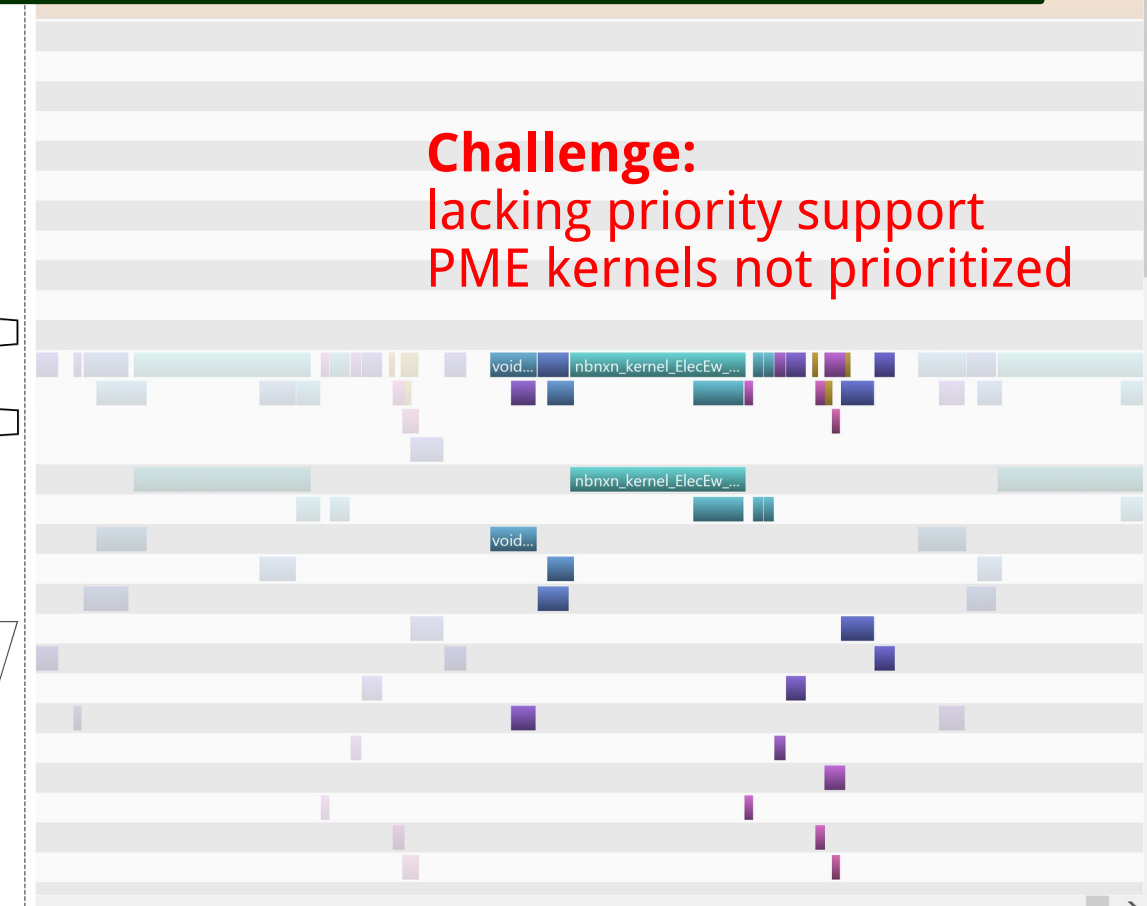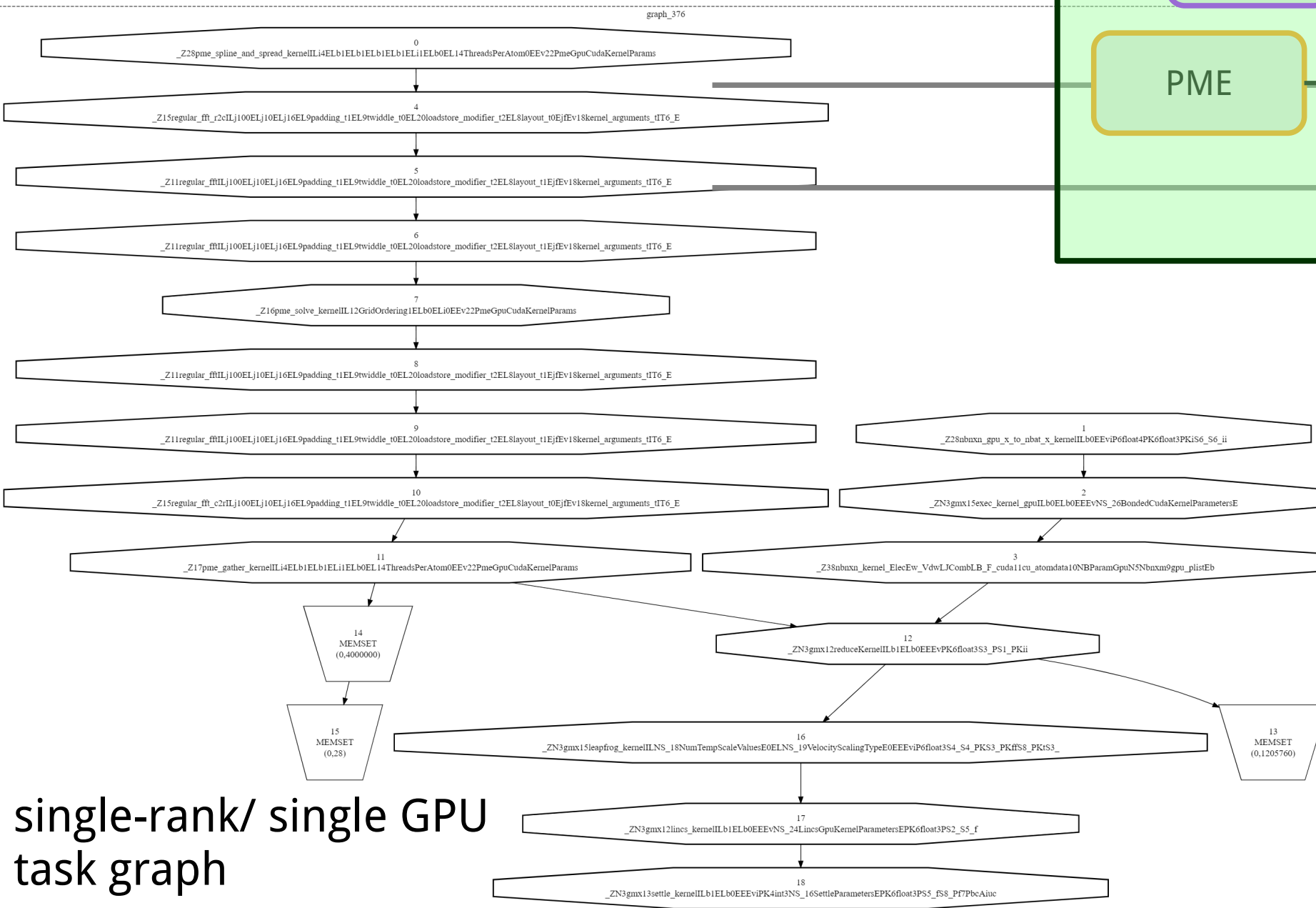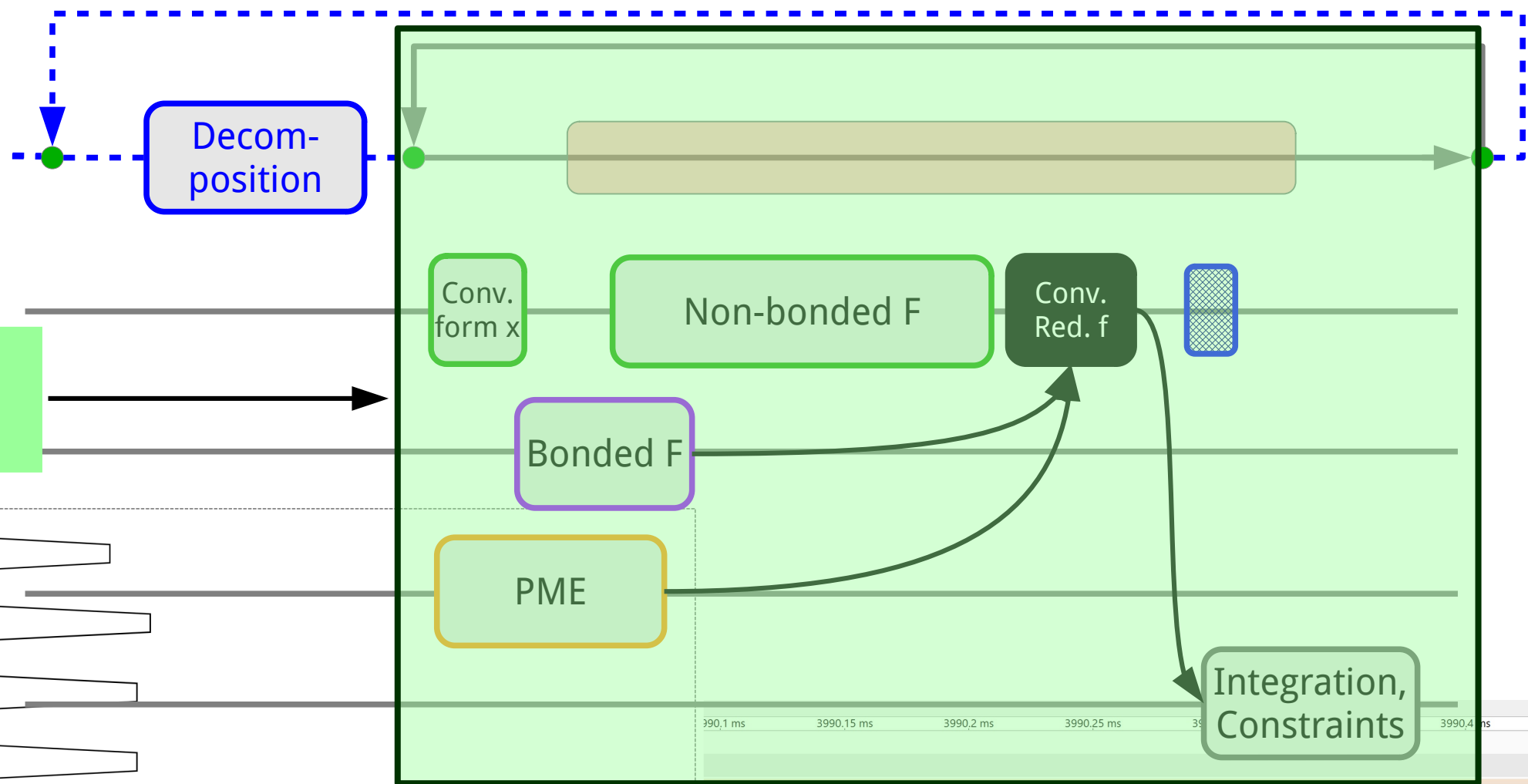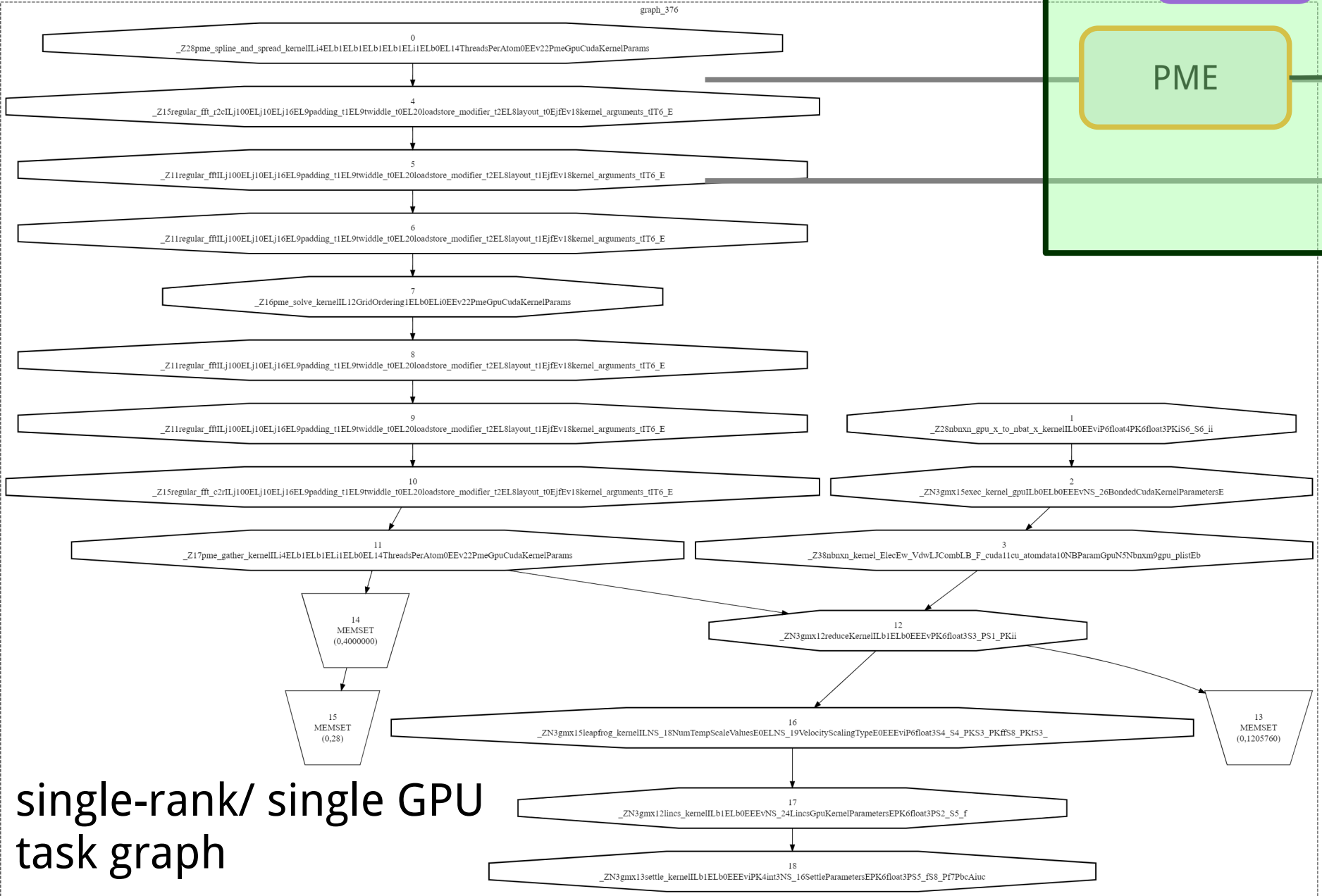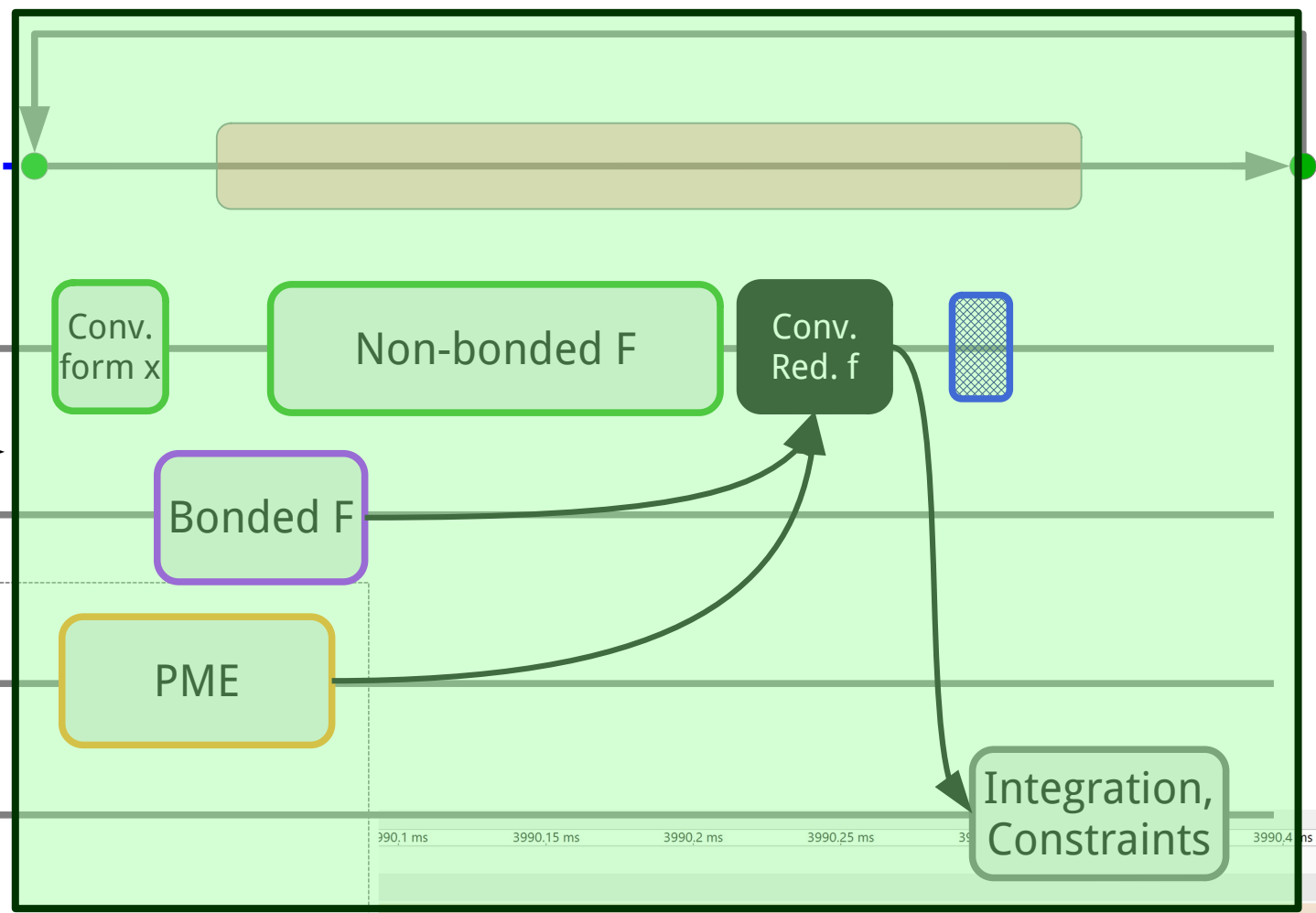
7 _Z16pme_solve_kernelIL12GridOrdering1ELb0ELi0EEv22PmeGpuCudaKernelParams

8 _Z11regular_fftILj100ELj10ELj16EL9padding_t1EL9twiddle_t0EL20loadstore_modifier_t2EL8layout_t1EjfEv18kernel_arguments_tIT6_E

9 _Z11regular_fftILj100ELj10ELj16EL9padding_t1EL9twiddle_t0EL20loadstore_modifier_t2EL8layout_t1EjfEv18kernel_arguments_tIT6_E

10 _Z15regular_fft_c2rILj100ELj10ELj16EL9padding_t1EL9twiddle_t0EL20loadstore_modifier_t2EL8layout_t0EjfEv18kernel_arguments_tIT6_E

11 _Z17pme_gather_kernelILi4ELb1ELb1ELb1ELb0EL14ThreadsPerAtom0EEv22PmeGpuCudaKernelParams

1 _Z28nbnxn_gpu_x_to_nbat_x_kernelILb0EEviP6float4PK6float3PKiS6_S6_ii

2 _ZN3gmx15exec_kernel_gpuILb0ELb0EEvNS_26BondedCudaKernelParametersE

3 _Z38nbnxn_kernel_ElecEw_VdwLJCombLB_F_cuda11cu_atomdata10NBParamGpuN5Nbnxm9gpu_plistEb

14 MEMSET (0,4000000)

12 _ZN3gmx12reduceKernelILb1ELb0EEEvPK6float3S3_PS1_PKii

15 MEMSET (0,28)

16 _ZN3gmx15leapfrog_kernelILNS_18NumTempScaleValuesE0ELNS_19VelocityScalingTypeE0EEEviP6float3S4_S4_PKS3_PKffS8_PKtS3_

13 MEMSET (0,1205760)

17 _ZN3gmx12lincs_kernelILb1ELb0EEEvNS_24LincsGpuKernelParametersEPK6float3PS2_S5_f

18 _ZN3gmx13settle_kernelILb1ELb0EEEviPK4int3NS_16SettleParametersEPK6float3PS5_fS8_Pf7PbcAiuc

single-rank/ single GPU
task graph

void...    nbnxn_kernel_ElecEw_...

nbnxn_kernel_ElecEw_...

void...

# Asynchronous scheduling: CUDA graphs



single-rank/ single GPU
task graph

# Multi-GPU graph scheduling

Pair-search & domain-decomposition step

remote MPI rank

**Graph launch for 10s-100s of iterations**

remote MPI rank

MPI comm: receive NLoc x

MPI comm: send NLoc F

DD comm

CPU

DD

Loc PS

Nloc PS

**H2D** Loc pair list

**H2D** NLoc pair list

D2H x

GPU

Local stream

List pruning

• • • preempted by non-local

Local non-bonded F

Rolling prune

Clear buffers

Reduce F

Non-local stream (high priority)

un-pack

Bonded F

Non-Local non-bonded F

pack

Spread

3D-FFT fwd

Solve

3D-FFT back

Gather

Integration Constraints

# Asynchronous scheduling: CUDA graphs

- multi-rank: leverage thread-MPI using pthreads for UVA direct async copies

Generated with
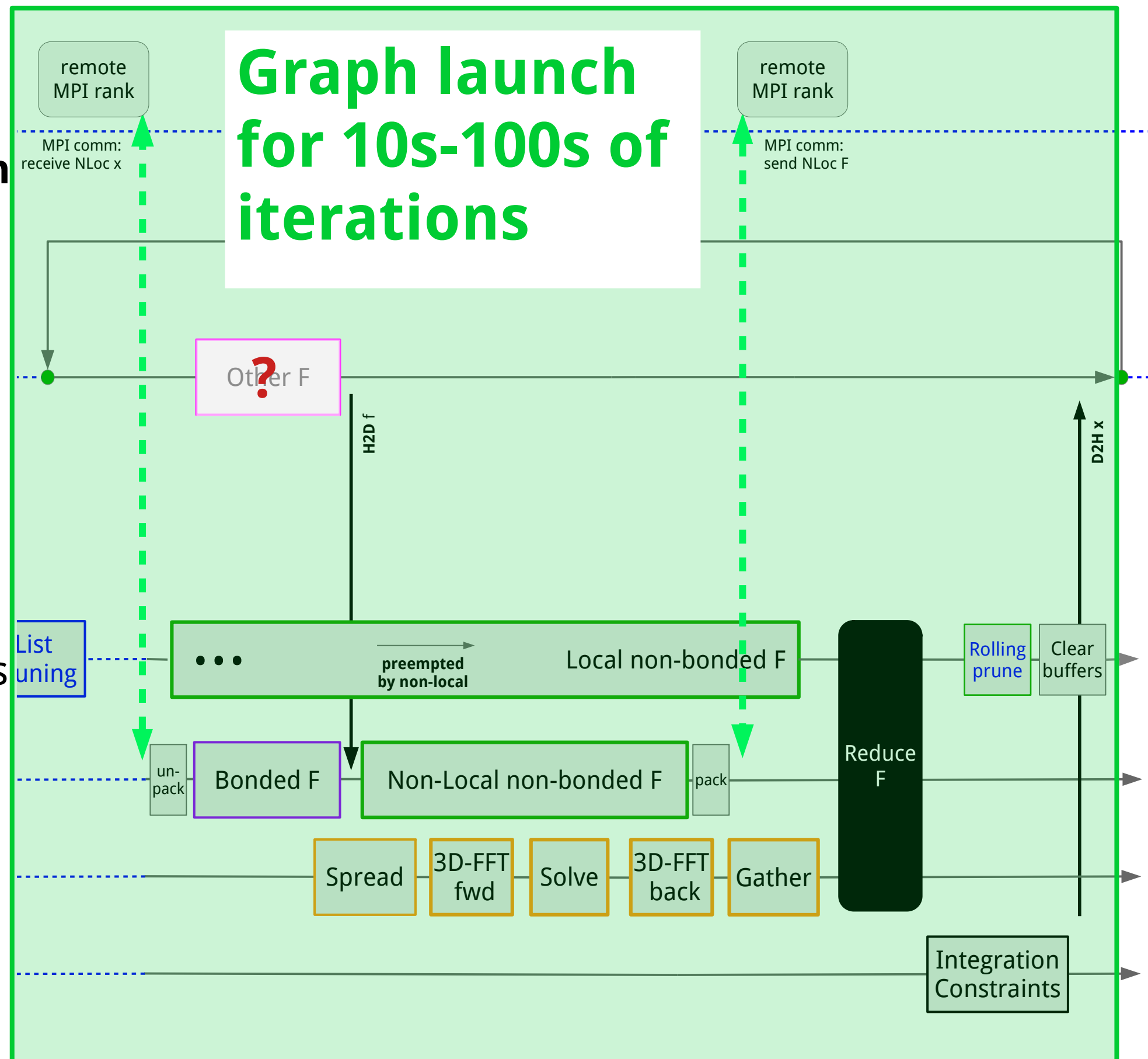**cudaGraphDebugDotPrint()**

# Multi-GPU graph scheduling

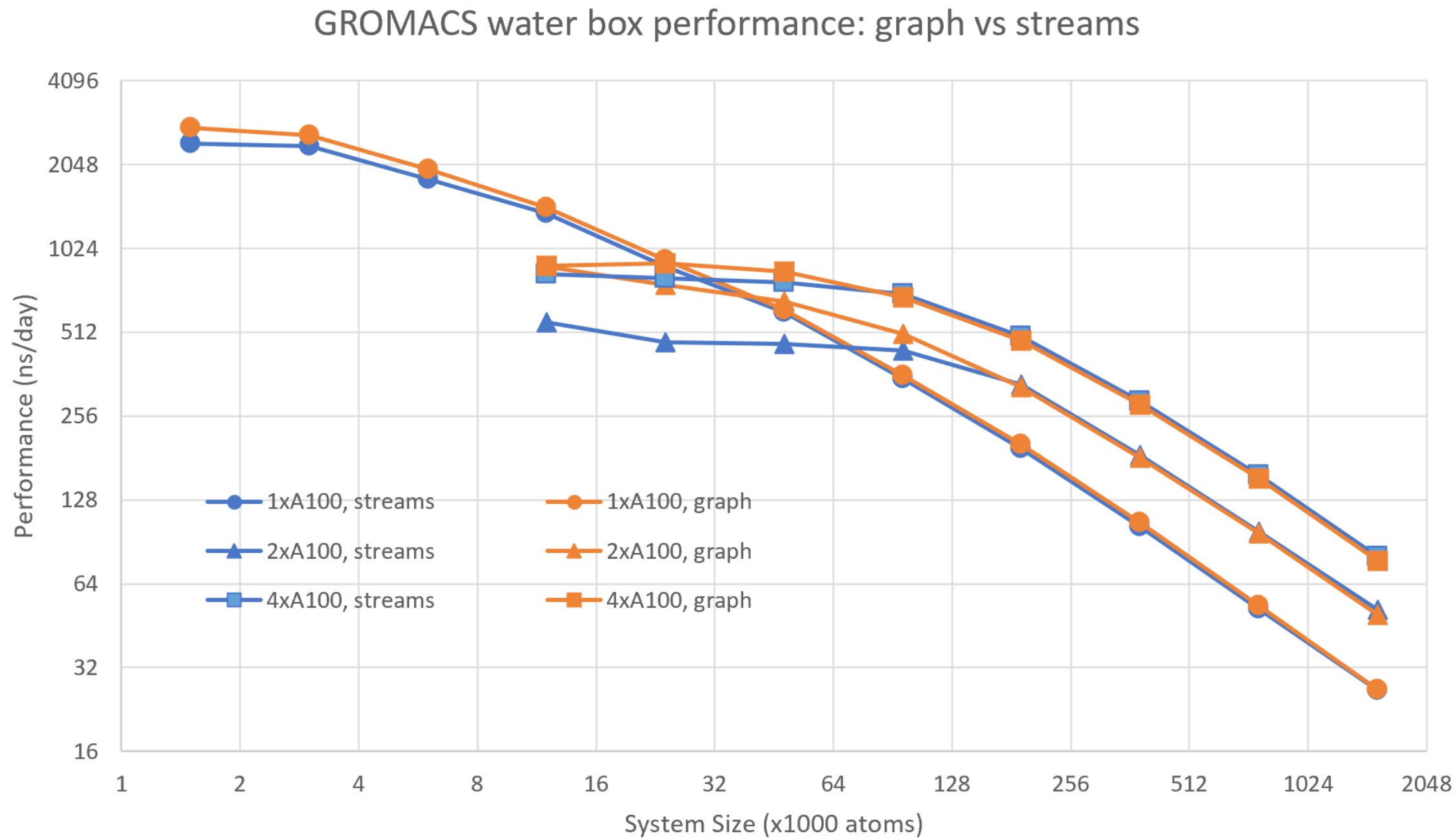- Inner loop compute + all intra-node communication **with a single graph launch**

  – Reduced overhead

  – Allow the runtime to optimize schedule

- Challenges / WIP:

  – integrating CPU tasks

  – heterogeneous iterations

  – inter-node communication

**Graph launch for 10s-100s of iterations**

remote MPI rank

MPI comm: receive NLoc x

remote MPI rank

MPI comm: send NLoc F

Other F

H2D f

D2H x

List tuning

preempted by non-local

Local non-bonded F

Rolling prune

Clear buffers

un-pack

Bonded F

Non-Local non-bonded F

pack

Reduce F

Spread

3D-FFT fwd

Solve

3D-FFT back

Gather

Integration Constraints

# CUDA graph scheduling performance



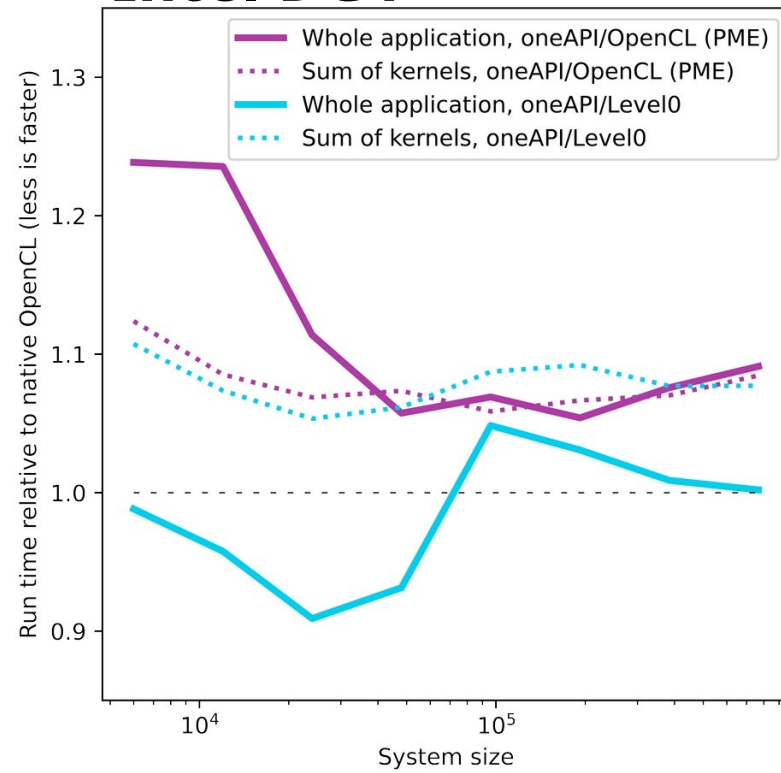GROMACS water box performance: graph vs streams
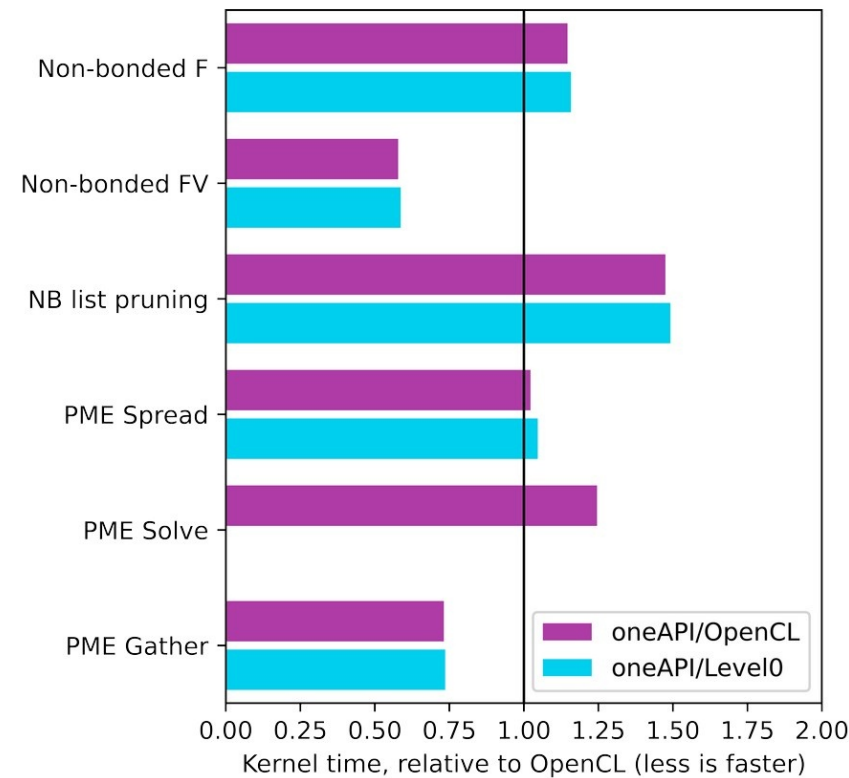
# SYCL for portability & performance

- SeRC & Intel: OneAPI CoE: Andrey Alekseenko (KTH)

- **1st GPU backend with DPC++**

  - early prototype released in GROMACS 2021

- added hipSYCL support as portability check first

- starting with the 2022 release:

  GROMACS adopted **hipSYCL for production AMD support**

- SYCL to replace OpenCL as portability GPU backend

  - already broader feature set coverage

  - broad vendor support: AMD, NVIDIA, Intel

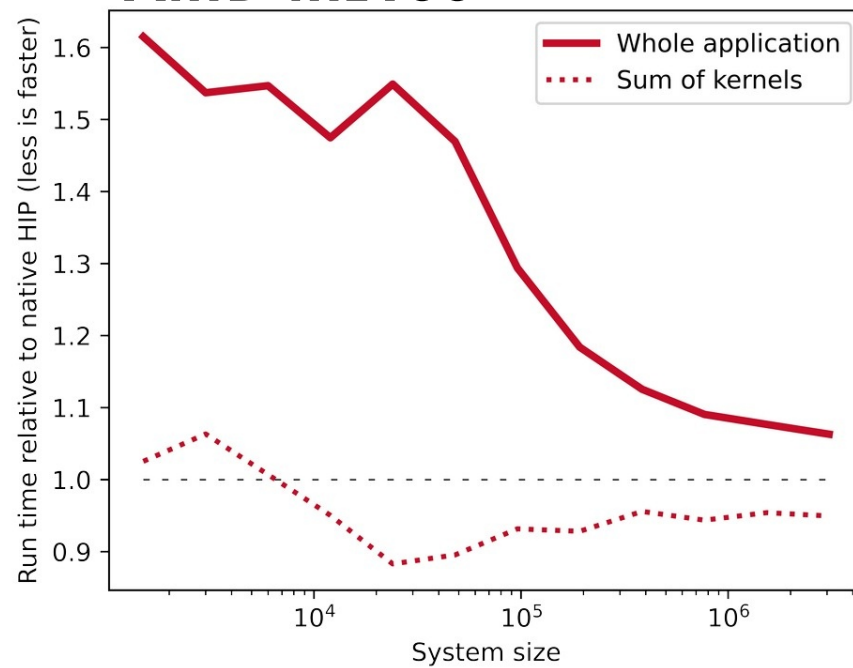# SYCL in GROMACS relative to native

**Intel DG1**



**AMD MI100**



oneAPI/DPC++ on OpenCL/L0 vs native OpenCL
(oneAPI 2022.0 except L0 2021.4)

hipSYCL vs native HIP on ROCm 4.5.2

Relative application perf and rel GPU kernels perf

Perf of individual GPU kernels

# Take-aways

- Codesign key for algorithm reformulation / redesign
  - enabled to keep up with hardware evolution
  - need to be forward-looking
  - disruptive vs constructive
- Long-term investment
- Interdisciplinarity helps but challenges too
- Collaboration needs long-term alignment
  - much easier intra-team/community
  - harder and often challenging cross-team
- Plan for the progress but allow for the incidental collaboration
  - accommodating SW design will allow new domain-science contribution

# Acknowledgments

## GROMACS

Andrey Alekseenko

Artem Zhmurov

Berk Hess

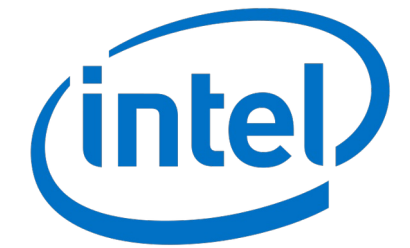Erik Lindahl

Magnus Lundborg

Paul Bauer

Mark Abraham (Intel)

Roland Schulz (Intel)

Alan Gray (NVIDIA)

Gaurav Garg (NVIDIA)

## HW / code contrib



## Funding