Heterogeneous parallelization in **GROMACS**

lessons learned from a accelerating molecular dynamics on CPUs & GPUs

Szilárd Páll

pszilard@kth.se

June 9, 2020







Acknowledgments

GROMACS

Berk Hess Artem Zhmurov

Erik Lindahl

Paul Bauer

Mark Abraham

Magnus Lundborg

Roland Schulz

Aleksei Yupinov

Alan Gray (NVIDIA) Gaurav Garg (NVIDIA)







CSCS Swiss National Supercomputing Centre

Funding





Swedish Foundation for STRATEGIC RESEARCH





European Research Council

Molecular dynamics

Newton's equation of motion

$$m_i \frac{d^2 \mathbf{r}_i}{dt^2} = \mathbf{F}_i = -\nabla_i U(\mathbf{r}),$$

• Integrate e.g. leap-frog

$$v_i(t + \frac{\Delta t}{2}) = v_i(t - \frac{\Delta t}{2})$$

$$r_i(t + \Delta t) = r_i(t) + v_i(t)$$

- Fundamental challenge:
 - sequential problem
 - **time-step limited** to 1-5 fs (by fastest atomic motions)



 $) + a_i(t)\Delta t$

 $(+\frac{\Delta t}{2})\Delta t$



Main cost: computing forces



Shared under CC BY 4.0: doi.org/10.6084/m9.figshare.12456314



Non-bonded

Over all atom-pairs!

Timescale challenge



Simulations:

- high spatial/temporal detail •
- sampling •
- model quality? •



Experiments:

- lower detail •
- higher efficiency
- •

high degree of averaging

Molecular simulation: use-cases

Biomolecular MD



Membrane protein: 10⁵ particles





DNA base-pair opening: 10⁴ particles

Shared under CC BY 4.0: doi.org/10.6084/m9.figshare





10¹⁰-10¹² particles



Materials MD

Contact line friction & wetting dynamics 10⁷-10⁹ particles



Nucleation in nano-crystals:



Molecular simulation: use-cases

Biomolecular MD

time-scale challenge 107 particles

strong scaling

→ latency sensitive often runs out of L2 cache

\rightarrow higher roofline

 \rightarrow strong benefit from high algorithm arithmetic intensity (SIMD, instruction tuning)





Shared under CC BY 4.0: doi.org/10.6084/m9.figshare

Materials MD time- & length-scale challenge strong / weak scaling → (can be) latency/BW sensitive might run out of main memory → lower roofline (less need for SIMD)



Nucleation in nano-crystals: 10¹⁰ - 10¹² particles





Molecular dynamics step



~ millisecond or less

Goal: do it as fast as possible!

MD: strong scaling challenge

Pair-search step every 50-200 iterations



~ millisecond or less

- Simluation vs real-world time-scale gap
 - Every simulation: 10⁸ 10¹⁵ steps
 - Every step: 10⁶ 10⁹ FLOPs
- MD codes at peak: ~100 µs / step (single node)
 - <100 atoms/core at peak (running in L1)</p>
 - <10000 atoms / GPU

Classical MD code

- supports all major force-fields
- broad algorithm support
- Development:

Stockholm Sweden & partners worldwide

• Large user base:

- 10k's academic & industry
- deployed on most HPC resources
- Open source: LGPLv2
- Open development:
 - code review & bug-tracker:https://gitlab.com/gromacs

GROMACS FAST. FLEXIBLE. FREE.



units cells





Eighth shell domain decomposition







virtual interaction sites



Triclinic unit cell with load balancing and staggered cell boundaries



Focus on high performance:

efficient algorithms & highly-tuned parallel code

- Bottom-up performance oriented design:
 - absolute performance over "just scaling"

Portability

- broad CI testing, Linux distro integration
- regular testing on all HPC arch
- **Code-base:** C++17, >1M LOC





arbitrary units cells

Eighth shell domain decomposition





virtual interaction sites



Triclinic unit cell with load balancing and staggered cell boundaries

GROMACS parallelization

Parallelism exploited on **multiple levels**: SIMD / threading / NUMA / async offload / MPI

- Hierarchical parallelization: target each level of hw parallelism
 - MPI: SPMD / MPMD; thread-MPI
 - OpenMP
 - SIMD: 14 flavors (SIMD library abstraction)
 - CUDA, OpenCL







Concurrency in the MD step



Parallelization of the MD step



Non-bonded interactions

- Lennard-Jones:
 - decays fast, use cutoff
- Coulomb
 - decays slowly, need long range
- Short-range:

calculate all interactions within a spherical cutoff

Long range:

PME grid-based solver using FFT



Pair interactions



Pair interactions: Verlet algorithm



Pair interactions: Verlet algorithm



Pair force calculation: SIMD-parallel traditional algorithm



Neighbor list = particle pair list: j-atoms in range = neighbors i-atom

- irregular data: non-contiguous layout •
- need to shuffle particles in registers
- low data reuse:

for each loadeded i-atom compute only once



Traditional algorithm: cache pressure, ill data reuse, register shuffle bound

Bad for SIMD!

Calculating zeros: serial vs SIMD

serial algoritm: no 0s



SIMD-style algoritm: **lots of 0s**



Source of zeros:

- need many squares/cubes to cover a sphere —
- extra particles considered: cells crossed by the cut-off sphere —

A better pair-interaction algorithm design objectves

- Need:
 - **Data reuse:** inner loop should calculate multiple interactions without having to load —
 - Calculate **as few 0s** as possible
 - Fixed pair list update interval
- Nice to have: ●
 - future-proof (can we predict the future?)
 - one algorithm, many architectures: narrow & wide SIMD, accelerators, FPGAs,...
- Keep in mind:
 - Flops are "frenered under CC BY 4.0: doi.org/10.6084/m9.fig
 - Data movement is expensive

Regularize the problem: cluster algorithm

grid cell **uniform in particle count** rather than dimensions!



cluster pair-list

Cluster algorithm



Cluster algorithm



- Regularization:
 - optimizes data layout efficient access
 - increases data **reuse**
- Flexible & adaptable:

SIMD width & reuse factor algorithmic parameters

Pair force calculation: SIMD-parallel cluster algorithm



Cluster sizes are the "knobs" to adjust for a specific arch:

- j-size: adjust to SIMD width
- i-reuse: data reuse & cache pressure

=> allows restoring the algorithm's arithmetic intensity

Parallel work-efficinency



- always >= 1 for a SIMD-optmimzed implementation
- benefits outweigh the cost ullet



pair thrpughput (pairs/kcycle)

Cluster algorithm: implicit interaction buffer

Extra particles "sticking out" of cut-off: implicit interaction buffer



Explicit interaction buffer: can.be reduced.

Parallel work-efficinency revisited



• Making use of the implicit buffer

gives a parallel efficiency gain:

=> **shorter explicit buffer** with the cluster

Example:

4x8 clusters allow

0.1 nm instead of 0.2 nm

(standard water-box with PME and 0.9 nm cutoff 40 steps list rebuild)

GPU super-cluster setup: avoid large tiling



Super-cluster setup for GPUs



Wider execution: use **hierarchical pair** list to avoid work-efficiency cliff

- joint list for multiple i-clusters
- swap i-j loop order



Cluster algorithm adaptability to hardware

- Flavors:
 - NVIDIA 32-wide 8x4
 - AMD 64-wide 8x8
 - Intel 8-wide 4x2
- Cluster-setup:

search 4x4, pruned to 4x2 for force

(parallel work efficiency advantage)



Pair interaction kernel throughput



CPUs insensitive to input size to 100s atoms/core cache effects at large inputs

Benchmark "showoff" regime:

This is where the **"free lunch" from new hardware** comes in full effect

Intra-GPU load balance



Workload per SMX: Tesla K20c, 1500 atoms

- GPU scheduling black-box
- Reshape irregular workload:
 - sort and split lists
 - improve load balance —
 - avoid kernel tail effects

Heterogeneous Parallelization

Heterogeneous offload design

- Heterogeneous vs homegeneous GPU offload
 - full port to multiople toolkits/APIs not an option for a large codebase
- Heterogeneity: both CPU and GPU
 - maintain the **versatility and feature** set
 - performance:
 - use the fastest compute unit
 - allow flexibility to balance load on various hardware
- Challenges:
 - short time/step
 - fast CPU code

se

Force offload



GPU:

- compute forces where most acceleration is to be had
- CPU:
- pair search
- other F / special algorithms
- integration

Force offload schemes



- Offloading different force components allows adjusting to hardware balance
- Seach / DD:
 - complex code
 - kept on the CPU
 - \rightarrow use algorithmic optimization to improve CPU—GPU overlap

Dual pair list

- Introduce two buffers
 - outer & inner
- Periodically re-prune
 outer → inner
- Build:

•

- outer list less frequently
- inner list more frequently
- Reduces search space to obtaining the inner list used for force compute



Heterogeneous scheme load balancing



Balance DD+ search & pruning cost

Force offload schemes



Integration on the CPU =>

CPU – GPU data movement needed

Amdahl: as GPUs get faster, **CPU integration time** increases

- Solutions:
 - use force decomp & pipeline update (PCIe bottleneck!)
 - offload integration

Multi-node force offload



Full step offload



- Trade GPU idling for CPU idling: ideal for GPU dense architectures
- CPU supporting role ("back-offload"):
 - non-offloaded per-step algorithms
 - infrequent tasks (search, DD)
- Major benefits with direct communication

Heterogeneous offload performance



- Hardware: AMD R 3900X
- Inputs:

 - 144k atoms
 - (AMBER FF)



NVIDIA 2080 SUPER

- RNAse: 24k atoms - GluCL ion channel:

Ensemble performance: small system



• Input: RNAse (24k atom AMBER99sb), uncoupled ensemble run

GPU direct communication

- CUDA-aware MPI: requires CPU sync
- thread-MPI: exchange CUDA events, fully offload communication
- NVIDIA co-design



A100 is coming

- Promising hardware and CUDA sw features
- Impressive performane at start: ~1.5x



*Disclaimer: the performance was measured on pre-production hardware

NVIDIA V100 vs A100*

Summary

- Heterogeneity is here to stay
- Strong scaling is hard, fast code is hard to accelerate
- Start with the algorithms
- Target each (most) levels of parallelism
- Open standards-based tools vs propritery tech

