

# Parallelization of molecular dynamics

## 2019年7月11日 Jaewoon Jung (RIKEN Center for Computational Science)

## **Overview of MD**

## **Molecular Dynamics (MD)**

- 1. Energy/forces are described by classical molecular mechanics force field.
- 2. Update state according to equations of motion



Equation of motion

Integration



Long time MD trajectory => Ensemble generation

Long time MD trajectories are important to obtain thermodynamic quantities of target systems.

## **Potential energy in MD**



## **Non-bonded interaction**

1. Non-bond energy calculation is reduced by introducing cutoff

$$\sum_{j=1}^{N-1} \sum_{i=j+1}^{N} \left\{ \varepsilon_{ij} \left[ \left( \frac{r_{0ij}}{r_{ij}} \right)^{12} - 2 \left( \frac{r_{0ij}}{r_{ij}} \right)^{6} \right] + \frac{q_i q_j}{4\pi \varepsilon r_{ij}} \right\} \xrightarrow{O(N^2)} \left\{ \sum_{|i-j| < R}^{N} \left\{ \varepsilon_{ij} \left[ \left( \frac{r_{0ij}}{r_{ij}} \right)^{12} - 2 \left( \frac{r_{0ij}}{r_{ij}} \right)^{6} \right] \right\} + U_{elec} \xrightarrow{O(N^1)} \left\{ \varepsilon_{ij} \left[ \left( \frac{r_{0ij}}{r_{ij}} \right)^{12} - 2 \left( \frac{r_{0ij}}{r_{ij}} \right)^{6} \right] \right\} + U_{elec} \xrightarrow{O(N^1)} \left\{ \varepsilon_{ij} \left[ \left( \frac{r_{0ij}}{r_{ij}} \right)^{12} - 2 \left( \frac{r_{0ij}}{r_{ij}} \right)^{6} \right] \right\} + U_{elec} \xrightarrow{O(N^1)} \left\{ \varepsilon_{ij} \left[ \left( \frac{r_{0ij}}{r_{ij}} \right)^{12} - 2 \left( \frac{r_{0ij}}{r_{ij}} \right)^{6} \right] \right\} + U_{elec} \xrightarrow{O(N^1)} \left\{ \varepsilon_{ij} \left[ \left( \frac{r_{0ij}}{r_{ij}} \right)^{12} - 2 \left( \frac{r_{0ij}}{r_{ij}} \right)^{6} \right] \right\} + U_{elec} \xrightarrow{O(N^2)} \left\{ \varepsilon_{ij} \left[ \left( \frac{r_{0ij}}{r_{ij}} \right)^{12} - 2 \left( \frac{r_{0ij}}{r_{ij}} \right)^{6} \right] \right\} + U_{elec} \xrightarrow{O(N^1)} \left\{ \varepsilon_{ij} \left[ \left( \frac{r_{0ij}}{r_{ij}} \right)^{12} - 2 \left( \frac{r_{0ij}}{r_{ij}} \right)^{6} \right] \right\} + U_{elec} \xrightarrow{O(N^1)} \left\{ \varepsilon_{ij} \left[ \left( \frac{r_{0ij}}{r_{ij}} \right)^{12} - 2 \left( \frac{r_{0ij}}{r_{ij}} \right)^{6} \right] \right\}$$

2. The electrostatic energy calculation beyond cutoff will be done in the reciprocal space with FFT

$$U_{elec} = \sum_{\substack{|i-j| < R}} \frac{q_i q_j}{4\pi\varepsilon_0} \frac{\operatorname{erfc}(\alpha_{ij})}{r_{ij}} + \frac{2\pi}{V} \sum_{\substack{G \neq 0}} \frac{\exp(-|\mathbf{G}|^2/4\alpha^2)}{|\mathbf{G}|^2} \sum_{ij} \frac{q_i q_j}{4\pi\varepsilon_0} \cos(\mathbf{G} \cdot r_{ij}) - \sum_{i} \frac{q_i q_i}{4\pi\varepsilon_0} \frac{\alpha}{\sqrt{\pi}}$$
Real part Reciprocal part Self energy

3. Further, it could be reduced by properly distributing over parallel processes, in particular good domain decomposition scheme.

## **Difficulty to perform long time MD simulation**

- 1. One time step length ( $\Delta t$ ) is limited to 1-2 fs due to vibrations.
- 2. On the other hand, biologically meaningful events occur on the time scale of milliseconds or longer.



## How to accelerate MD simulations? => Parallelization

**X16?** 







Parallel



16 cpus

#### **Good Parallelization ;**

- 1) Small amount of
  - computation in one process
- 2) Small communicational cost

## Parallelization

## **Shared memory parallelization (OpenMP)**



- All processes share data in memory
- For efficient parallelization, processes should not access the same memory address.
- It is only available for multi-processors in a physical node

## **Distributed memory parallelization (MPI)**



- Processors do not share data in memory
- We need to send/receive data via communications
- For efficient parallelization, the amount of communication data should be minimized

## Hybrid parallelization (MPI+OpenMP)



- Combination of shared memory and distributed memory parallelization.
- It is useful for minimizing communicational cost with very large number of processors

## **SIMD** (Single instruction, multiple data)



- Same operation on multiple data points simultaneously
- Usually applicable to common tasks like adjusting graphic image or volume
- In most MD programs, SIMD becomes the one of the important topics to increase the performance

## **SIMT (Single instruction, multiple threads)**

- SIMD combined with multithreading
- SIMT execution model is usually implemented on GPUs and related with GPGPU (General Purpose computing on Graphics Processing Units)
- Currently, CUDA allows 32 threads for SIMT (warp size =32)

Grid 1			Block (1,0	))		
Block (0,0)	Block (1,0)		Thread (0,0) Thread	Thread (1,0) Thread	Thread (2,0)	Thread (3,0) Thread
Block (0,1)	Błock (1,1)	*****	(0,1) Thread (0,2)	(1,1) Thread (1,2)	(2,1) Thread (2,2)	(3,1) Thread (3,2)

## MPI Parallelization of MD (non-bonded interaction in real space)

## **Parallelization scheme 1: Replicated data approach**

- 1. Each process has a copy of all particle data.
- 2. Each process works only part of the whole works by proper assign in do loops.





## **Parallelization scheme 1: Replicated data approach**



Perfect load balance is not guaranteed in this parallelization scheme

## Hybrid (MPI+OpenMP) parallelization of the Replicated data approach

- 1. Works are distributed over MPI and OpenMP threads.
- 2. Parallelization is increased by reducing the number of MPIs involved in communications.

my_rank = MPI_Rank proc = total MPI
<pre>do i = my_rank+1, N, prod do j = i+1, N</pre>
energy(i,j) force(i,j)
end do end do
MPI reduction (energy,force)

```
my rank = MPI Rank
proc = total MPI
nthread = total OMP thread
!$omp parallel
id = omp thread id
my id = my rank*nthread + id
do i = my id+1,N,proc*nthread
 do j = i+1, N
   energy(i,j)
    force(i,j)
  end do
end do
Openmp reduciton
!$omp end parallel
```

MPI reduction (energy, force)

## **Pros and Cons of the Replicated data approach**

- 1. **Pros** : easy to implement
- 2. Cons
  - Parallel efficiency is not good
    - No perfect load balance
    - Communication cost is not reduced by increasing the number of processes
    - We can parallelize only for energy calculation (with MPI, parallelization of integration is not so much efficient)
  - Needs a lot of memory
    - Usage of global data

## **Parallelization scheme 2: Domain decomposition**



- The simulation space is divided into subdomains according to MPI (different colors for different MPIs).
- 2. Each MPI only considers the corresponding subdomain.
- 3. MPI communications only among neighboring processes.



#### **Parallelization scheme 2: Domain decomposition**



- 1. For the interaction between different subdomains, it is necessary to have the data of the buffer region of each subdomain.
- 2. The size of the buffer region is dependent on the cutoff values.
- 3. The interaction between particles in different subdomains should be considered very carefully.
- 4. The size of the subdomain and buffer region is decreased by increasing the number of processes.

## Pros and Cons of the domain decomposition approach

#### 1. Pros

- Good parallel efficiency
  - Reduced computational cost by increasing the number of processes
  - We can easily parallelize not only energy but also integration
- Availability of a huge system
  - Data size is decreased by increasing the number of processes

#### 2. Cons

- Implementation is not easy
  - Domain decomposition scheme is highly depend on the potential energy type, cutoff and so on
- Good performance cannot be obtained for nonuniform particle distributions

## Special treatment for nonuniform distributions

#### **Tree method**

: Subdomain size is adjusted to have the same number of particles



#### Hilbert space filling curve

: a map that relates multi-dimensional space to one-dimensional curve



## **Comparison of two parallelization scheme**

	Computation	Communication	Memory
Replicated data	<i>O</i> ( <i>N</i> / <i>P</i> )	O(N)	O(N)
Domain decomposition	<i>O</i> ( <i>N</i> / <i>P</i> )	$O((N/P)^{2/3})$	<i>O</i> ( <i>N</i> / <i>P</i> )

N: system size P: number of processes MPI Parallelization of MD (reciprocal space)

### **Smooth particle mesh Ewald method**



The structure factor in the reciprocal part is approximated as

$$S(k_1, k_2, k_3) = b_1(k_1)b_2(k_2)b_3(k_3)F(Q)(k_1, k_2, k_3)$$
  
Using Cardinal B-splines of order n   
**Fourier Transform of charge**

It is important to parallelize the Fast Fourier transform efficiently in PME!!

Ref : U. Essmann et al, J. Chem. Phys. 103, 8577 (1995)

**Overall procedure of reciprocal space calculation** 



## Simple note of MPI\_alltoall communications



**MPI\_alltoall is same as matrix transpose!!** 

## Parallel 3D FFT – slab (1D) decomposition



- Each process is assigned a slob of size N ×
   N × N/P for computing FFT of N × N × N
   grids on P processes.
- 2. The scalability is limited by *N*
- 3. *N* should be divisible by *P*

## Parallel 3D FFT – slab (1D) decomposition



## Parallel 3D FFT – slab (1D) decomposition (continued)

#### 1. Slab decomposition of 3D FFT has three steps

- 2D FFT (or two 1D FFT) along the two local dimension
- Global transpose (communication)
- 1D FFT along third dimension

#### 2. Pros

- fast using small number of processes
- 3. Cons
  - Limitation of the number of processes

## **Parallel 3D FFT – 2D decomposition**



Each process is assigned a data of size  $N \times N/P \times N/Q$  for computing a FFT of  $N \times N \times N$  grids on  $P \times Q$  processes.

## Parallel 3D FFT –2D decomposition (continued)

- 1. 2D decomposition of 3D FFT has five steps
  - 1D FFT along the local dimension
  - Global transpose
  - 1D FFT along the second dimension
  - Global transpose
  - 1D FFT along the third dimension
- 2. The global transpose requires communication only between subgroups of all nodes
- 3. Cons: Slower than 1D decomposition for a small number of processes
- 4. Pros : Maximum parallelization is increased

## Pseudo code of reciprocal space calculation with 2D decomposition of 3D FFT

! compute Q factor do i = 1, natom/P compute Q\_orig end do call mpi\_alltoall(Q\_orig, Q\_new, ...) accumulate Q from Q\_new

```
!FFT : F(Q)
do iz = 1, zgrid(local)
 do iy = 1, ygrid(local)
   work local = Q(my rank)
   call fft(work local)
   Q(my rank) = work local
  end do
end do
call mpi alltoall(Q, Q new,...)
do iz = 1, zgrid(local)
 do ix = 1, xgrid(local)
   work_local = Q_new(my_rank)
   call fft(work local)
   Q(my rank) = work local
  end do
end do
call mpi alltoall(Q,Q new,..)
```

do iy = 1, ygrid(local)
 do ix = 1, xgrid(local)
 work\_local = Q(my\_rank)
 call fft(work\_local)
 Q(my\_rank) = work\_local
 end do
end do

! compute energy and virial do iz = 1, zgrid do iy = 1, ygrid(local) do ix = 1, xgrid(local) energy = energy + sum(Th\*Q) virial = viral + .. end do end do end do

! X=F\_1(Th)\*F\_1(Q)

! FFT (F(X))

do iy = 1, ygrid(local)do ix = 1, xgrid(local) work local = Q(my rank)call fft(work local) Q(my rank) = work local end do end do call mpi alltoall(Q,Q new,..) do iz = 1, zgrid(local)do ix = 1, xgrid(local) work local = Q(my rank)call fft(work local) Q(my rank) = work local end do end do call mpi alltoall(Q,Q\_new) do iz = 1, zgrid(local)do iy = 1, ygrid(local)work local = Q(my rank)call fft(work local) Q(my rank) = work\_local end do end do

!compute force

## **Parallel 3D FFT – 3D decomposition**



- More communications than existing FFT
- MPI\_Alltoall communications only in one dimensional space
- Reduce communicational cost for large number of processes

Ref) J. Jung et al. Comput. Phys. Comm. 200, 57-65 (2016)

## **Case Study: Parallelization in GENESIS MD software**

## **Domain decomposition in GENESIS SPDYN**



- 1. The simulation space is divided into subdomain according to the number of MPI processes.
- 2. Each subdomain is further divided into the unit domain (named cell).
- 3. Particle data are grouped together in each cell.
- 4. Communications are considered only between neighboring processes.
- 5. In each subdomain, OpenMP parallelization is used.

## Efficient shared memory calculation =>Cell-wise particle data

- 1. Each cell contains an array with the data of the particles that reside within it
- 2. This improves the locality of the particle data for operations on individual cells or pairs of cells.







cell-wise arrays of particle data

Ref : P. Gonnet, JCC 33, 76-81 (2012)

## Midpoint cell method



Subdomain assigned by MPI (Computation is assigned to one CPU)

Cell

Boundary cell

- For every cell pair, midpoint cell is decided.
- If midpoint cell is not uniquely decided, we decide it by considering load balance.
- For every particle pairs, interaction subdomain is decided from the midpoint cell.
- With this scheme, we can only communicate data of each subdomain only considering the adjacent cells of each subdomain.

## Subdomain structure with midpoint cell method

- 1. Each MPI has local subdomain which is surrounded by boundary cells.
- 2. At least, the local subdomain has at least two cells in each direction



## **Communication pattern (example of coordinates)**

- 1. Each subdomain can evaluate interactions by obtaining the data of boundary region from other processes
- 2. Using the communication pattern shown below, we can minimize the frequency of communication.



- ① Receive from the upper process in x dimension
- 2 Receive from the lower process in x dimension
- ③ Receive from the upper process in y dimension
- A Receive from the lower process in y dimension

## **Communication pattern (example of forces)**

- 1. Communication pattern of forces is opposite to the communication pattern of coordinates.
- 2. Each process send the force data in boundary cells.



- ① Send to the lower process in y dimension
- 2 Send to the upper process in y

dimension

- ③ Send to the lower process in x dimension
- ④ Send to the upper process in x dimension

## **OpenMP parallelization in GENESIS**

- 1. In the case of integrator, every cell indices are divided according to the thread id.
- 2. As for the non-bond interaction, cell-pair lists are first identified and cell-pair lists are distributed to each thread

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



## **Parallelization of FFT in GENESIS**



## FFT in GENESIS (2 dimensional view)



## Parallelization based on CPU/GPU calculation



Computationally expensive part: GPU

Communicationally expensive part: CPU

## **SIMD in GENESIS (developer version)**

• Array of Structure (AoS)

<i>x</i> <sub>1</sub>	$\mathcal{Y}_1$	<i>z</i> <sub>1</sub>	<i>x</i> <sub>2</sub>	<i>Y</i> <sub>2</sub>	<i>z</i> <sub>2</sub>									$x_N$	$\mathcal{Y}_N$	$Z_N$
-----------------------	-----------------	-----------------------	-----------------------	-----------------------	-----------------------	--	--	--	--	--	--	--	--	-------	-----------------	-------

In GENESIS, it is expressed as coord\_pbc(1:3,1:natom,1:ncell)



• Structure of Array (SoA)

<i>x</i> <sub>1</sub>	<i>x</i> <sub>2</sub>	<i>x</i> <sub>3</sub>	<i>x</i> <sub>4</sub>											<i>Z</i> <sub><i>N</i>-2</sub>	<i>z<sub>N-1</sub></i>	$Z_N$
-----------------------	-----------------------	-----------------------	-----------------------	--	--	--	--	--	--	--	--	--	--	--------------------------------	------------------------	-------

In updated GENESIS source code for KNL, it will be expressed as coord\_pbc(1:natom, 1:3, 1:ncell)

• For Haswell/Broadwell/KNL machines, SoA shows better performance than AoS due to efficient vectorization (SIMD)

## **GENESIS** developments on KNL

#### Simple pair-list scheme with better performance



• Pair-list scheme in GENESIS 1.X does not show good performance due to noncontiguous memory access (Figure (a))

## **GENESIS** developments on KNL

	Memory usage	Number of interactions	CPU time (pair-list)	CPU time (non-bonded)	CPU time (total)
Algorithm 1 (AoS)	305.88 MiB	47364770	9.95	36.42	46.37
Algorithm 1 (SoA)	305.88 MiB	47364770	6.60	33.71	40.31
Algorithm 2 (AoS)	48.4 MiB	122404436	4.91	102.46	107.37
Algorithm 2 (SoA)	48.4 MiB	122404436	4.15	27.86	32.01

- Algorithm 1: pair-list algorithm used in GENESIS 1.X
- Algorithm 2: New pair-list scheme
- Target system : ApoA1 (about 90,000 atoms, 2000 time steps)

## **Benchmark performance of GENESIS**

## **GENESIS 1.0 performance on K (1)**

• atom size:

#### 11737298

- macro molecule size:
- metabolites size:
- ion size:
- water size:
- PBC size:

216 (43 type) 4212 (76 type) 23049 (Na<sup>+</sup>, Cl<sup>-</sup>, Mg<sup>2+</sup>, K<sup>+</sup>) 2944143 480 x 480 x 480(Å<sup>3</sup>)









## **GENESIS** performance on KNL (Oak-foreast)



## **GENESIS** performance on KNL (LANL Trinity)



## Efficient parallelization enable us to perform the MD simulations of the world largest system



## Summary

- Parallelization: MPI and OpenMP
  - MPI: Distributed memory parallelization
  - OpenMP: Shared memory parallelization
- Parallelization of MD: mainly by domain decomposition with FFT parallelization
- Key issues in parallelization: minimizing communication cost to maximize the parallel efficiency