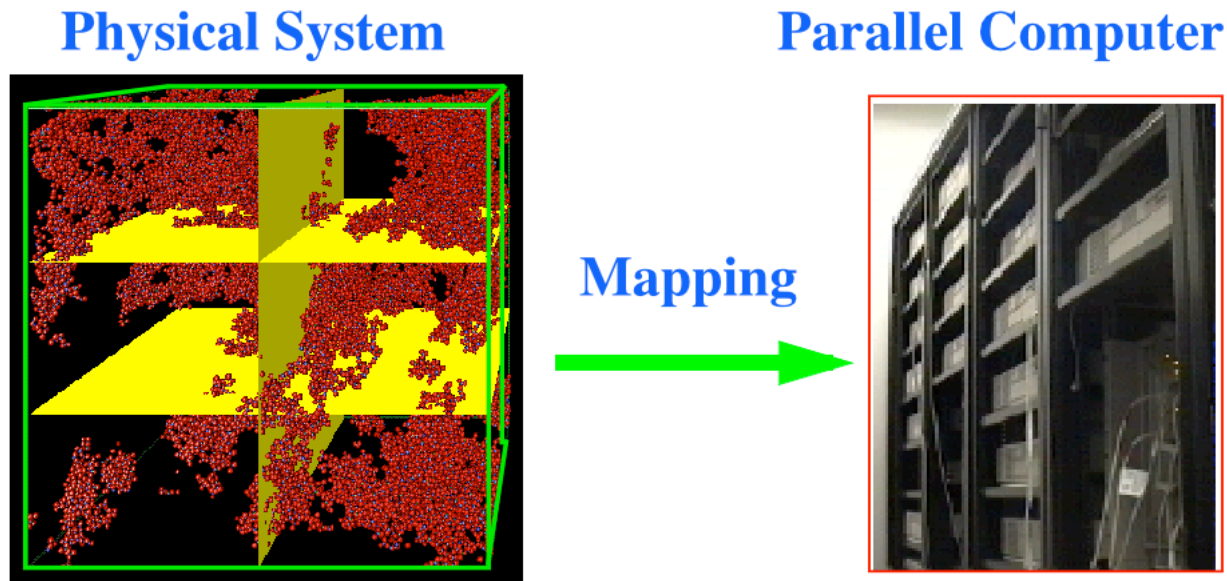


# Parallel Molecular Dynamics

This chapter explains the example parallel MD program, `pmd.c`, in detail.

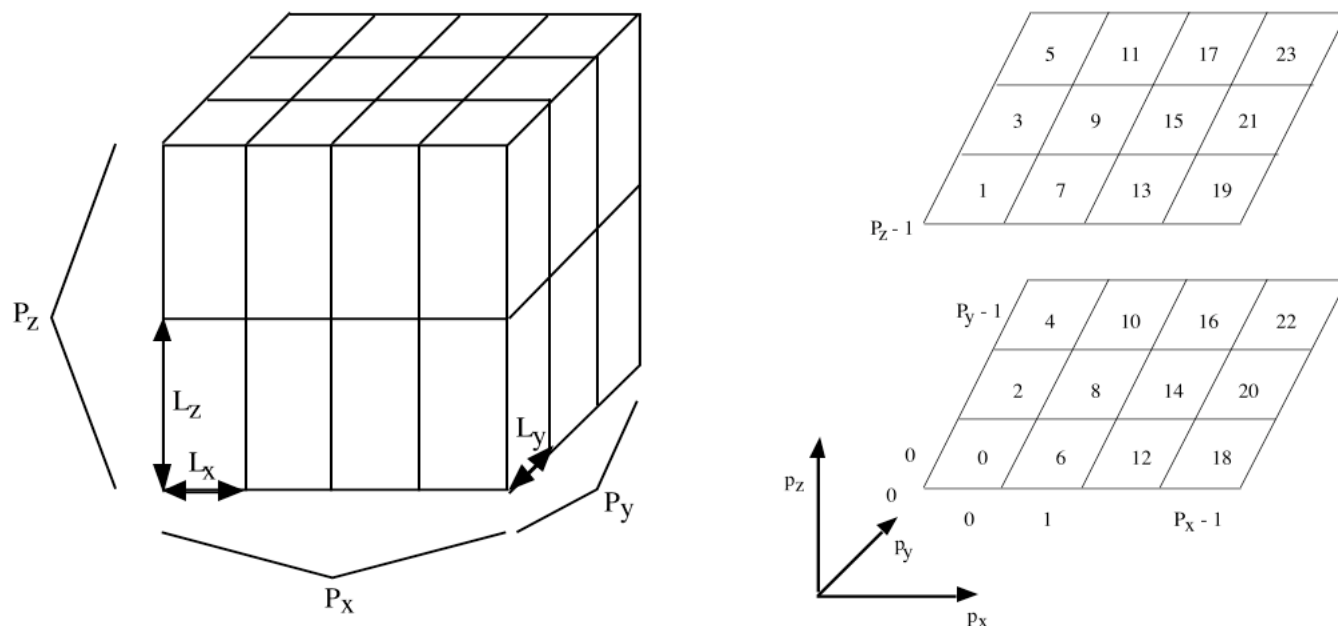
## Spatial Decomposition

- Spatial decomposition:** The physical system to be simulated is partitioned into subsystems of equal volume. Processors in a parallel computer are logically arranged according to the topology of physical subsystems. Atoms that are located in a particular subsystem are assigned to the corresponding processor.



We use a simple 3D mesh (or torus because of the periodic boundary conditions) decomposition. Subsystems are (and accordingly processors are logically) arranged in a 3D array of dimensions  $P_x \times P_y \times P_z$  (`int vproc[3]` in `pmd.c`). Each subsystem is a parallel-piped of size  $L_x \times L_y \times L_z$  (`double al[3]`).

## PROCESSOR ID



Each processor is given a unique processor ID,  $p$ , the range of which is  $[0, P-1]$  (`int sid`) where  $P = P_x P_y P_z$  is the total number of processors (`int nproc`). We also define a vector processor ID  $\vec{p} = (p_x, p_y, p_z)$  (`int vid[3]`), where  $p_x = 0, \dots, P_x-1$ ;  $p_y = 0, \dots, P_y-1$ ; and  $p_z = 0, \dots, P_z-1$ .

The relation between the sequential and vector ID's is

$$\begin{aligned} p_x &= p / (P_y P_z) \\ p_y &= (p / P_z) \bmod P_y \\ p_z &= p \bmod P_z \end{aligned}$$

or

$$p = p_x \times P_y P_z + p_y \times P_z + p_z.$$

(Example)  $(P_x, P_y, P_z) = (4, 3, 2)$

Sequential processor ID, $p$	$P_x$	Vector ID, $\vec{p}$	
		$p_x$	$p_z$
0	0	0	0
1	0	0	1
2	0	1	0
3	0	1	1
4	0	2	0
5	0	2	1
6	1	0	0
7	1	0	1
8	1	1	0
9	1	1	1
10	1	2	0
11	1	2	1
...	...	...	...

## NEIGHBOR PROCESSOR ID

For each processor, the six face-shared neighbors are identified by a sequential index ( $\kappa = 0, \dots, 5$ ). For each neighbor,  $\kappa$ , the shift-length vector  $\vec{\Delta} = (\Delta_x, \Delta_y, \Delta_z)$  denotes the position of the neighbor subsystem relative to itself (`double sv[6][3]`). The following table lists the six neighbors, where an integer vector  $\vec{\delta} = (\delta_x, \delta_y, \delta_z)$  specifies the relative location of each neighbor.

Neighbor ID, $\kappa$	$\vec{\delta} = (\delta_x, \delta_y, \delta_z)$	$\vec{\Delta} = (\Delta_x, \Delta_y, \Delta_z)$
0 (east)	(-1, 0, 0)	(- $L_x$ , 0, 0)
1 (west)	(1, 0, 0)	( $L_x$ , 0, 0)
2 (north)	(0, -1, 0)	(0, - $L_y$ , 0)
3 (south)	(0, 1, 0)	(0, $L_y$ , 0)
4 (up)	(0, 0, -1)	(0, 0, - $L_z$ )
5 (down)	(0, 0, 1)	(0, 0, $L_z$ )

The sequential processor ID,  $p'(\kappa)$  (`int nn[6]`), is obtained by

$$p'_\alpha(\kappa) = [p_\alpha + \delta_\alpha(\kappa) + P_\alpha] \bmod P_\alpha \quad (\alpha = x, y, z)$$

and

$$p'(\kappa) = p'_x(\kappa) \times P_y P_z + p'_y(\kappa) \times P_z + p'_z(\kappa)$$

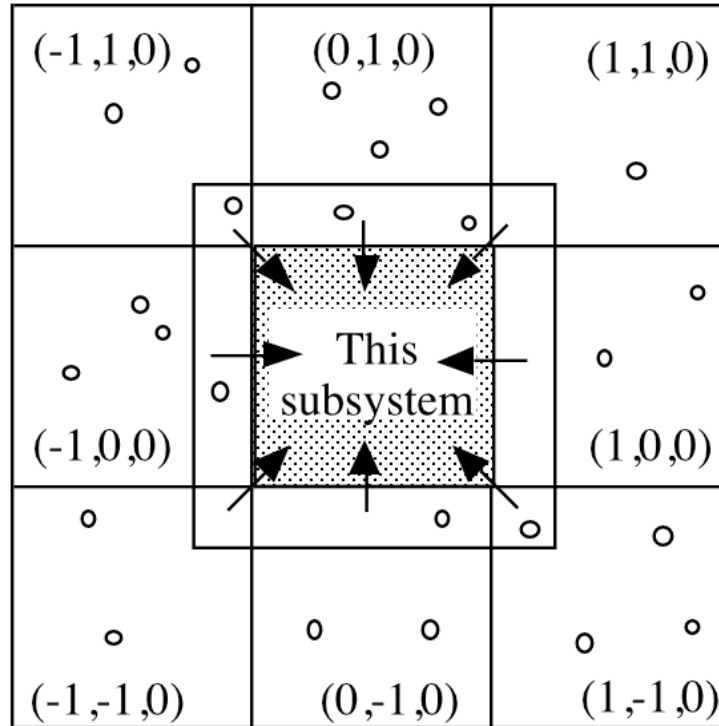
In an SPMD program using the MPI, each process knows its process ID through the returned rank value from an `MPI_Comm_rank(MPI_Comm comm, int *rank)` call. We identify the rank with `sid`.

## Parallel MD Concepts

Using the periodic boundary conditions, every subsystem has 26 neighbor subsystems, which share either a corner, an edge, or a face with it. (Or there are 6 face-sharing neighbor subsystems.) We express atomic coordinates relative to the origin of each subsystem, i.e.,  $0 < x_\alpha < L_\alpha$  ( $\alpha = x, y, z$ ). (Every subsystem thinks it is the center of the world.)

### ATOM CACHING

- **Augmented system:** In order to compute interatomic interaction with cut-off length  $r_c$  (RCUT), atomic coordinates of 26 neighbor subsystems which are located within  $r_c$  from the subsystem boundary are copied to “this” processor (**atom caching**).



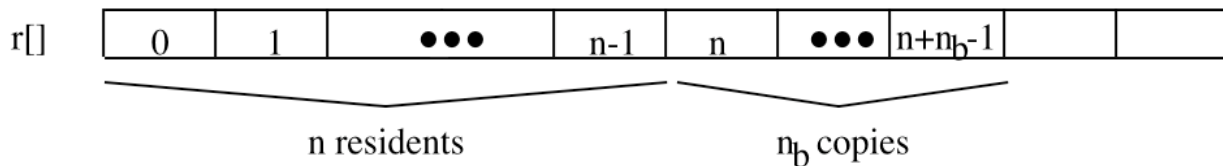
- **Cache coherence** is maintained by copying the latest neighbor surface atoms every time before atomic accelerations are computed.

### DATA STRUCTURES

$n$ : The number of resident atoms which reside in “this” subsystem.

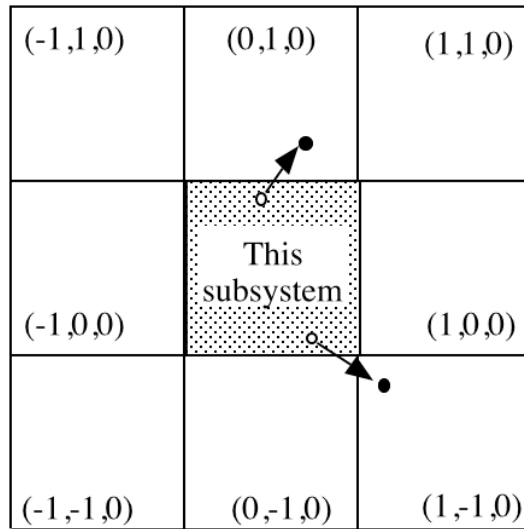
$n_b$ : The number of copied surface atoms from the neighbors.

$r[\text{NEMAX}][3]$ : Atomic coordinates relative to the subsystem origin.  $r[0:n-1][ ]$  hold the atoms reside in this subsystem, while  $r[n:n+n_b-1][ ]$  hold the cached neighbor atoms.



### ATOM MIGRATION

After the atomic coordinates are updated according to the velocity Verlet algorithm, some resident atoms may have moved out of the subsystem boundary. Such **migrated** atoms must be **moved** to proper processors.



## PARALLEL MD ALGORITHM

In the parallel MD program, `pmd.c`, a single time step of the velocity-Verlet algorithm consists of the following steps.

1. First half kick to obtain  $v_i(t+Dt/2)$
2. Update atomic coordinates to obtain  $r_i(t+Dt)$
3. `atom_move()`: Migrate the moved-out atoms to the neighbor processors
4. `atom_copy()`: Copy the surface atoms within distance  $r_c$  from the neighbors
5. `compute_accel()`: Compute new accelerations,  $a_i(t+Dt)$ , including the contributions from the cached atoms
6. Second half kick to obtain  $v_i(t+Dt)$

## Linked-List Cell Method for Interatomic Force Computation

### **function compute\_accel()**

Recall that a naive double-loop implementation to compute pair interactions scales as  $O(N^2)$ . In fact, with a finite cut-off length,  $r_c$ , an atom interacts with only a limited number of atoms  $\sim (4\pi/3) r_c^3 (N/V)$ . The linked-list cell algorithm computes the entire interaction with  $O(N)$  operations.

## CELLS

First divide the system consisting of the resident and cached atoms into small cells of equal size. The edge lengths of each cell,  $(r_{cx}, r_{cy}, r_{cz})$  (`double rc[3]`), must be at least  $r_c$ ; we use  $r_{c\alpha} = L_\alpha / L_{c\alpha}$ , where  $L_{c\alpha} = \lfloor L_\alpha / r_{c\alpha} \rfloor$ . An atom in a cell interacts with only other atoms in the same cell and its 26 neighbor cells. Including the cached atoms, the range of atomic coordinates is  $[-r_c, L_\alpha + r_c]$  ( $\alpha = x, y, z$ ). The number of cells to accommodate all these atoms is  $(L_{cx}+2) \times (L_{cy}+2) \times (L_{cz}+2)$ , where  $L_{c\alpha} = L_\alpha / r_{c\alpha}$  ( $\alpha = x, y, z$ ) (`int 1c[3]`). We identify a cell with a vector cell index,  $\vec{c} = (c_x, c_y, c_z)$  ( $0 \leq c_x \leq L_{cx}+1$ ;  $0 \leq c_y \leq L_{cy}+1$ ;  $0 \leq c_z \leq L_{cz}+1$ ), and a sequential cell index,

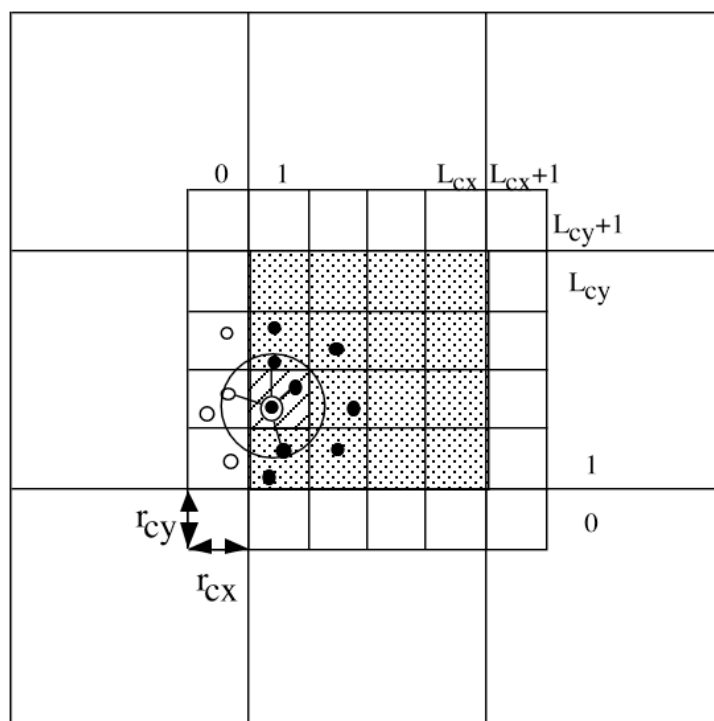
$$c = c_x(L_{cy}+2)(L_{cz}+2) + c_y(L_{cz}+2) + c_z$$

or

$$\begin{aligned} c_x &= c / [(L_{cy}+2)(L_{cz}+2)] \\ c_y &= [c / (L_{cz}+2)] \bmod (L_{cy}+2) \\ c_z &= c \bmod (L_{cz}+2). \end{aligned}$$

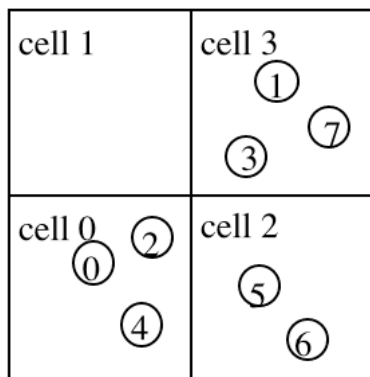
An atom with coordinate  $\vec{r}$  belongs to a cell with the vector cell index,

$$c_\alpha = \lfloor (r_\alpha + r_{c\alpha}) / r_{c\alpha} \rfloor \quad (\alpha = x, y, z).$$



## LISTS

The atoms belonging to a cell is organized as a linked list.



	0	1	2	3					
head	4	E	6	7					
	0	1	2	3	4	5	6	7	
lsc1	E	E	0	1	2	E	5	3	



## DATA STRUCTURES

$lsc1[NEMAX]$ : An array implementation of the linked list.  $lsc1[i]$  holds the atom index to which the  $i$ -th atom points.

$head[NCLMAX]$ :  $head[c]$  holds the index of the first atom in the  $c$ -th cell, or  $head[c] = EMPTY (= -1)$  if there is no atom in the cell.

## ALGORITHM 1: LIST CONSTRUCTOR

```
/* Reset the headers, head */
for (c=0; c<lcxyz2; c++) head[c] = EMPTY;
/* Scan atoms to construct headers, head, & linked lists, lscl */
for (i=0; i<n+nb; i++) {
  /* Vector cell index to which this atom belongs */
  for (a=0; a<3; a++) mc[a] = (r[i][a]+rc[a])/rc[a];
  /* Translate the vector cell index, mc, to a scalar cell index */
  c = mc[0]*lcyz2+mc[1]*lc2[2]+mc[2];
  /* Link to the previous occupant (or EMPTY if you're the 1st) */
  lscl[i] = head[c];
  /* The last one goes to the header */
  head[c] = i;
}
lcyz2 = lc2[1]*lc2[2] where lc2[a] = lc[a]+2 (a = 0,1,2); lcxyz2 = lcyz2*lc2[0]
```

## ALGORITHM 2: INTERACTION COMPUTATION

```
/* Scan inner cells */
for (mc[0]=1; mc[0]<=lc[0]; (mc[0])++)
for (mc[1]=1; mc[1]<=lc[1]; (mc[1])++)
for (mc[2]=1; mc[2]<=lc[2]; (mc[2])++) {
  /* Calculate a scalar cell index */
  c = mc[0]*lcyz2+mc[1]*lc2[2]+mc[2];
  /* Scan the neighbor cells (including itself) of cell c */
  for (mc1[0]=mc[0]-1; mc1[0]<=mc[0]+1; (mc1[0])++)
  for (mc1[1]=mc[1]-1; mc1[1]<=mc[1]+1; (mc1[1])++)
  for (mc1[2]=mc[2]-1; mc1[2]<=mc[2]+1; (mc1[2])++) {
    /* Calculate the scalar cell index of the neighbor cell */
    c1 = mc1[0]*lcyz2+mc1[1]*lc2[2]+mc1[2];
    /* Scan atom i in cell c */
    i = head[c];
    while (i != EMPTY) {
      /* Scan atom j in cell c1 */
      j = head[c1];
      while (j != EMPTY) {
        ...
        if (i < j && rij < rc2) Process pair (i, j)
        ...
        j = lscl[j];
      }
      i = lscl[i];
    }
  }
}
}
```

Note that the central cells,  $c$ , and hence the central atoms,  $i$ , are residents, whereas the pair atoms,  $j$ , are either residents or cached. The parallelization strategy in `pmd.c` is that each MPI process is responsible for computing the forces on its resident atoms. Accordingly, the force contribution,  $\vec{a}_{ij}$ , between pair  $(i, j)$  is added on the force on atom  $j$ , only if  $j$  is resident. Since all cached atoms are appended after the  $n$  resident atoms, atom  $j$  is resident if  $j < n$ :

```
if (j < n) ra[j][a] -= the a-th component of  $\vec{a}_{ij}$ 
```

Regarding the potential energy calculation, note that if  $j$  is cached, the same pair,  $(i, j)$ , is also considered by the process of which  $j$  is a resident. Our parallelization strategy is that each process adds 1/2 of the pair contribution:

```
if (j < n) potential energy +=  $u(r_{ij})$ ; else potential energy +=  $u(r_{ij})/2$ ;
```

Note that each process first computes the local contribution to the potential energy, and subsequently `MPI_Allreduce()` is used to compute the global summation of the local potential energies.

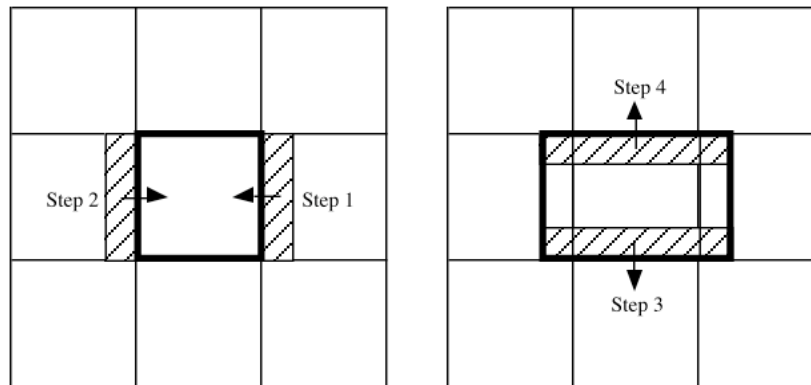
## Atom Caching—function `atom_copy()`

### DATA STRUCTURES

`lsb[6][NBMAX]`: `lsb[ku][0]` is the total number of boundary atoms to be sent to neighbor `ku`; `lsb[ku][k]` is the atom ID, used in `r`, of the `k`-th atom to be sent.

logical function `bbd(r,ku)`: true if the coordinate vector `r` is within distance `rc` from the face this subsystem shares with the `ku`-th neighbor.

The coordinates of the boundary atoms to the six face-sharing neighbors are sent to these nodes in six steps. Copies to the non-face-sharing neighbors ( $26 - 6 = 20$ ) are done by message forwarding.



### ALGORITHM

```

Reset the number of received cache atoms, nbnew = 0
for x, y, and z directions
  Make boundary-atom lists, lsb, for lower and higher directions including both
  resident, n, and cache, nbnew, atoms
  for lower and higher directions*
    Send/receive boundary-atom coordinates to/from the neighbor#
    Increment nbnew
  endfor
endfor
nb = nbnew

```

#### \*Copying condition

```

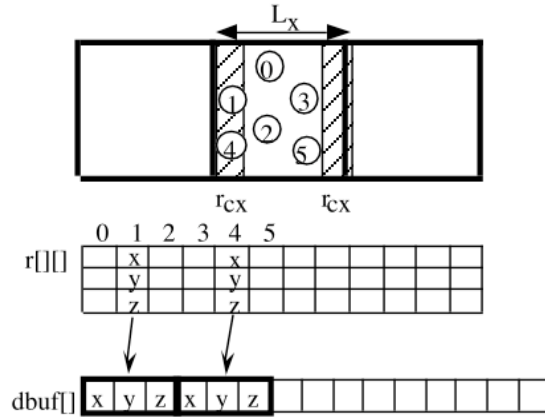
bbd(ri[],ku) {
  kd = ku / 2 (= 0|1|2)
  kdd = ku % 2 (= 0|1)
  if (kdd == 0)
    return ri[kd] < RCUT
  else
    return al[kd] - RCUT < ri[kd]
}

```

#### #Three-phase message passing

1. Message buffering: `dbuf ← r-sv` (shift positions), gather
2. Message passing: `dbuf ← dbufr`  
Send `dbuf`  
Receive `dbuf`
3. Message storing: `r ← dbufr`, append after the residents

Variable `nsd` is the number of atoms to be copied to the neighbor currently considered. The  $3 \cdot nsd$  coordinates are packed in a 1D array, `dbuf`.

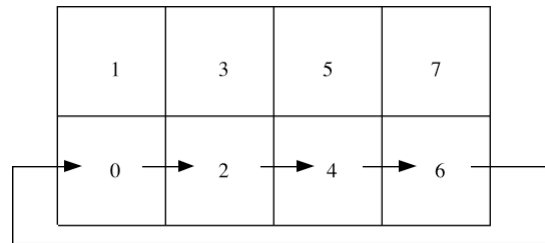


### DEADLOCK AVOIDANCE

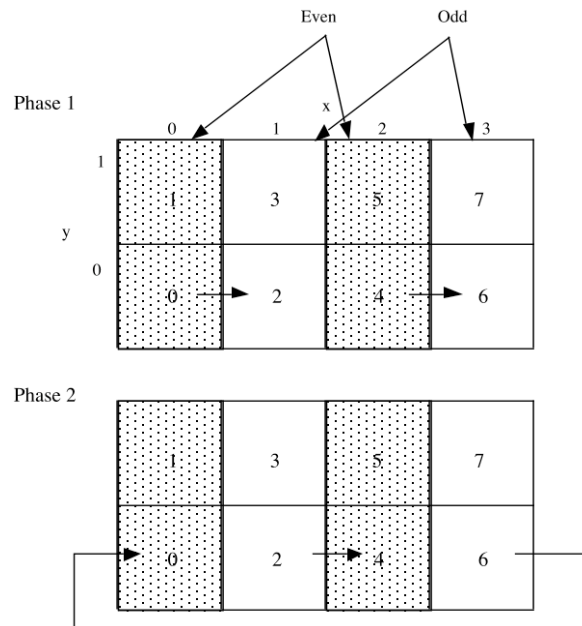
A circular send and receive relation could cause a deadlock. (Circular wait can occur if blocking.) In the following code, suppose that the destination node, `inode`, is defined as in the figure below.

```
MPI_Send(dbuf, 3*nsd, MPI_DOUBLE_PRECISION, inode, 10, MPI_COMM_WORLD);
MPI_Recv(dbuf, 3*nrc, MPI_DOUBLE_PRECISION, MPI_ANY_SOURCE, 10, MPI_COMM_WORLD, &status);
```

With a finite buffer size in the receiver's communication system, a sender blocks until the receiver's buffer is cleared. However in the example above, each processor cannot start receiving until its send operation is completed.



To avoid this deadlock, we classify the processors into even and odd processors in each of the x, y, and z directions.





## DATA STRUCTURES

$myparity[0|1|2]$  = Parity of vector processor ID in the x|y|z direction.  $myparity[a]$  is 0|1 if the vector processor ID,  $vid[a]$ , is even|odd. If there is only one processor in the a direction,  $myparity[a] = 2$ . In this case, no message passing is required.

## ALGORITHM

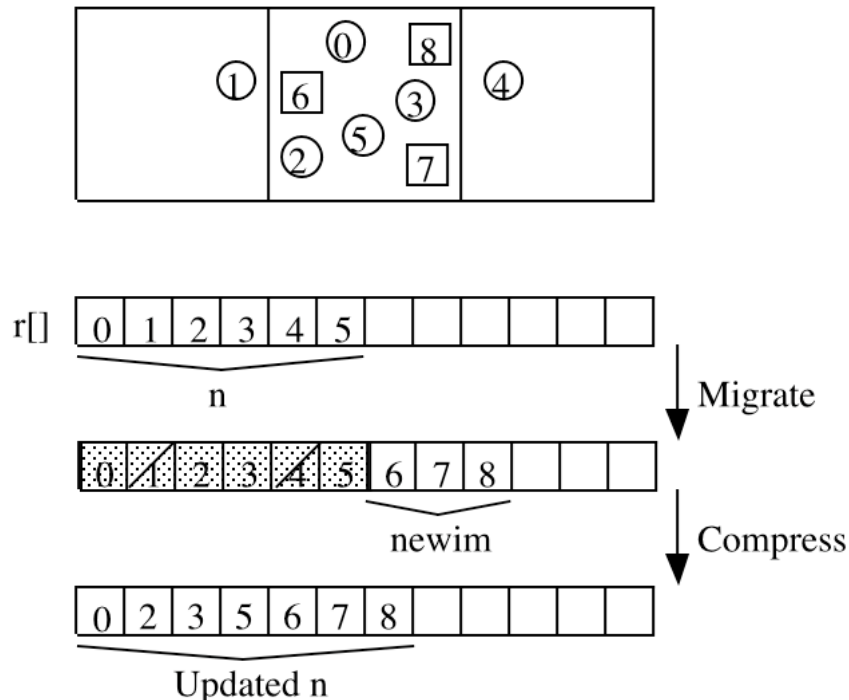
1. Message buffering:  $dbuf \leftarrow r$ , gather
2. Message passing:  $dbuffr \leftarrow dbuf$   
 /\* Even node: send & recv, if not empty \*/  
 if ( $myparity[kd] == 0$ ) {  
   MPI\_Send( $dbuf, 3*nsd, MPI\_DOUBLE, inode, 120, MPI\_COMM\_WORLD$ );  
   MPI\_Recv( $dbuffr, 3*nrc, MPI\_DOUBLE, MPI\_ANY\_SOURCE, 120,$   
     MPI\_COMM\_WORLD, &status);  
 }  
 /\* Odd node: recv & send, if not empty \*/  
 else if ( $myparity[kd] == 1$ ) {  
   MPI\_Recv( $dbuffr, 3*nrc, MPI\_DOUBLE, MPI\_ANY\_SOURCE, 120,$   
     MPI\_COMM\_WORLD, &status);  
   MPI\_Send( $dbuf, 3*nsd, MPI\_DOUBLE, inode, 120, MPI\_COMM\_WORLD$ );  
 }  
 /\* Single layer: Exchange information with myself \*/  
 else  
   for ( $i=0; i<3*nrc; i++$ )  $dbuffr[i] = dbuf[i]$ ;
3. Message storing:  $r \leftarrow dbufr$ , append

## Atom Migration—function `atom_move()`

### DATA STRUCTURES

$mvque[6][NBMAX]$ :  $mvque[ku][0]$  is the # of to-be-moved atoms to neighbor  $ku$ ;  $MVQUE[ku][k>0]$  is the atom ID, used in  $r$ , of the  $k$ -th atom to be moved.

logical function  $bmv(r, ku)$ : .true. if the coordinate vector  $r$  is out of the subsystem boundary plane which faces the  $ku$ -th neighbor.



The algorithm is basically the same as the one in `atom_copy`. During the 6-step loop, variable `newim` keeps the number of received new immigrants. These atoms are appended in the `r` array after the current resident. Suppose atom `i` has moved-out. When sending out `r[i]` to a proper neighbor, `r[i]` is marked as moved-out by putting constant `MOVED_OUT` (large negative value,  $-1e10$ , defined in `pmd.h`) in `r[i][0]`. When scanning `n+newim` atoms to sort out the moved-out atoms, these already-moved atoms are not considered. After the six move steps, the `r` array is compressed to remove the moved-out atoms.

## ALGORITHM

```

Reset the number of received new immigrants, newim = 0
for x, y, and z directions
  Make moving-atom lists, mvque, for lower and higher directions including both
  resident, n, and immigrant, newim, atoms but excluding those already moved out
  for lower and higher directions*
    Send/receive moving-atom coordinates to/from the neighbor#
    (When moving, r[][0] ← MOVED_OUT = -1010)
    Increment newim
  endfor
endfor
Compress the r array to eliminate the moved-out atoms

```

### \*Moving condition

```

bmv(ri[],ku) {
  kd = ku / 2 (= 0|1|2)
  kdd = ku % 2 (= 0|1)
  if (kdd == 0)
    return ri[kd] < 0.0
  else
    return al[kd] < ri[kd]
}

```

### #Three-phase message passing

1. Message buffering: `dobuf ← r & rv`, gather  
Mark `MOVED_OUT` in `r`
2. Message passing: `dobuf ← dobuf`  
Send `dobuf`  
Receive `dobuf`
3. Message storing: `r & rv ← dobuf`, append after the residents

