# *Communication with Double Buffering*

## Issaku Kanamori (RIKEN)

December 12, 2019 at R-CCS
Fugaku QCD Coding workshop

RIKEN
Center for
Computational Science

# Outline

1. Introduction
2. Algorithm with Double Buffering
3. Benchmark (A64FX 1 node)
4. Conclusions

# Outline

1. Introduction
2. Algorithm with Double Buffering
3. Benchmark (A64FX 1 node)
4. Conclusions

Acknowledgments
this talk is based on discussion with the codesign team for LQCD

# Introduction

Performance Bottle Neck on Lattice QCD

- the most time consuming: mult of $D$ in the solver
- memory bandwidth
- communication bandwidth

  - neighboring communication in $D$: need to wait for boundary data comes
  - overlapping communication and computation: as computation becomes faster, it becomes more difficult to hide communication

$\boxed{\text{double buffering algorithm}}$ may reduce the comm. overhead
(implementation: RDMA through the uTofu interface)

# Algorithm with Double Buffering

## process 1 (send)

```
 1    // 1st iter.                    send buffer
 2    start receiving
 3    pack the boundary data
 4    start sending
 5     computation: bulk
 6    wait for the boundary data comes
 7    computation: boundary
 8    wait for sending is done
 9    // 2nd iter.
10     start receiving
11    pack the boundary data
12    start sending
13     computation: bulk
14     wait for the boundary data comes
15     computation: boundary
16     wait for sending is done
17     ...
```

# Neighboring Communication

## process 1 (send)

```
              send buffer
1    // 1st iter.
2    start receiving
3    pack the boundary data
4    start sending
5     computation: bulk
6    wait for the boundary data comes
7    computation: boundary
8    wait for sending is done
9    // 2nd iter.
10    start receiving
11   pack the boundary data
12   start sending
13    computation: bulk
14    wait for the boundary data comes
15    computation: boundary
16    wait for sending is done
17    ...
```

## process 2 (recv.)

```
              recv. buffer
1    // 1st iter.
2    start receiving
3     pack the boundary data
4    start sending
5     computation: bulk
6    wait for the boundary data comes
7    computation: boundary
8    wait for sending is done
9    // 2nd iter.
10    start receiving
11    pack the boundary data
12    start sending
13    computation: bulk
14    wait for the boundary data comes
15   computation: boundary
16    wait for sending is done
17    ...
```

# Neighboring Communication

## process 1 (send)

| | |
|---|---|
| 1 | *// 1st iter.* |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | *// 2nd iter.* |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

send buffer

P

## process 2 (recv.)

| | |
|---|---|
| 1 | *// 1st iter.* |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | *// 2nd iter.* |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

recv. buffer

R

send buf.   P : Packed   P : Sending

recv. buf.   R : Receiving,   R : Receiving done,   U : being Used

# Neighboring Communication

## process 1 (send)

|   |   |
|---|---|
| 1 | // 1st iter. |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | // 2nd iter. |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

send buffer

## process 2 (recv.)

|   |   |
|---|---|
| 1 | // 1st iter. |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | // 2nd iter. |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

recv. buffer

send buf.　P : Packed　P : Sending

recv. buf.　R : Receiving,　R : Receiving done,　U : being Used

# Neighboring Communication

## process 1 (send)

```
1    // 1st iter.
2    start receiving
3    pack the boundary data
4    start sending
5     computation: bulk
6    wait for the boundary data comes
7    computation: boundary
8    wait for sending is done
9    // 2nd iter.
10    start receiving
11   pack the boundary data
12    start sending
13    computation: bulk
14    wait for the boundary data comes
15    computation: boundary
16    wait for sending is done
17    ...
```

send buffer

## process 2 (recv.)

```
1    // 1st iter.
2    start receiving
3    pack the boundary data
4    start sending
5     computation: bulk
6    wait for the boundary data comes
7    computation: boundary
8    wait for sending is done
9    // 2nd iter.
10    start receiving
11   pack the boundary data
12    start sending
13    computation: bulk
14    wait for the boundary data comes
15    computation: boundary
16    wait for sending is done
17    ...
```

recv. buffer

overlap btw. comm. and comp.

send buf.  **P** : Packed  **P** : Sending

recv. buf.  R : Receiving,  **R** : Receiving done,  **U** : being Used

# Neighboring Communication

## process 1 (send)

```
 1    // 1st iter.                                    send buffer
 2    start receiving                                      ☐
 3    pack the boundary data                               P
 4    start sending                                        S
 5     computation: bulk
 6    wait for the boundary data comes                     S
 7    computation: boundary
 8    wait for sending is done
 9    // 2nd iter.
10     start receiving
11    pack the boundary data
12    start sending
13     computation: bulk
14    wait for the boundary data comes
15    computation: boundary
16    wait for sending is done
17    ...
```

## process 2 (recv.)

```
 1    // 1st iter.                                    recv. buffer
 2    start receiving                                      R
 3    pack the b...                   overlap btw. comm. and comp.
 4    start send...
 5    computation: bulk
 6    wait for the boundary data comes                     R
 7    computation: boundary                                U
 8    wait for sending is done
 9    // 2nd iter.
10    start receiving
11    pack the boundary data
12    start sending
13    computation: bulk
14    wait for the boundary data comes
15    computation: boundary
16    wait for sending is done
17    ...
```
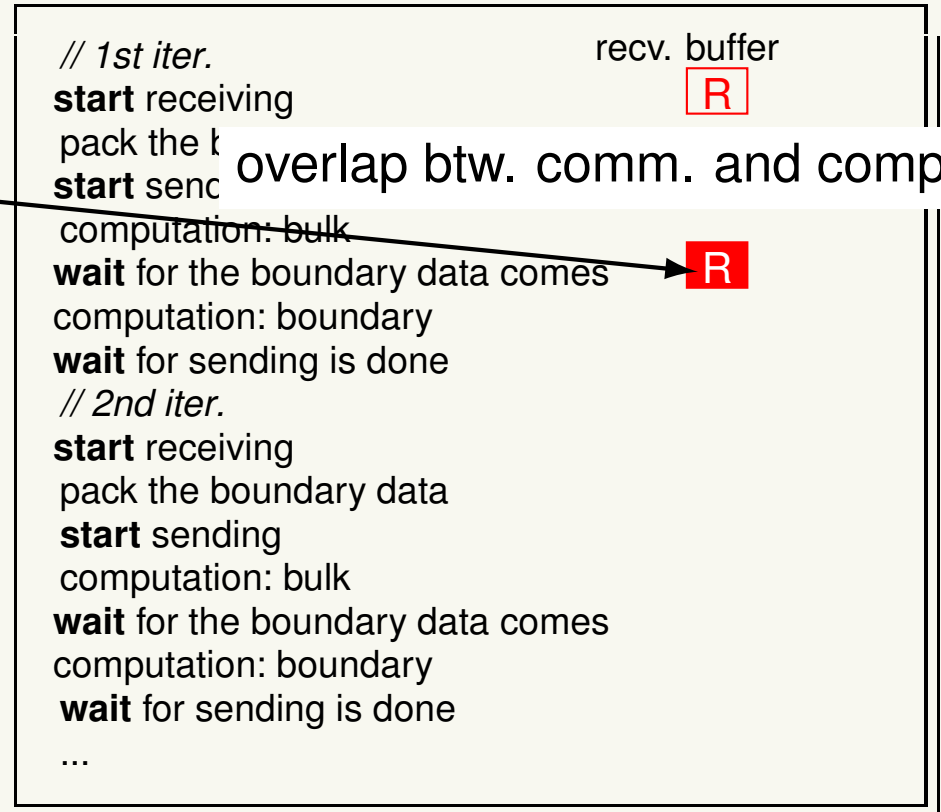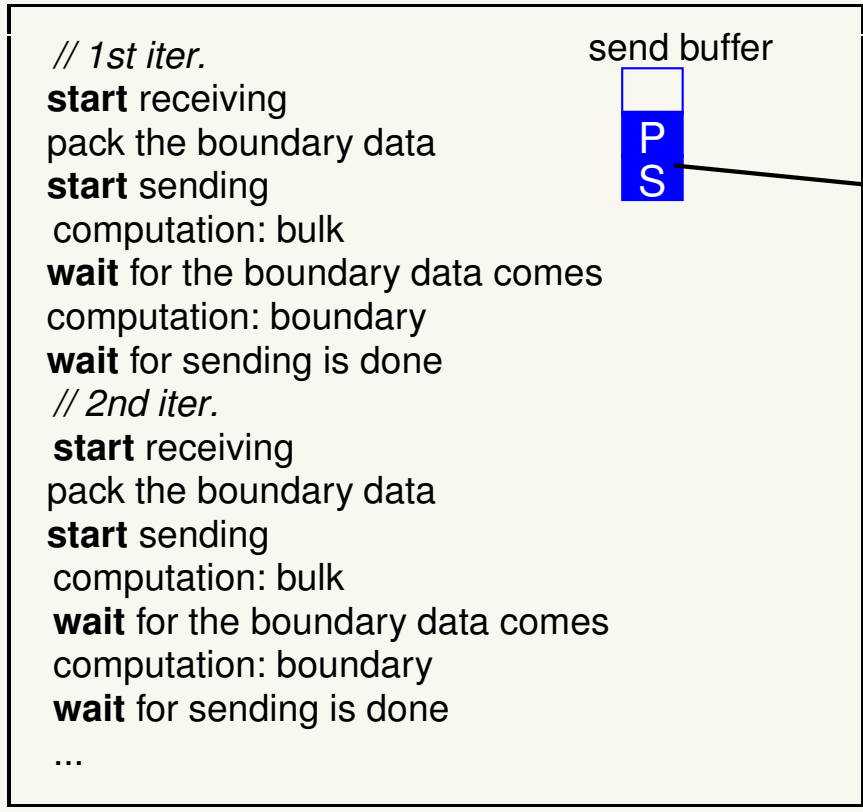
send buf.     **P** : Packed   **P** : Sending

recv. buf.    **R** : Receiving,   **R** : Receiving done,   **U** : being Used

# Neighboring Communication

## process 1 (send)

send buffer

| | |
|---|---|
| 1 | *// 1st iter.* |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | *// 2nd iter.* |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

P
S

S

## process 2 (recv.)

recv. buffer

| | |
|---|---|
| 1 | *// 1st iter.* |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | *// 2nd iter.* |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

R

overlap btw. comm. and comp.

R
U

send buf.   P : Packed   P : Sending

recv. buf.   R : Receiving,   R : Receiving done,   U : being Used

# Neighboring Communication

## process 1 (send)

| | | send buffer |
|---|---|---|
| 1 | // 1st iter. | |
| 2 | **start** receiving | |
| 3 | pack the boundary data | **P** |
| 4 | **start** sending | **S** |
| 5 | computation: bulk | |
| 6 | **wait** for the boundary data comes | |
| 7 | computation: boundary | **S** |
| 8 | **wait** for sending is done | |
| 9 | // 2nd iter. | |
| 10 | **start** receiving | |
| 11 | pack the boundary data | **P** |
| 12 | **start** sending | **S** |
| 13 | computation: bulk | |
| 14 | **wait** for the boundary data comes | |
| 15 | computation: boundary | |
| 16 | **wait** for sending is done | |
| 17 | ... | |

## process 2 (recv.)

| | | recv. buffer |
|---|---|---|
| 1 | // 1st iter. | |
| 2 | **start** receiving | R |
| 3 | pack the b... | |
| 4 | **start** send... | |
| 5 | computation: bulk | |
| 6 | **wait** for the boundary data comes | **R** / **U** |
| 7 | computation: boundary | |
| 8 | **wait** for sending is done | |
| 9 | // 2nd iter. | |
| 10 | **start** receiving | R |
| 11 | pack the boundary data | |
| 12 | **start** sending | |
| 13 | computation: bulk | |
| 14 | **wait** for the boundary data comes | **R** / **U** |
| 15 | computation: boundary | |
| 16 | **wait** for sending is done | |
| 17 | ... | |

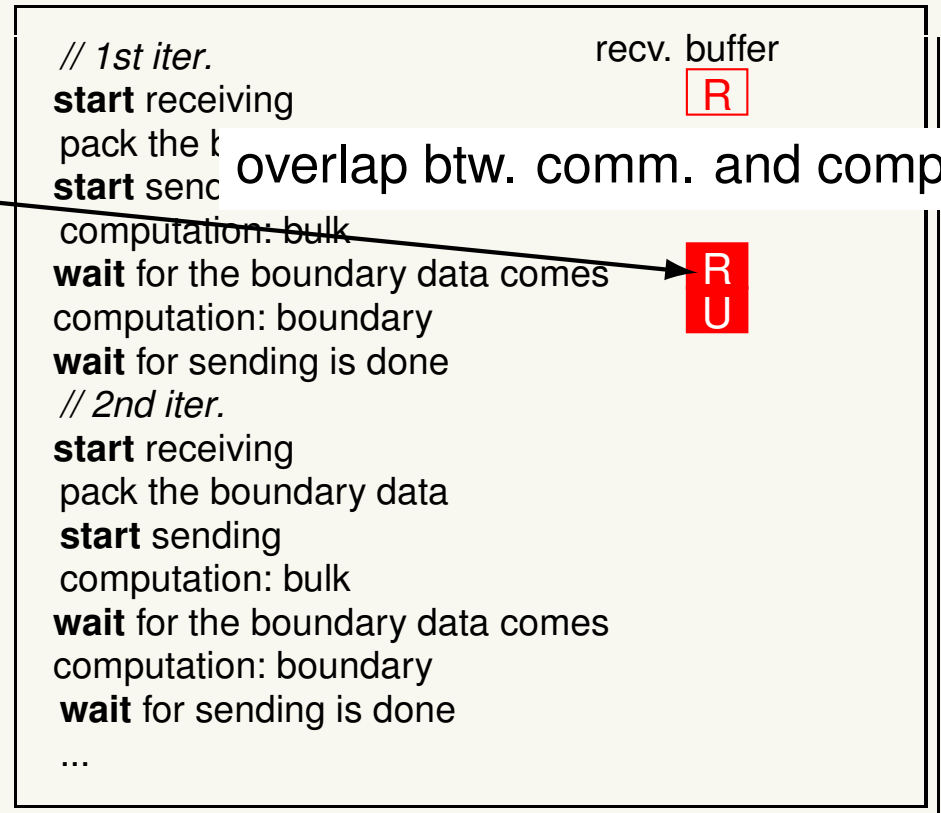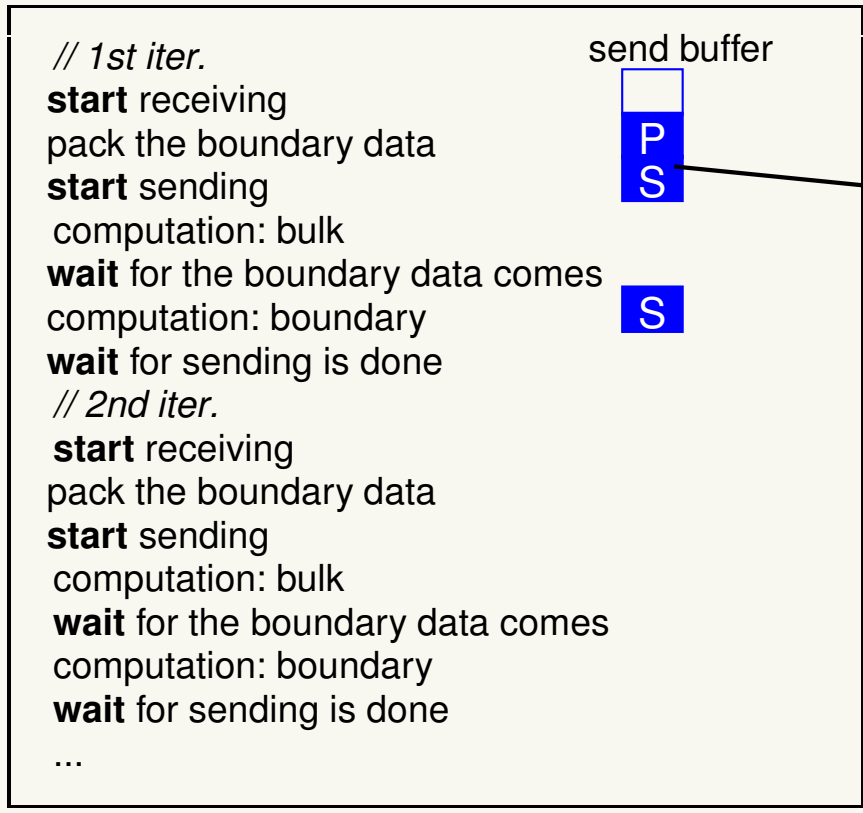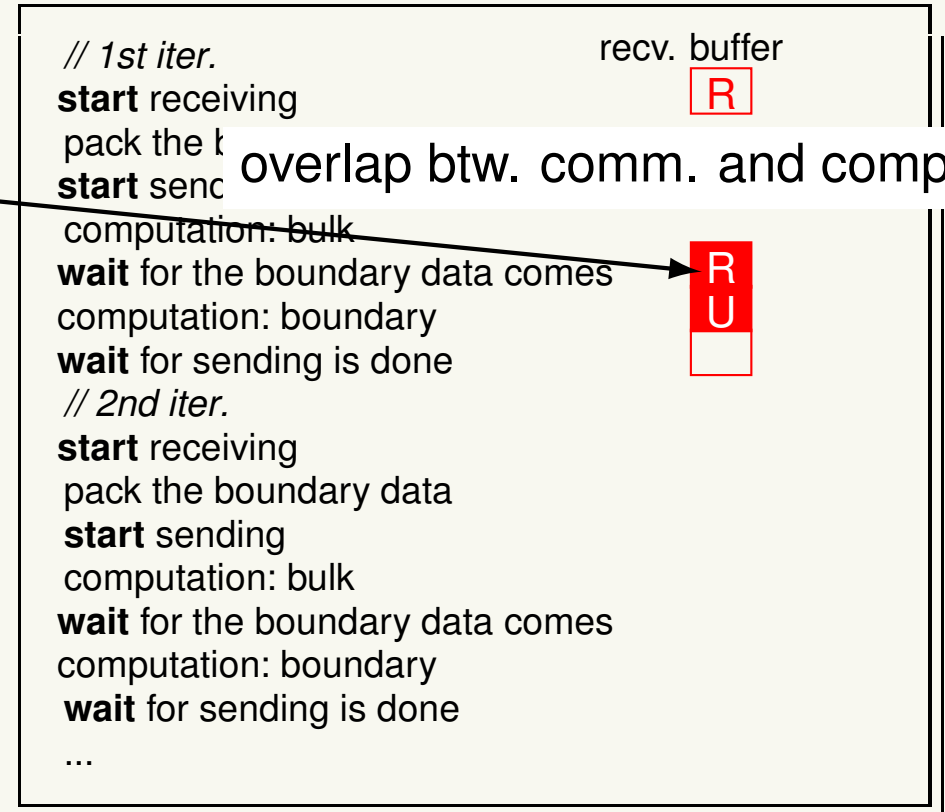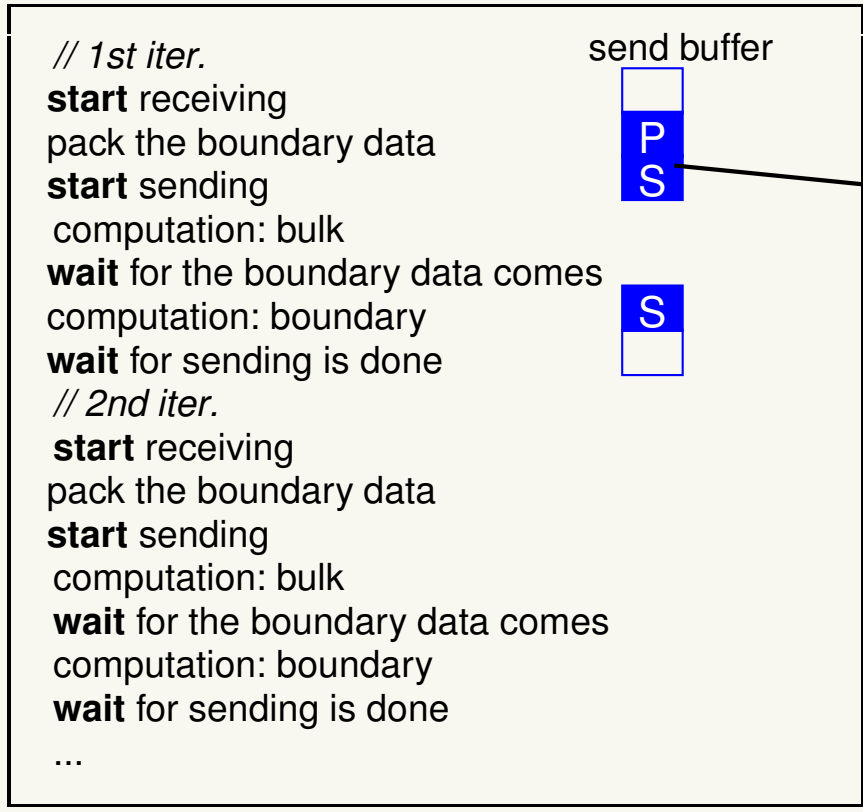overlap btw. comm. and comp.

send buf.    **P** : Packed   **P** : Sending

recv. buf.    R : Receiving,  **R** : Receiving done,  **U** : being Used

# Neighboring Communication

## process 1 (send)

```
1    // 1st iter.                         send buffer
2    start receiving                          [ ]
3    pack the boundary data                    P
4    start sending                             S
5     computation: bulk
6    wait for the boundary data comes
7    computation: boundary                     S
8    wait for sending is done                 [ ]
9    // 2nd iter.
10    start receiving
11   pack the boundary data                    P
12   start sending                             S
13    computation: bulk
14   wait for the boundary data comes
15   computation: boundary
16   wait for sending is done
17   ...
```

## process 2 (recv.)

```
1    // 1st iter.                         recv. buffer
2    start receiving                          [R]
3    pack the b
4    start sen                                       overlap btw. comm. and comp.
5     computation: bulk
6    wait for the boundary data comes          R
7    computation: boundary                      U
8    wait for sending is done                 [ ]
9    // 2nd iter.
10    start receiving                         [R]
11   pack the boundary data
12   start sending
13    computation: bulk
14   wait for the boundary data comes          R
15   computation: boundary                      U
16    wait for sending is done
17   ...
```
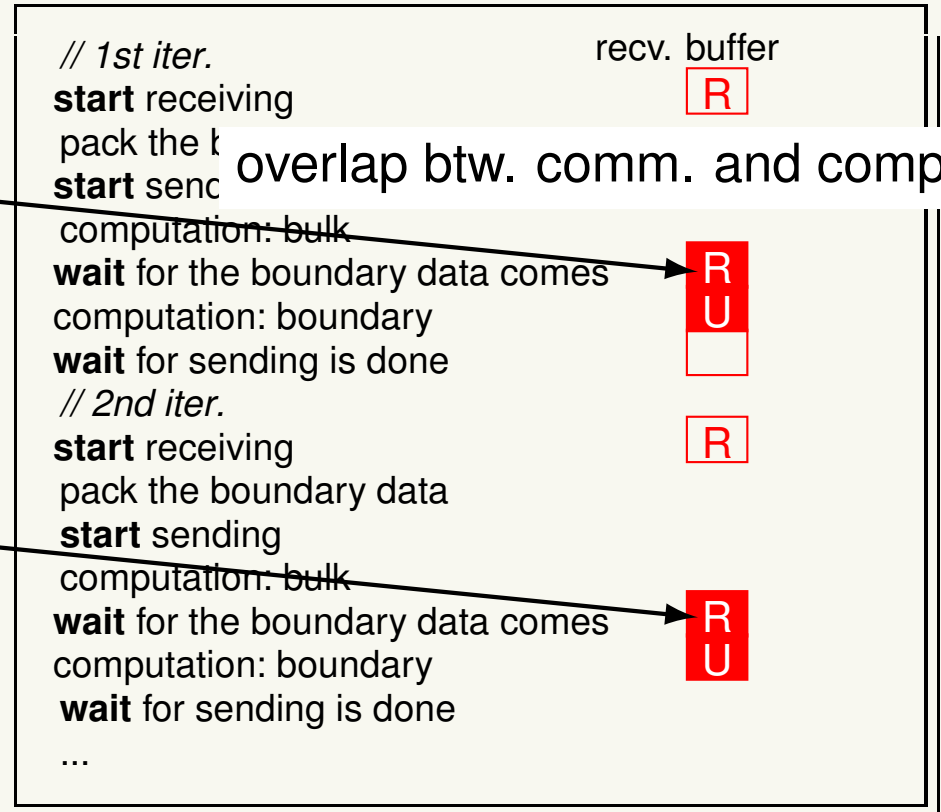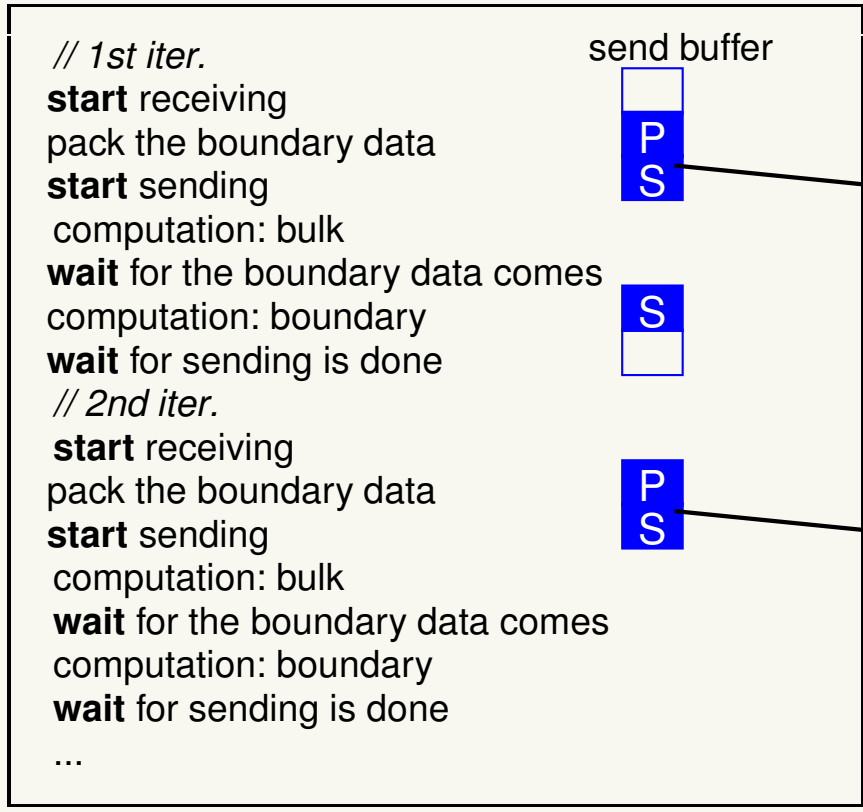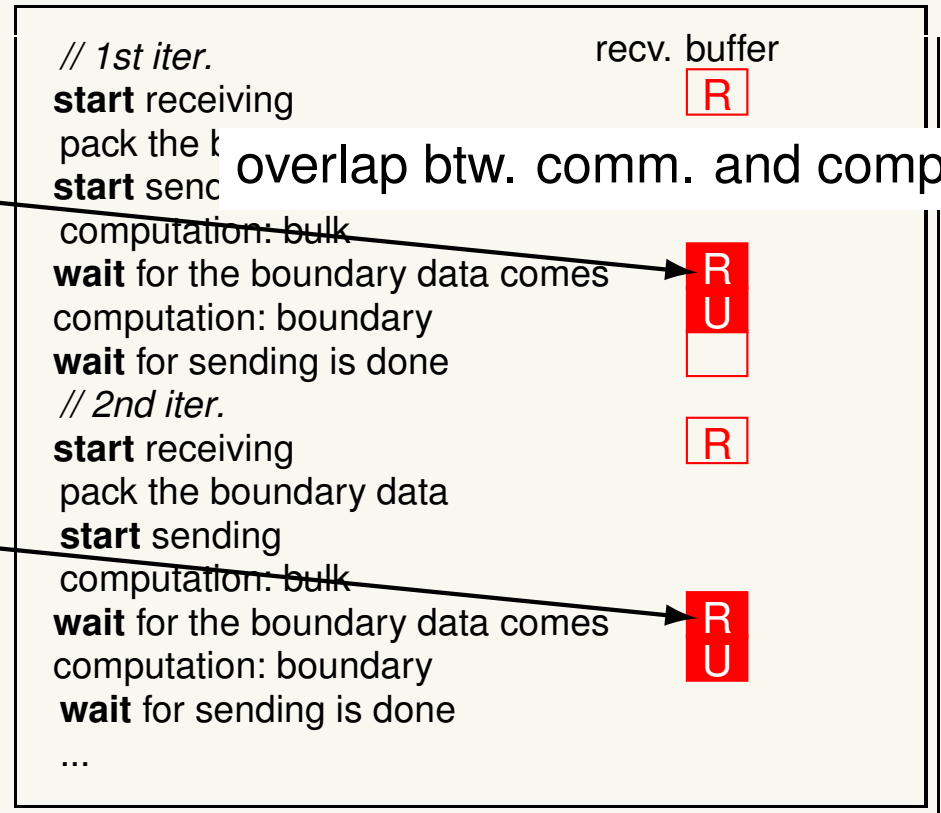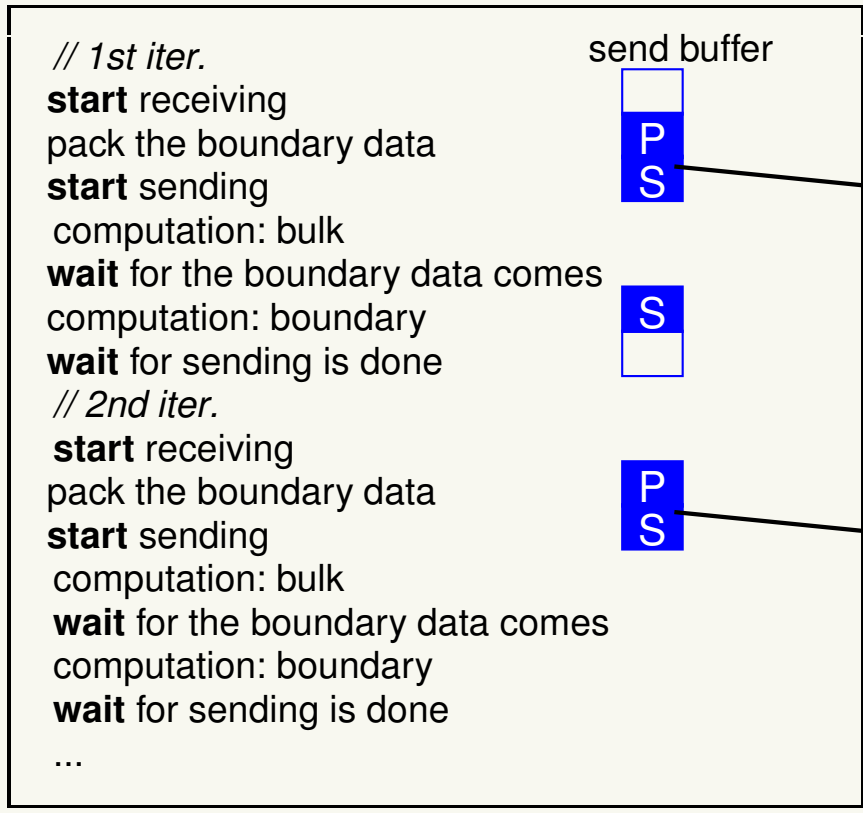
send buf.    **P** : Packed   **P** : Sending

recv. buf.   [R] : Receiving,   **R** : Receiving done,   **U** : being Used

# Neighboring Communication, cont'd

## If the receiving process delays....

<div style="border: 1px solid black; padding: 10px;">

send buffer

```
1   // 1st iter.
2   start receiving
3   pack the boundary data
4   start sending
5    computation: bulk
6   wait for the boundary data comes
7   computation: boundary
8   wait for sending is done
9   // 2nd iter.
10   start receiving
11  pack the boundary data
12  start sending
13   computation: bulk
14   wait for the boundary data comes
15   computation: boundary
16  wait for sending is done
17   ...
```

</div>

<div style="border: 1px solid black; padding: 10px;">

recv. buffer

```
1   // 1st iter.
2   start receiving
3    pack the boundary data
4   start sending
5    computation: bulk
6
7   wait for the boundary data comes
8
9   computation: boundary
10
11
12  wait for sending is done
13
14
15   // 2nd iter.
16  start receiving
17   pack the boundary data
18   start sending
19   computation: bulk
20  wait for the boundary data comes
21  computation: boundary
22   wait for sending is done
23   ...
```

</div>

# Neighboring Communication, cont'd

## If the receiving process delays....

| | | send buffer |
|---|---|---|
| 1 | *// 1st iter.* | |
| 2 | **start** receiving | |
| 3 | pack the boundary data | |
| 4 | **start** sending | |
| 5 | computation: bulk | |
| 6 | **wait** for the boundary data comes | |
| 7 | computation: boundary | |
| 8 | **wait** for sending is done | |
| 9 | *// 2nd iter.* | |
| 10 | **start** receiving | |
| 11 | pack the boundary data | |
| 12 | **start** sending | |
| 13 | computation: bulk | |
| 14 | **wait** for the boundary data comes | |
| 15 | computation: boundary | |
| 16 | **wait** for sending is done | |
| 17 | ... | |

| | | recv. buffer |
|---|---|---|
| 1 | *// 1st iter.* | |
| 2 | **start** receiving | |
| 3 | pack the boundary data | |
| 4 | **start** sending | |
| 5 | computation: bulk | |
| 6 | | |
| 7 | **wait** for the boundary data comes | |
| 8 | | |
| 9 | computation: boundary | |
| 10 | | |
| 11 | | |
| 12 | **wait** for sending is done | |
| 13 | | |
| 14 | | |
| 15 | *// 2nd iter.* | |
| 16 | **start** receiving | |
| 17 | pack the boundary data | |
| 18 | **start** sending | |
| 19 | computation: bulk | |
| 20 | **wait** for the boundary data comes | |
| 21 | computation: boundary | |
| 22 | **wait** for sending is done | |
| 23 | ... | |

# Neighboring Communication, cont'd

## If the receiving process delays....

| | send buffer |
|---|---|
| 1 | *// 1st iter.* |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | *// 2nd iter.* |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

send buffer: P

| | recv. buffer |
|---|---|
| 1 | *// 1st iter.* |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | |
| 7 | **wait** for the boundary data comes |
| 8 | |
| 9 | computation: boundary |
| 10 | |
| 11 | |
| 12 | **wait** for sending is done |
| 13 | |
| 14 | |
| 15 | *// 2nd iter.* |
| 16 | **start** receiving |
| 17 | pack the boundary data |
| 18 | **start** sending |
| 19 | computation: bulk |
| 20 | **wait** for the boundary data comes |
| 21 | computation: boundary |
| 22 | **wait** for sending is done |
| 23 | ... |

recv. buffer: R

# Neighboring Communication, cont'd

If the receiving process delays....

```
 1    // 1st iter.                           send buffer
 2    start receiving
 3    pack the boundary data                      P
 4    start sending                               S
 5     computation: bulk
 6    wait for the boundary data comes
 7    computation: boundary
 8    wait for sending is done
 9    // 2nd iter.
10     start receiving
11    pack the boundary data
12    start sending
13     computation: bulk
14     wait for the boundary data comes
15     computation: boundary
16    wait for sending is done
17     ...
```

```
 1    // 1st iter.                           recv. buffer
 2    start receiving                             R
 3     pack the boundary data
 4    start sending
 5     computation: bulk
 6
 7    wait for the boundary data comes
 8
 9    computation: boundary
10
11
12    wait for sending is done
13
14
15    // 2nd iter.
16    start receiving
17     pack the boundary data
18     start sending
19     computation: bulk
20    wait for the boundary data comes
21    computation: boundary
22     wait for sending is done
23     ...
```

# Neighboring Communication, cont'd

If the receiving process delays....

| | |
|---|---|
| 1 | // 1st iter. |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | // 2nd iter. |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

send buffer

P
S

| | |
|---|---|
| 1 | // 1st iter. |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | |
| 7 | **wait** for the boundary data comes |
| 8 | |
| 9 | computation: boundary |
| 10 | |
| 11 | |
| 12 | **wait** for sending is done |
| 13 | |
| 14 | |
| 15 | // 2nd iter. |
| 16 | **start** receiving |
| 17 | pack the boundary data |
| 18 | **start** sending |
| 19 | computation: bulk |
| 20 | **wait** for the boundary data comes |
| 21 | computation: boundary |
| 22 | **wait** for sending is done |
| 23 | ... |

recv. buffer

R

R

# Neighboring Communication, cont'd

If the receiving process delays....

| | |
|---|---|
| 1 | // 1st iter. |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | // 2nd iter. |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

send buffer
P
S
S

| | |
|---|---|
| 1 | // 1st iter. |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | |
| 7 | **wait** for the boundary data comes |
| 8 | |
| 9 | computation: boundary |
| 10 | |
| 11 | |
| 12 | **wait** for sending is done |
| 13 | |
| 14 | |
| 15 | // 2nd iter. |
| 16 | **start** receiving |
| 17 | pack the boundary data |
| 18 | **start** sending |
| 19 | computation: bulk |
| 20 | **wait** for the boundary data comes |
| 21 | computation: boundary |
| 22 | **wait** for sending is done |
| 23 | ... |

recv. buffer
R
R

# Neighboring Communication, cont'd

## If the receiving process delays....

| | |
|---|---|
| 1 | // 1st iter. |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | **wait** for the boundary data comes |
| 7 | computation: boundary |
| 8 | **wait** for sending is done |
| 9 | // 2nd iter. |
| 10 | **start** receiving |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | **wait** for sending is done |
| 17 | ... |

send buffer

P
S

S

| | |
|---|---|
| 1 | // 1st iter. |
| 2 | **start** receiving |
| 3 | pack the boundary data |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | |
| 7 | **wait** for the boundary data comes |
| 8 | |
| 9 | computation: boundary |
| 10 | |
| 11 | |
| 12 | **wait** for sending is done |
| 13 | |
| 14 | |
| 15 | // 2nd iter. |
| 16 | **start** receiving |
| 17 | pack the boundary data |
| 18 | **start** sending |
| 19 | computation: bulk |
| 20 | **wait** for the boundary data comes |
| 21 | computation: boundary |
| 22 | **wait** for sending is done |
| 23 | ... |

recv. buffer

R

R

If the receiving process delays....

```
 1    // 1st iter.                          send buffer
 2    start receiving                       [  ]
 3    pack the boundary data                [P]
 4    start sending                         [S]
 5     computation: bulk
 6    wait for the boundary data comes      [S]
 7    computation: boundary
 8    wait for sending is done              [  ]
 9    // 2nd iter.
10     start receiving
11    pack the boundary data                [P]
12    start sending                         [S]
13     computation: bulk
14     wait for the boundary data comes
15     computation: boundary
16    wait for sending is done
17    ...
```
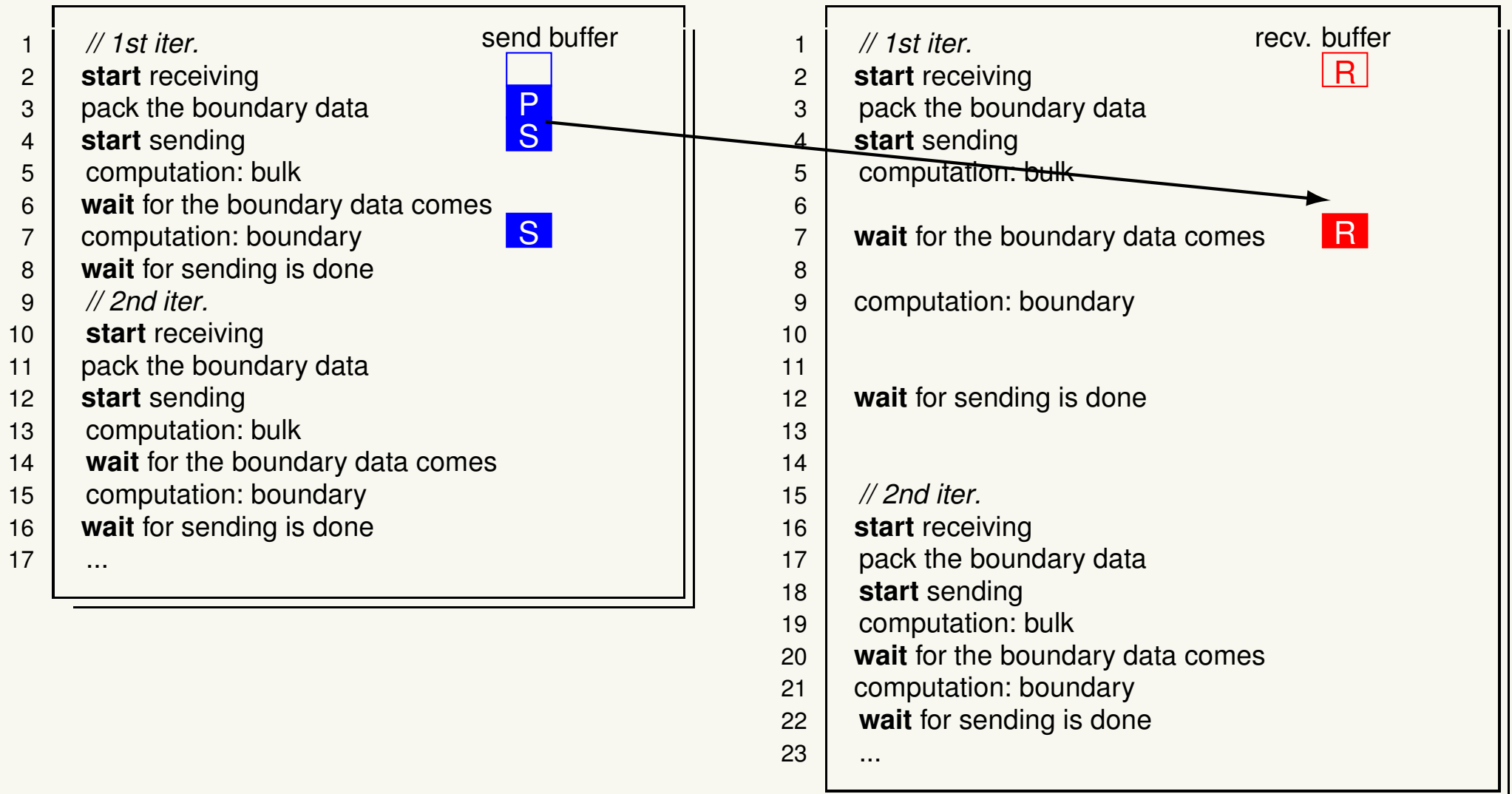
```
 1    // 1st iter.                          recv. buffer
 2    start receiving                       [R]
 3     pack the boundary data
 4    start sending
 5     computation: bulk
 6
 7    wait for the boundary data comes      [R]
 8
 9    computation: boundary                 [U]
                                            [U]
10                                          [U]
11                                          [U]
12    wait for sendi...                     cannot receive yet
13
14
15    // 2nd iter.
16    start receiving
17     pack the boundary data
18     start sending
19     computation: bulk
20    wait for the boundary data comes
21    computation: boundary
22     wait for sending is done
23    ...
```

If the receiving process delays....



```
                                      send buffer
1    // 1st iter.
2    start receiving                      □
3    pack the boundary data               P
4    start sending                        S
5     computation: bulk
6    wait for the boundary data comes     S
7    computation: boundary
8    wait for sending is done             □
9    // 2nd iter.
10    start receiving
11   pack the boundary data               P
12   start sending                        S
13    computation: bulk
14    wait for the boundary data comes
15    computation: boundary
16   wait for sending is done             S
17    ...
```

```
                                      recv. buffer
1    // 1st iter.
2    start receiving                      R
3     pack the boundary data
4    start sending
5     computation: bulk
6
7    wait for the boundary data comes     R
8
9    computation: boundary                U
10                                        U
11                                        U
12   wait for sendi...
13
14
15    // 2nd iter.
16   start receiving
17    pack the boundary data
     t sending
     putation: bulk
20   wait for the boundary data comes
21   computation: boundary
22    wait for sending is done
23    ...
```

cannot receive yet

takes time before finish sending

If the receiving process delays....

|   | | send buffer |
|---|---|---|
| 1 | // 1st iter. | |
| 2 | **start** receiving | |
| 3 | pack the boundary data | P |
| 4 | **start** sending | S |
| 5 | computation: bulk | |
| 6 | **wait** for the boundary data comes | |
| 7 | computation: boundary | S |
| 8 | **wait** for sending is done | |
| 9 | // 2nd iter. | |
| 10 | **start** receiving | |
| 11 | pack the boundary data | P |
| 12 | **start** sending | S |
| 13 | computation: bulk | |
| 14 | **wait** for the boundary data comes | |
| 15 | computation: boundary | |
| 16 | **wait** for sending is done | S |
| 17 | ... | |

|   | | recv. buffer |
|---|---|---|
| 1 | // 1st iter. | |
| 2 | **start** receiving | R |
| 3 | pack the boundary data | |
| 4 | **start** sending | |
| 5 | computation: bulk | |
| 6 | | |
| 7 | **wait** for the boundary data comes | R |
| 8 | | |
| 9 | computation: boundary | U |
| 10 | | U |
| 11 | | U |
| 12 | **wait** for sending is done | |
| 13 | | |
| 14 | | |
| 15 | // 2nd iter. | |
| 16 | **start** receiving | R |
| 17 | pack the boundary data | |
| | start sending | |
| | putation: bulk | |
| 20 | **wait** for the boundary data comes | R |
| 21 | computation: boundary | |
| 22 | **wait** for sending is done | |
| 23 | ... | |

cannot receive yet

takes time before finish sending

If the receiving process delays....

| | | send buffer |
|---|---|---|
| 1 | // 1st iter. | |
| 2 | **start** receiving | |
| 3 | pack the boundary data | P |
| 4 | **start** sending | S |
| 5 | computation: bulk | |
| 6 | **wait** for the boundary data comes | |
| 7 | computation: boundary | S |
| 8 | **wait** for sending is done | |
| 9 | // 2nd iter. | |
| 10 | **start** receiving | |
| 11 | pack the boundary data | P |
| 12 | **start** sending | S |
| 13 | computation: bulk | |
| 14 | **wait** for the boundary data comes | |
| 15 | computation: boundary | |
| 16 | **wait** for sending is done | S |
| 17 | ... | |

| | | recv. buffer |
|---|---|---|
| 1 | // 1st iter. | |
| 2 | **start** receiving | R |
| 3 | pack the boundary data | |
| 4 | **start** sending | |
| 5 | computation: bulk | |
| 6 | | |
| 7 | **wait** for the boundary data comes | R |
| 8 | | |
| 9 | computation: boundary | U |
| 10 | | U |
| 11 | | U |
| 12 | **wait** for sendi... | |
| 13 | | |
| 14 | | |
| 15 | // 2nd iter. | |
| 16 | **start** receiving | R |
| 17 | pack the boundary data | |
| | t sending | |
| | putation: bulk | |
| 20 | **wait** for the boundary data comes | R |
| 21 | computation: boundary | U |
| 22 | **wait** for sending is done | |
| 23 | ... | |

cannot receive yet

takes time before finish sending

# Neighboring Communication: Double Buffering (RDMA)

```
1    // 1st iter.                    send buffer
2    pack the boundary data
3    start sending
4     computation: bulk
5    wait for the boundary data comes
6    computation: boundary
7     clear the received flag
8    wait for sending is done
9    switch the buffer to send
10    // 2nd iter.
11   pack the boundary data
12   start sending
13    computation: bulk
14    wait for the boundary data comes
15    computation: boundary
16    clear the received flag
17    wait for sending is done
18   switch the buffer to send
19    ...
```
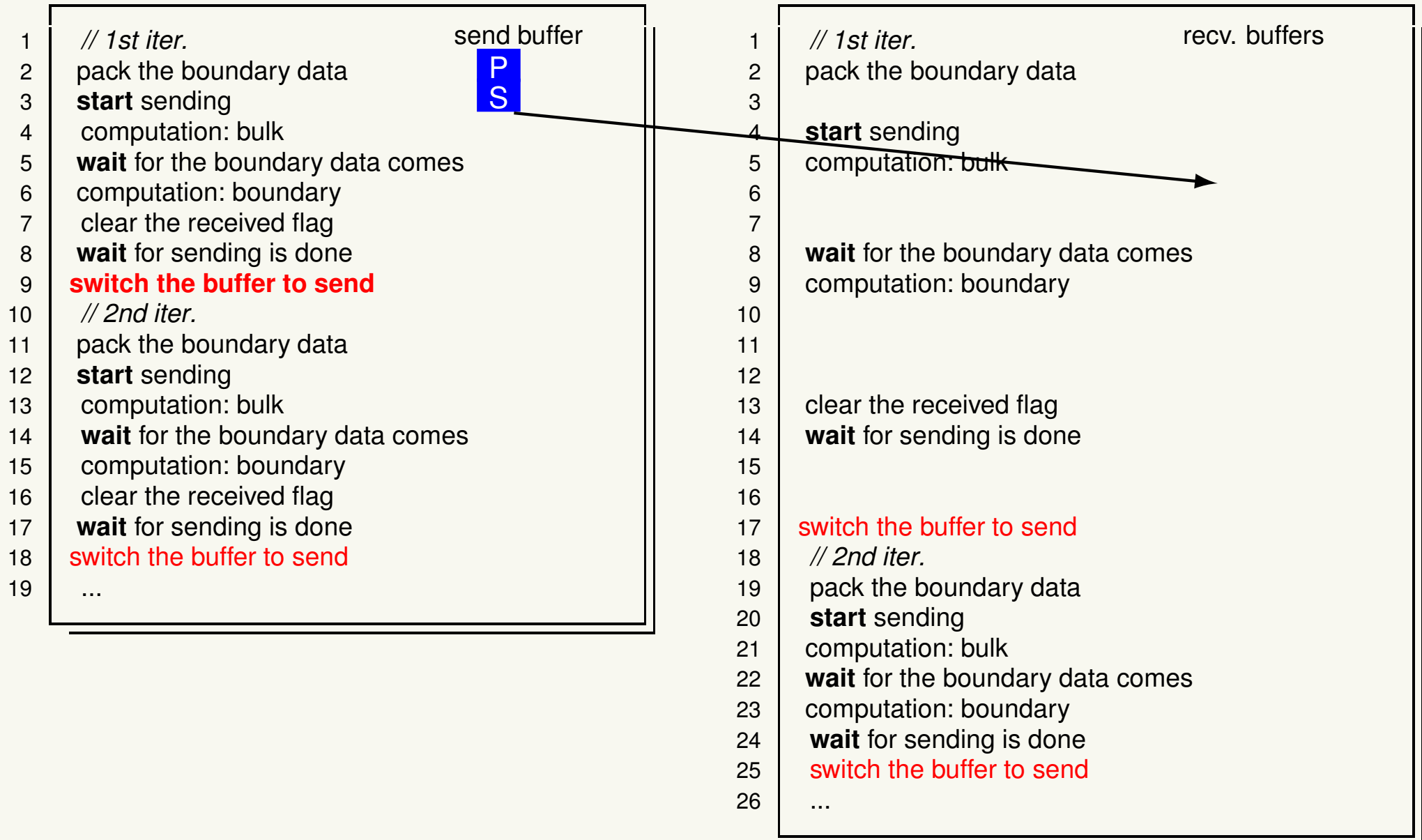
# Neighboring Communication: Double Buffering (RDMA)

send buffer

```
1    // 1st iter.
2    pack the boundary data
3    start sending
4     computation: bulk
5    wait for the boundary data comes
6    computation: boundary
7     clear the received flag
8    wait for sending is done
9    switch the buffer to send
10    // 2nd iter.
11   pack the boundary data
12   start sending
13    computation: bulk
14    wait for the boundary data comes
15    computation: boundary
16    clear the received flag
17    wait for sending is done
18   switch the buffer to send
19    ...
```
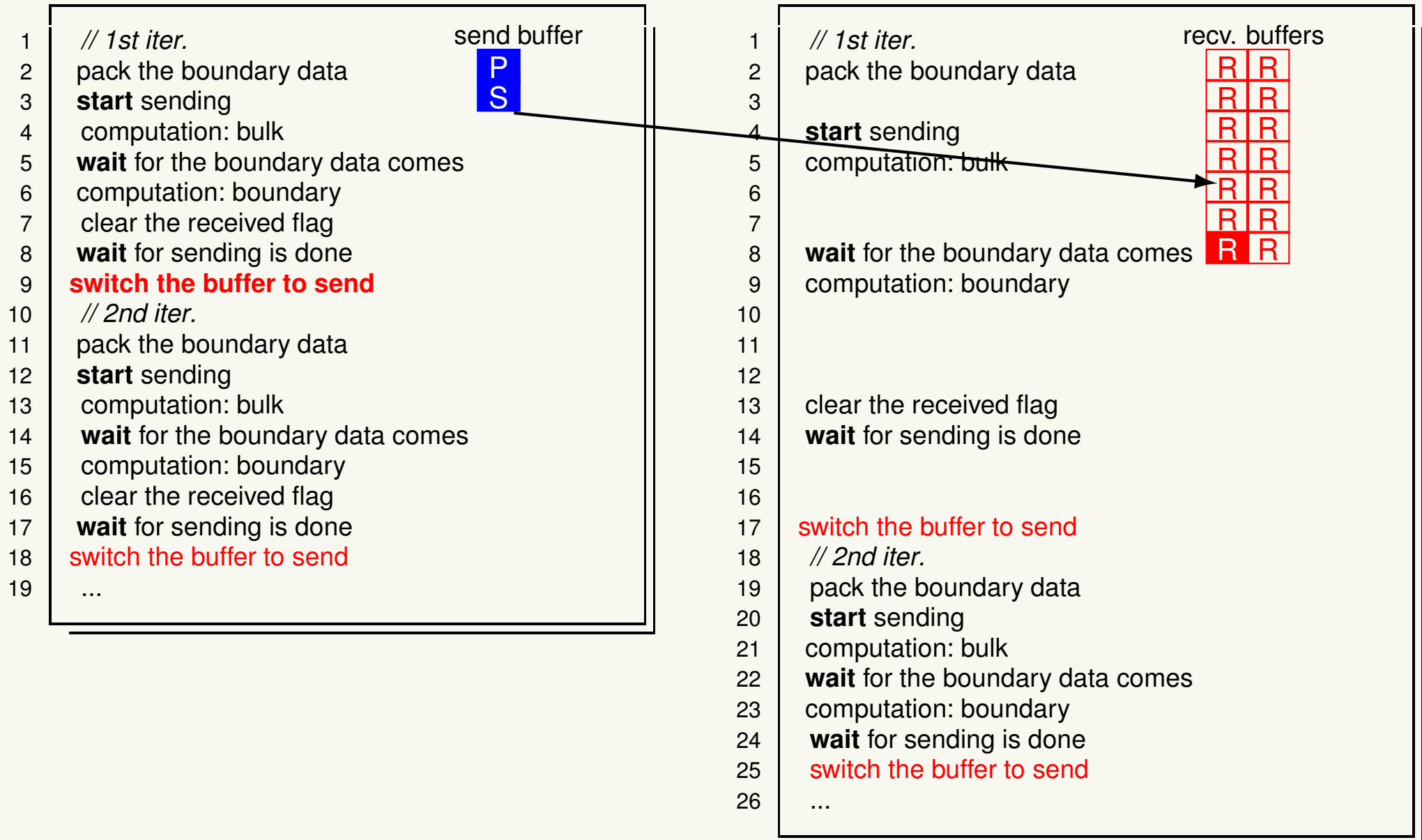
recv. buffers

```
1    // 1st iter.
2    pack the boundary data
3
4    start sending
5    computation: bulk
6
7
8    wait for the boundary data comes
9    computation: boundary
10
11
12
13    clear the received flag
14    wait for sending is done
15
16
17   switch the buffer to send
18    // 2nd iter.
19    pack the boundary data
20    start sending
21   computation: bulk
22   wait for the boundary data comes
23   computation: boundary
24    wait for sending is done
25    switch the buffer to send
26    ...
```
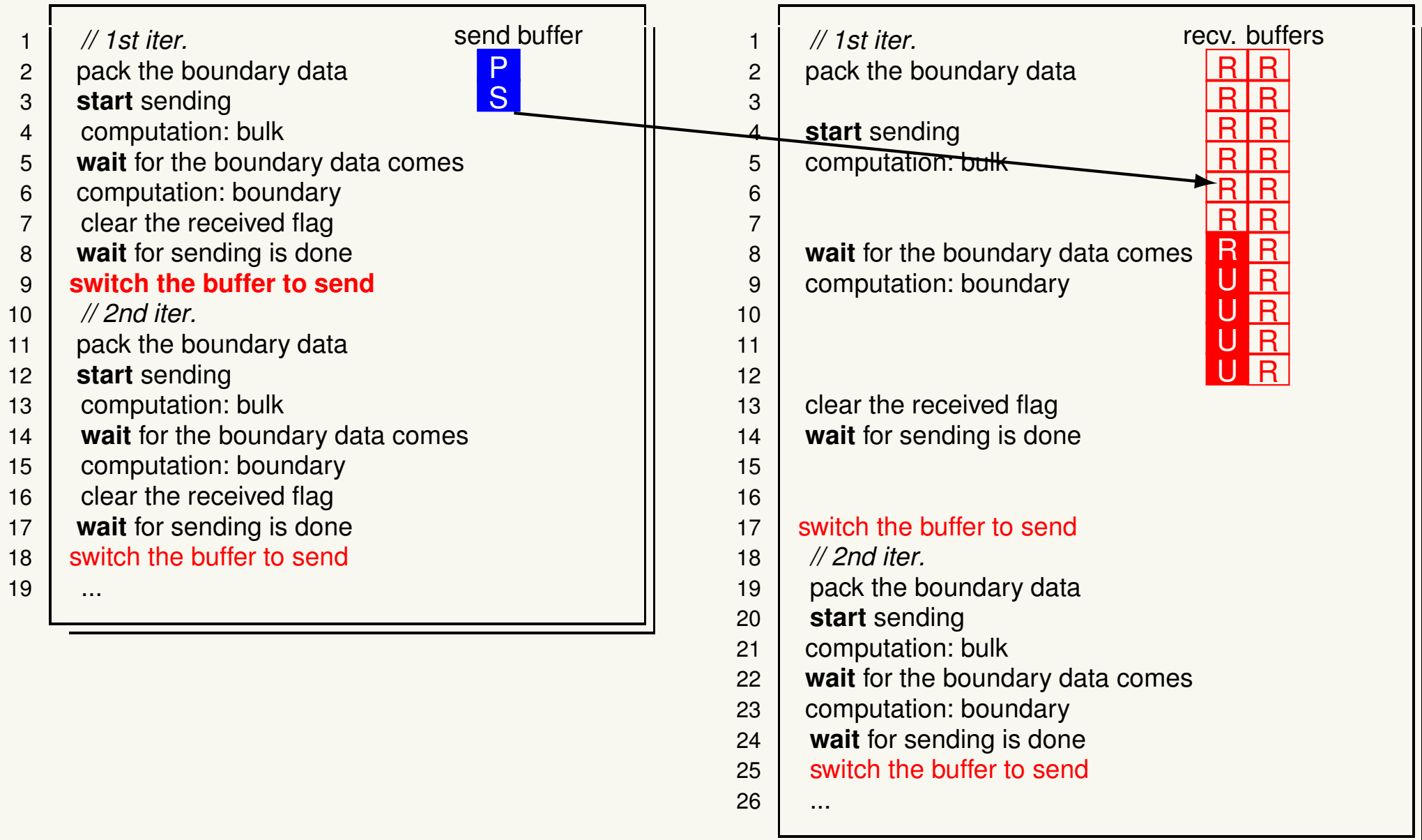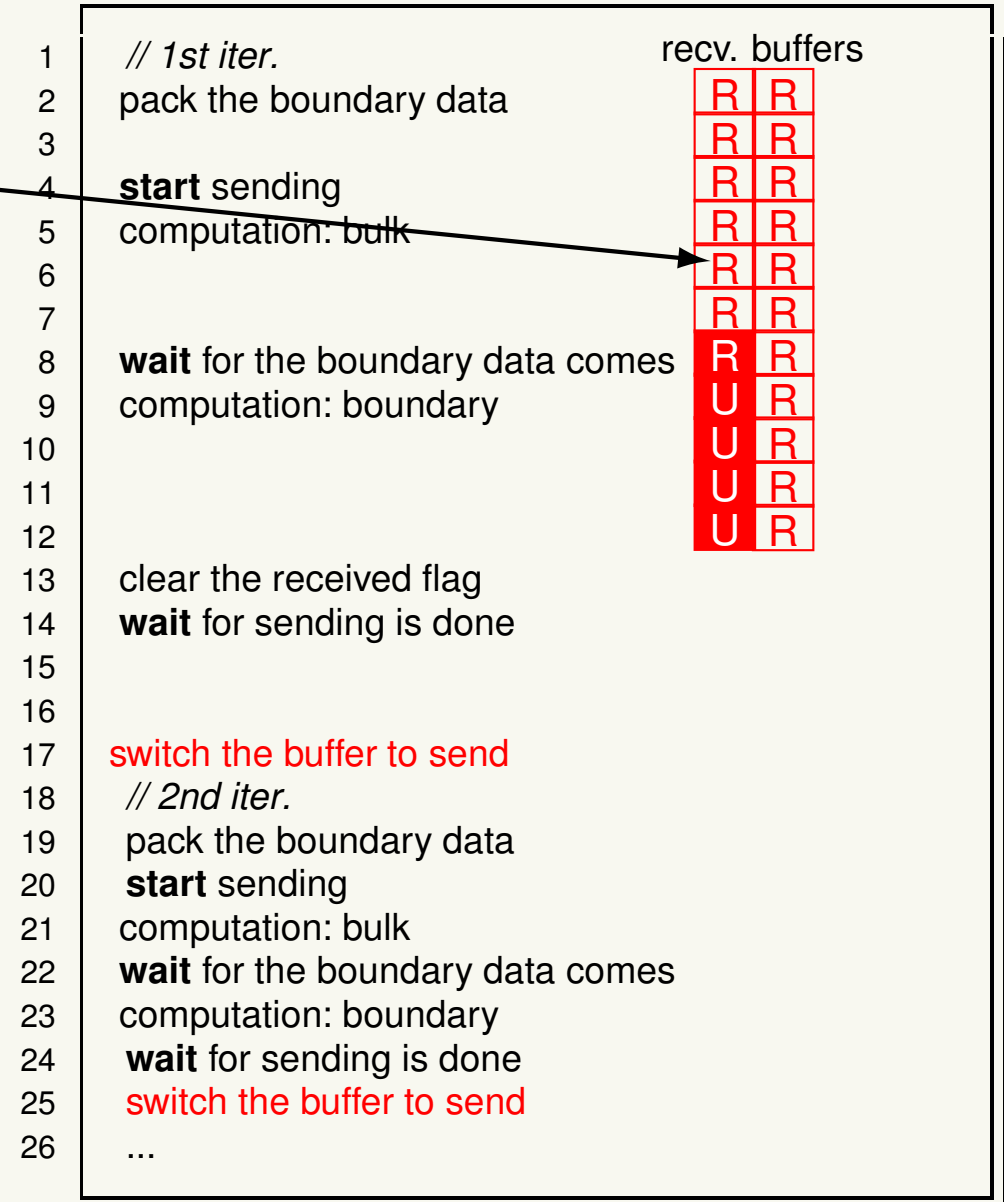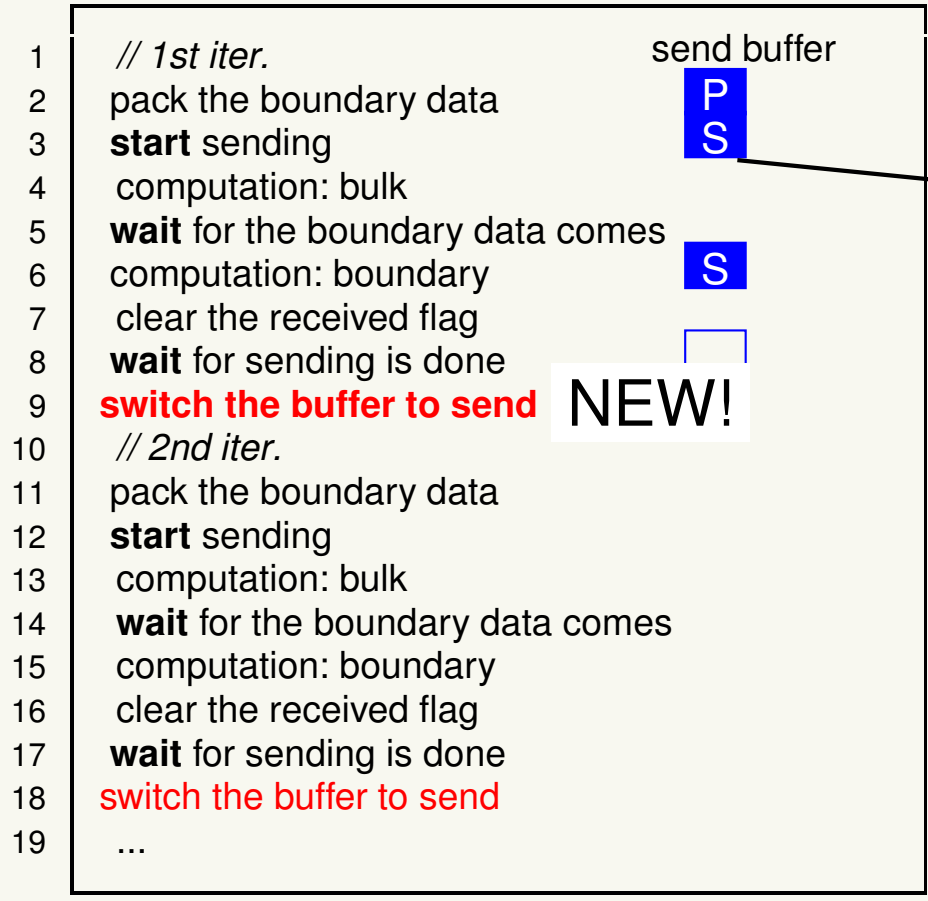
# Neighboring Communication: Double Buffering (RDMA)

| | | |
|---|---|---|
| 1 | // 1st iter. | send buffer |
| 2 | pack the boundary data | **P** |
| 3 | **start** sending | |
| 4 | computation: bulk | |
| 5 | **wait** for the boundary data comes | |
| 6 | computation: boundary | |
| 7 | clear the received flag | |
| 8 | **wait** for sending is done | |
| 9 | **switch the buffer to send** | |
| 10 | // 2nd iter. | |
| 11 | pack the boundary data | |
| 12 | **start** sending | |
| 13 | computation: bulk | |
| 14 | **wait** for the boundary data comes | |
| 15 | computation: boundary | |
| 16 | clear the received flag | |
| 17 | **wait** for sending is done | |
| 18 | switch the buffer to send | |
| 19 | ... | |

| | | |
|---|---|---|
| 1 | // 1st iter. | recv. buffers |
| 2 | pack the boundary data | |
| 3 | | |
| 4 | **start** sending | |
| 5 | computation: bulk | |
| 6 | | |
| 7 | | |
| 8 | **wait** for the boundary data comes | |
| 9 | computation: boundary | |
| 10 | | |
| 11 | | |
| 12 | | |
| 13 | clear the received flag | |
| 14 | **wait** for sending is done | |
| 15 | | |
| 16 | | |
| 17 | switch the buffer to send | |
| 18 | // 2nd iter. | |
| 19 | pack the boundary data | |
| 20 | **start** sending | |
| 21 | computation: bulk | |
| 22 | **wait** for the boundary data comes | |
| 23 | computation: boundary | |
| 24 | **wait** for sending is done | |
| 25 | switch the buffer to send | |
| 26 | ... | |

# Neighboring Communication: Double Buffering (RDMA)

send buffer

P
S

| | |
|---|---|
| 1 | // 1st iter. |
| 2 | pack the boundary data |
| 3 | **start** sending |
| 4 | computation: bulk |
| 5 | **wait** for the boundary data comes |
| 6 | computation: boundary |
| 7 | clear the received flag |
| 8 | **wait** for sending is done |
| 9 | **switch the buffer to send** |
| 10 | // 2nd iter. |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | clear the received flag |
| 17 | **wait** for sending is done |
| 18 | switch the buffer to send |
| 19 | ... |

recv. buffers

| | |
|---|---|
| 1 | // 1st iter. |
| 2 | pack the boundary data |
| 3 | |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | |
| 7 | |
| 8 | **wait** for the boundary data comes |
| 9 | computation: boundary |
| 10 | |
| 11 | |
| 12 | |
| 13 | clear the received flag |
| 14 | **wait** for sending is done |
| 15 | |
| 16 | |
| 17 | switch the buffer to send |
| 18 | // 2nd iter. |
| 19 | pack the boundary data |
| 20 | **start** sending |
| 21 | computation: bulk |
| 22 | **wait** for the boundary data comes |
| 23 | computation: boundary |
| 24 | **wait** for sending is done |
| 25 | switch the buffer to send |
| 26 | ... |

# Neighboring Communication: Double Buffering (RDMA)

send buffer

| 1 | // 1st iter. |
|---|---|
| 2 | pack the boundary data |
| 3 | **start** sending |
| 4 | computation: bulk |
| 5 | **wait** for the boundary data comes |
| 6 | computation: boundary |
| 7 | clear the received flag |
| 8 | **wait** for sending is done |
| 9 | **switch the buffer to send** |
| 10 | // 2nd iter. |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | clear the received flag |
| 17 | **wait** for sending is done |
| 18 | switch the buffer to send |
| 19 | ... |

recv. buffers

| 1 | // 1st iter. |
|---|---|
| 2 | pack the boundary data |
| 3 | |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | |
| 7 | |
| 8 | **wait** for the boundary data comes |
| 9 | computation: boundary |
| 10 | |
| 11 | |
| 12 | |
| 13 | clear the received flag |
| 14 | **wait** for sending is done |
| 15 | |
| 16 | |
| 17 | switch the buffer to send |
| 18 | // 2nd iter. |
| 19 | pack the boundary data |
| 20 | **start** sending |
| 21 | computation: bulk |
| 22 | **wait** for the boundary data comes |
| 23 | computation: boundary |
| 24 | **wait** for sending is done |
| 25 | switch the buffer to send |
| 26 | ... |

# Neighboring Communication: Double Buffering (RDMA)

|    |                                   | send buffer |
|----|-----------------------------------|-------------|
| 1  | // 1st iter.                      |             |
| 2  | pack the boundary data            | **P**       |
| 3  | **start** sending                 | **S**       |
| 4  | computation: bulk                 |             |
| 5  | **wait** for the boundary data comes |          |
| 6  | computation: boundary             |             |
| 7  | clear the received flag           |             |
| 8  | **wait** for sending is done      |             |
| 9  | **switch the buffer to send**     |             |
| 10 | // 2nd iter.                      |             |
| 11 | pack the boundary data            |             |
| 12 | **start** sending                 |             |
| 13 | computation: bulk                 |             |
| 14 | **wait** for the boundary data comes |          |
| 15 | computation: boundary             |             |
| 16 | clear the received flag           |             |
| 17 | **wait** for sending is done      |             |
| 18 | switch the buffer to send         |             |
| 19 | ...                               |             |

|    |                                      | recv. buffers |
|----|--------------------------------------|---------------|
| 1  | // 1st iter.                         | R R           |
| 2  | pack the boundary data               | R R           |
| 3  |                                      | R R           |
| 4  | **start** sending                    | R R           |
| 5  | computation: bulk                    | R R           |
| 6  |                                      | R R           |
| 7  |                                      | R R           |
| 8  | **wait** for the boundary data comes | R R           |
| 9  | computation: boundary                | U R           |
| 10 |                                      | U R           |
| 11 |                                      | U R           |
| 12 |                                      | U R           |
| 13 | clear the received flag              |               |
| 14 | **wait** for sending is done         |               |
| 15 |                                      |               |
| 16 |                                      |               |
| 17 | switch the buffer to send            |               |
| 18 | // 2nd iter.                         |               |
| 19 | pack the boundary data               |               |
| 20 | **start** sending                    |               |
| 21 | computation: bulk                    |               |
| 22 | **wait** for the boundary data comes |               |
| 23 | computation: boundary                |               |
| 24 | **wait** for sending is done         |               |
| 25 | switch the buffer to send            |               |
| 26 | ...                                  |               |

# Neighboring Communication: Double Buffering (RDMA)

send buffer

| 1 | // 1st iter. |
|---|---|
| 2 | pack the boundary data |
| 3 | **start** sending |
| 4 | computation: bulk |
| 5 | **wait** for the boundary data comes |
| 6 | computation: boundary |
| 7 | clear the received flag |
| 8 | **wait** for sending is done |
| 9 | **switch the buffer to send** NEW! |
| 10 | // 2nd iter. |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | clear the received flag |
| 17 | **wait** for sending is done |
| 18 | switch the buffer to send |
| 19 | ... |

P S

S

recv. buffers

| 1 | // 1st iter. |
|---|---|
| 2 | pack the boundary data |
| 3 | |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | |
| 7 | |
| 8 | **wait** for the boundary data comes |
| 9 | computation: boundary |
| 10 | |
| 11 | |
| 12 | |
| 13 | clear the received flag |
| 14 | **wait** for sending is done |
| 15 | |
| 16 | |
| 17 | switch the buffer to send |
| 18 | // 2nd iter. |
| 19 | pack the boundary data |
| 20 | **start** sending |
| 21 | computation: bulk |
| 22 | **wait** for the boundary data comes |
| 23 | computation: boundary |
| 24 | **wait** for sending is done |
| 25 | switch the buffer to send |
| 26 | ... |

# Neighboring Communication: Double Buffering (RDMA)

send buffer

| 1 | // 1st iter. |
|---|---|
| 2 | pack the boundary data |
| 3 | **start** sending |
| 4 | computation: bulk |
| 5 | **wait** for the boundary data comes |
| 6 | computation: boundary |
| 7 | clear the received flag |
| 8 | **wait** for sending is done |
| 9 | **switch the buffer to send** |
| 10 | // 2nd iter. |
| 11 | pack the boundary data |
| 12 | **start** sending |
| 13 | computation: bulk |
| 14 | **wait** for the boundary data comes |
| 15 | computation: boundary |
| 16 | clear the received flag |
| 17 | **wait** for sending is done |
| 18 | switch the buffer to send |
| 19 | ... |

NEW!

recv. buffers

| 1 | // 1st iter. |
|---|---|
| 2 | pack the boundary data |
| 3 | |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | |
| 7 | |
| 8 | **wait** for the boundary data comes |
| 9 | computation: boundary |
| 10 | |
| 11 | |
| 12 | |
| 13 | clear the receiv |
| 14 | **wait** for sendin |
| 15 | |
| 16 | |
| 17 | switch the buffer to send |
| 18 | // 2nd iter. |
| 19 | pack the boundary data |
| 20 | **start** sending |
| 21 | computation: bulk |
| 22 | **wait** for the boundary data comes |
| 23 | computation: boundary |
| 24 | **wait** for sending is done |
| 25 | switch the buffer to send |
| 26 | ... |

sending to the 2nd buf.

# Neighboring Communication: Double Buffering (RDMA)

send buffer

| # | Code | |
|---|------|---|
| 1 | // 1st iter. | |
| 2 | pack the boundary data | P S |
| 3 | **start** sending | |
| 4 | computation: bulk | |
| 5 | **wait** for the boundary data comes | S |
| 6 | computation: boundary | |
| 7 | clear the received flag | |
| 8 | **wait** for sending is done | |
| 9 | **switch the buffer to send** | |
| 10 | // 2nd iter. | |
| 11 | pack the boundary data | P S |
| 12 | **start** sending | |
| 13 | computation: bulk | |
| 14 | **wait** for the boundary data comes | |
| 15 | computation: boundary | |
| 16 | clear the received flag | |
| 17 | **wait** for sending is done | S |
| 18 | switch the buffer to send | |
| 19 | ... | |

NEW!

recv. buffers

| # | Code |
|---|------|
| 1 | // 1st iter. |
| 2 | pack the boundary data |
| 3 | |
| 4 | **start** sending |
| 5 | computation: bulk |
| 6 | |
| 7 | |
| 8 | **wait** for the boundary data comes |
| 9 | computation: boundary |
| 10 | |
| 11 | |
| 12 | |
| 13 | clear the receiv... |
| 14 | **wait** for sendin... |
| 15 | |
| 16 | |
| 17 | switch the buffer to send |
| 18 | // 2nd iter. |
| 19 | pack the boundary data |
| 20 | **start** sending |
| 21 | computation: bulk |
| 22 | **wait** for the boundary data comes |
| 23 | computation: boundary |
| 24 | **wait** for sending is done |
| 25 | switch the buffer to send |
| 26 | ... |

sending to the 2nd buf.

# Remote Direct Memory Access (RDMA)

we use "put" in sending
 it directly writes to the memory on the remote process

- send: put (directly memory on the remote process)

 boundary data + watchdog flag
- Wait (recv.): check the flag is updated

 after the boundary computation, the flag is reset

# Remote Direct Memory Access (RDMA)

we use "put" in sending
    it directly writes to the memory on the remote process

- send: put (directly memory on the remote process)
                                               boundary data + watchdog flag
- Wait (recv.): check the flag is updated
                          after the boundary computation, the flag is reset

NOTE 1: 2 buffers are enough: "sending proc." also receives data from "receiving proc." $\Rightarrow$ automatic synchronization

# Remote Direct Memory Access (RDMA)

we use "put" in sending
    it directly writes to the memory on the remote process

- send: put (directly memory on the remote process)
                                                    boundary data + watchdog flag
- Wait (recv.): check the flag is updated
                                    after the boundary computation, the flag is reset

NOTE 1: 2 buffers are enough: "sending proc." also receives data
from "receiving proc." $\Rightarrow$ automatic synchronization

NOTE 2: one can alternatively use MPI (persistent) communication
to implement double buffering.

# Benchmark

# Test detail: Jacobi method for 2-dim system

target system: $Mx = b$ with

$$(Mx)(i,j) =$$
$$\underbrace{(4 + m^2)x(i,j)}_{\equiv Dx} \underbrace{-x(i+1,j) - x(i-1,j) - x(i,j+1) - x(i,j-1)}_{\equiv Hx}$$

$$\xrightarrow{\text{cont. limit}} (m^2 - \partial^2)x$$

Jacobi method

$$x^{(k)} \to x^{(k+1)} = D^{-1}(b - Hx^{(k)})$$

Only the hopping $H$ contains the communication

- fixed number of iterations: 10
- local lattice size: $60 \times 60$
- communication buffer: needed size + dummy (+ flag)

$$s = \frac{\text{needed} + \text{dummy}}{\text{needed}}, 1 \leq s \leq 8192$$

# Hopping (Mult of $H$)

1. packing the boundary data
2. start sending/receiving the boundary data
3. calculate: internal area
4. wait for receiving
5. calculate: boundary area
6. wait for sending finished

1. packing the boundary data
2. start sending/receiving the boundary data
3. calculate: internal area
4. wait for receiving
5. calculate: boundary area
6. wait for sending finished

comm.

# Hopping (Mult of $H$)

1. packing the boundary data
2. start sending/receiving the boundary data
3. calculate: internal area        overlap
4. wait for receiving                    comm.
5. calculate: boundary area
6. wait for sending finished

# Hopping (Mult of $H$)

1. packing the boundary data
2. start sending/receiving the boundary data
3. calculate: internal area    overlap
4. wait for receiving                comm.
5. calculate: boundary area
6. wait for sending finished    send wait

# Hopping (Mult of $H$)

1. packing the boundary data
2. start sending/receiving the boundary data
3. calculate: internal area ⎯⎯ overlap
4. wait for receiving ⎯⎯ comm.
5. calculate: boundary area
6. wait for sending finished ⎯⎯ send wait

non-overlap = comm. − overlap
= start sending/receiving + wait for receiving

# Estimated Bandwidth and uTofu Interface

## uTofu

- Low level interface to use Tofu Interconnect
- It allows to specify Tofu Network Interface (TNI) to use
  tuning with the optimal TNI assignment for QCD
- 6 TNI/node, 6.8GB/s for each TNI

## Bandwidth estimate

- 1 node with 4 MPI proc.

$$\Rightarrow (4 \text{ directions}) \times (4 \text{ ranks}) = 16 \text{ comm.}$$
$$(\text{each comm. has the same size})$$

- Using 4 TNI: each TNI is used 4 times
  $$6.8 \times 4 = 27.2 \text{ GB/s}$$
- Using 6 TNI: each TNI is used 2 or 3 times
  $$6.8 \times 6 \times \frac{16}{18} = 36.3 \text{ GB/s}$$

16 comm.

1st.　　　2nd.　　　3rd.

# Performance: mult of $H$ on A64fx 1 node

NOTE: result on the evaluation environment, it does not guarantee the performance on the actual Fugaku
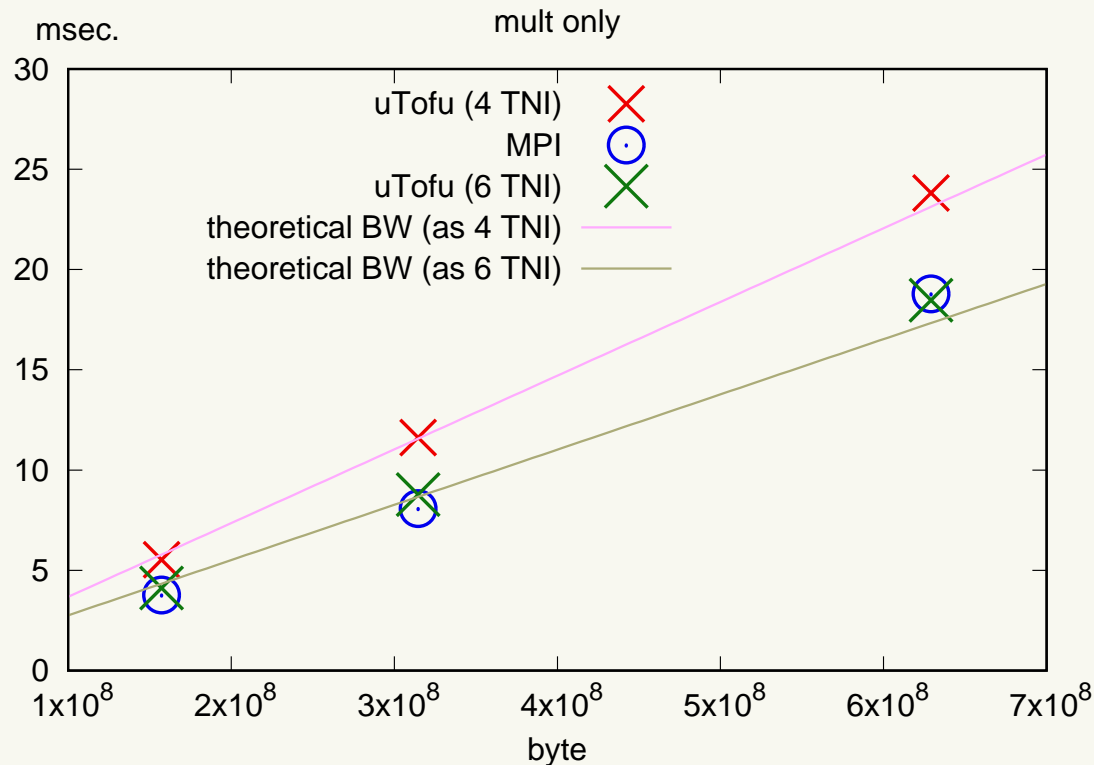


- comm. overlap $\simeq$ 0.61 msec. (replaced with a line in the plots )
- both show good scaling for large communication data size
- uTofu interface has a smaller overhead
- TNI for uTofu — 0:$+x$, 1:$-x$, 2:$+y$, 3:$-y$

# Performance: mult of $H$ on A64fx 1 node

NOTE: result on the evaluation environment, it does not guarantee the performance on the actual Fugaku



- comm. overlap $\simeq$ 0.61 msec. (replaced with a line in the plots )
- both show good scaling for large communication data size
- uTofu interface has a smaller overhead
- TNI for uTofu

NOTE: result on the evaluation environment, it does not guarantee the performance on the actual Fugaku



msec.

mult only

- uTofu (4 TNI) ✕
- MPI ⊙
- uTofu (6 TNI) ✕
- theoretical BW (as 4 TNI)
- theoretical BW (as 6 TNI)

large comm. size
$\Rightarrow$ time for mult
$\simeq$ time for comm.

saturation of the network bandwidth

- MPI: 32.4 GB/s
- uTofu (4 TNI): 25.8 GB/s          cf. $6.8 \times 4 = 27.2$ GB/s
- uTofu (6 TNI): 33.0 GB/s          cf. $6.8 \times 16/18 = 36.3$ GB/s

# Conclusions

## Conclusions

to accelerate neighboring communication, we have implemented double buffering algorithm

test with a simple 2-dim system

- using uTofu interface seems promising

Future: to do (or on going) for Fugaku

- implement double buffering+uTofu to QCD code

  qws: Nakamura-san's talk

- multi nodes, proper TNI settings,...
- official predicted performance for LQCD (vs. K-computer): x25+

  "+" will be how much????