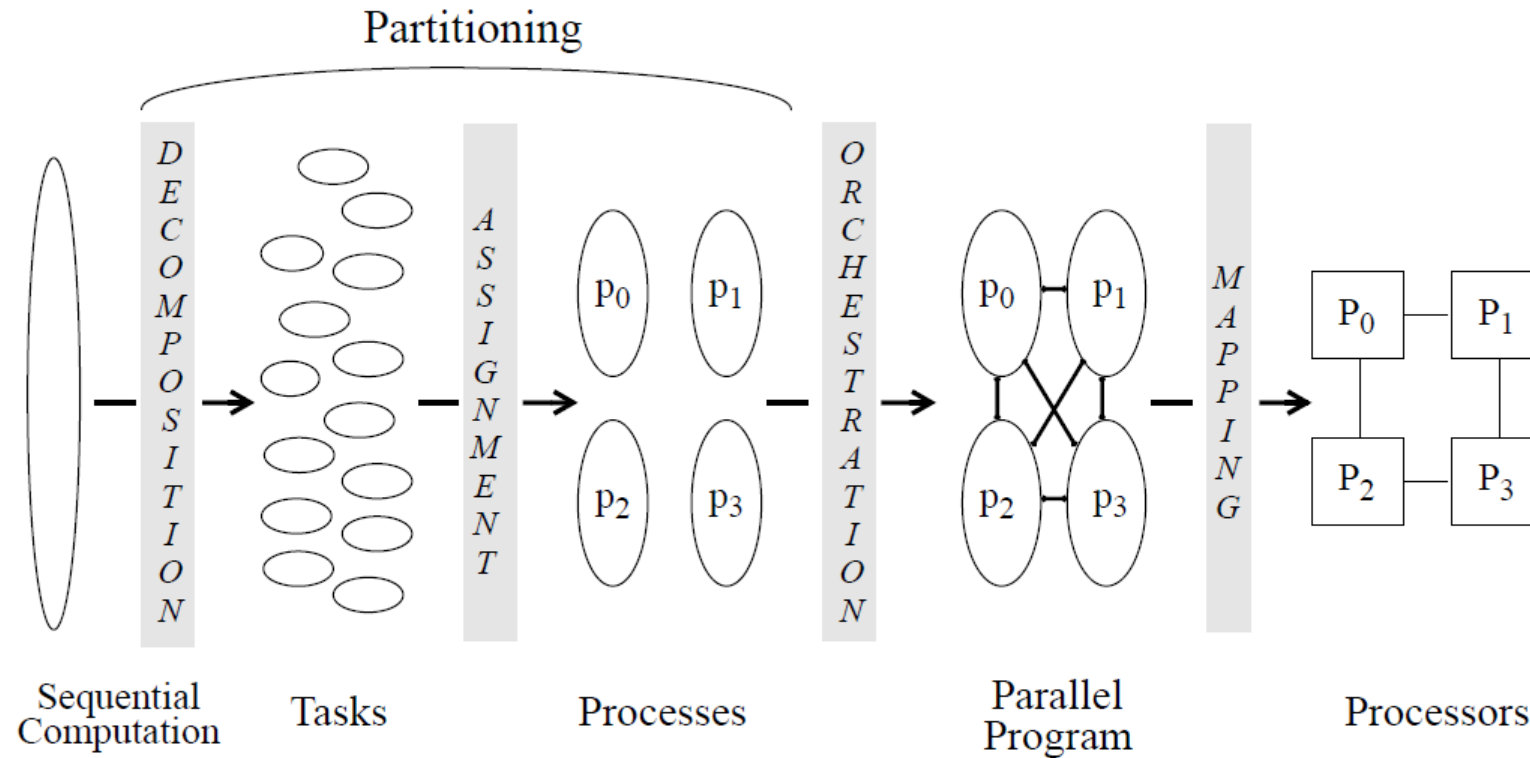


# Parallelization

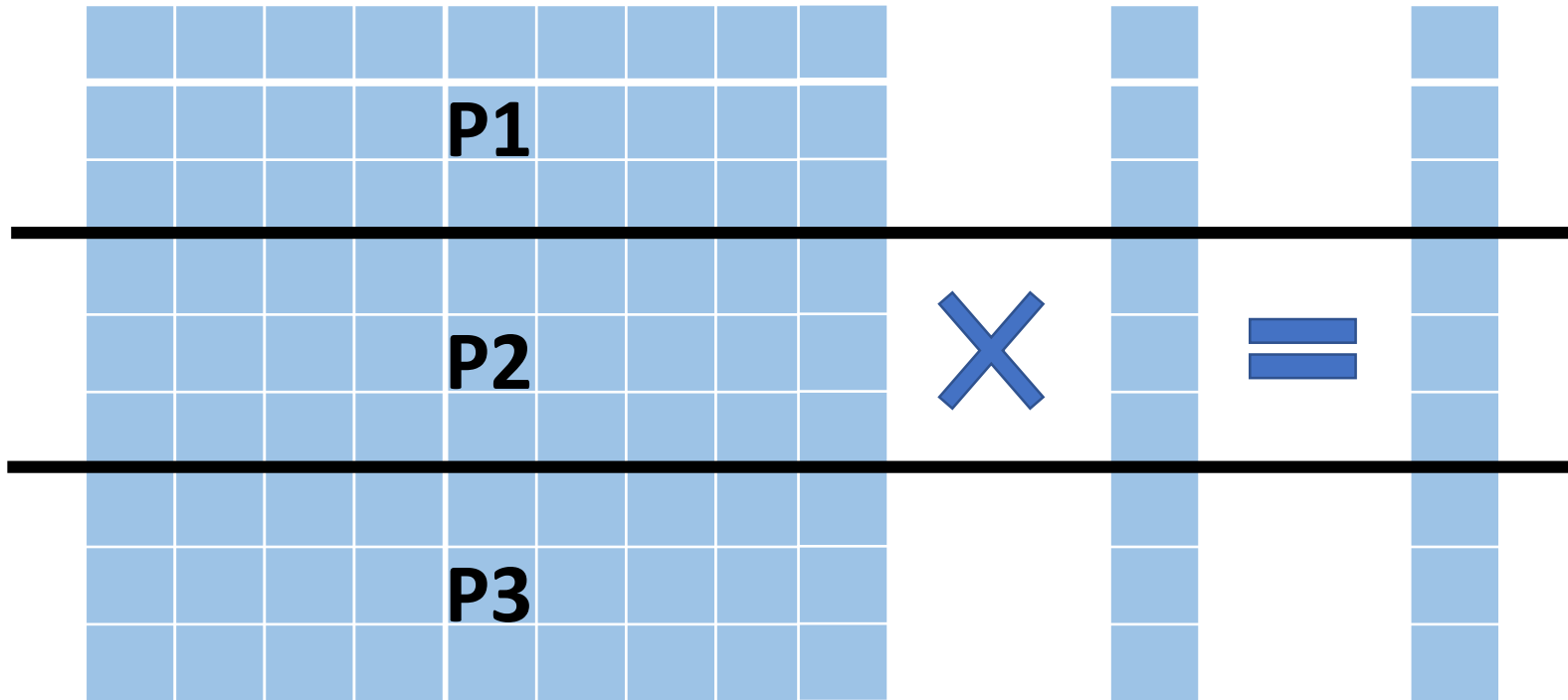
Jan 25, 2019

# Illustration of Parallelization Steps



Source: Culler et al.

# Parallelization – Matrix Vector Multiplication



$P = 3$

Decomposition

Assignment

Orchestration

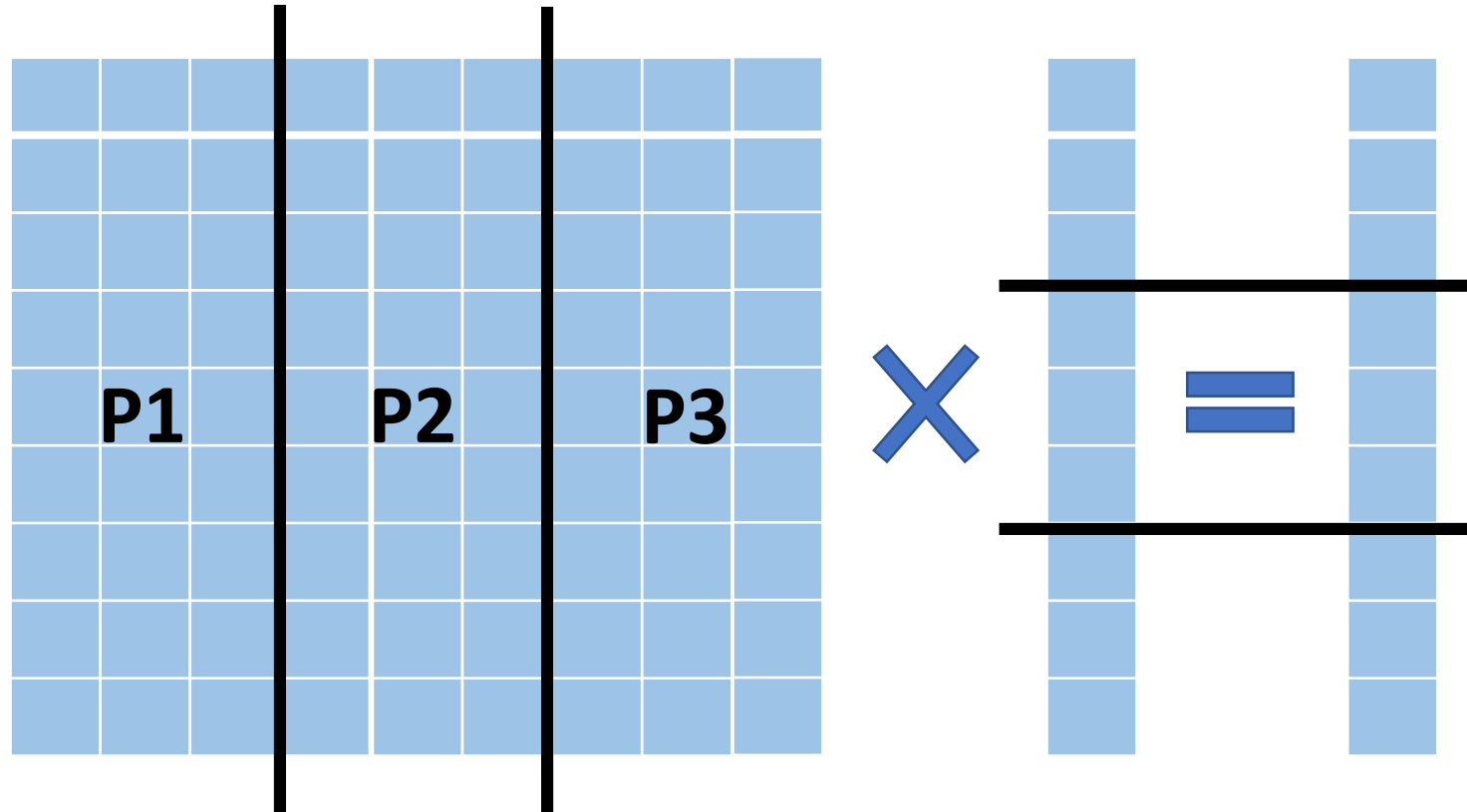
- Allgather

- Gather

Mapping

- What is the initial communication step ?
- Ways to distribute vector ?
- What are the differences between distribution and parallel reads?

# Parallelization – Matrix Vector Multiplication



$P = 3$

Decomposition  
Assignment  
Orchestration  
Mapping

What is the advantage of column-wise partitioning ?

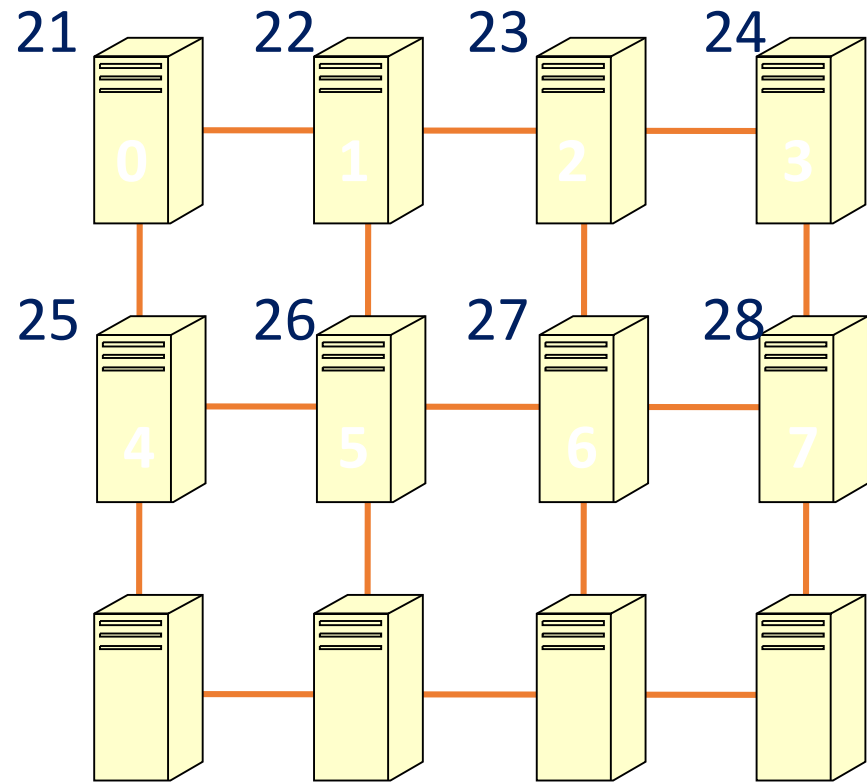
# Virtual Topology

<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>
<b>8</b>	<b>9</b>	<b>10</b>	<b>11</b>	<b>12</b>	<b>13</b>	<b>14</b>	<b>15</b>
<b>16</b>	<b>17</b>	<b>18</b>	<b>19</b>	<b>20</b>	<b>21</b>	<b>22</b>	<b>23</b>
<b>24</b>	<b>25</b>	<b>26</b>	<b>27</b>	<b>28</b>	<b>29</b>	<b>30</b>	<b>31</b>

**8 x 4 2D virtual process topology**

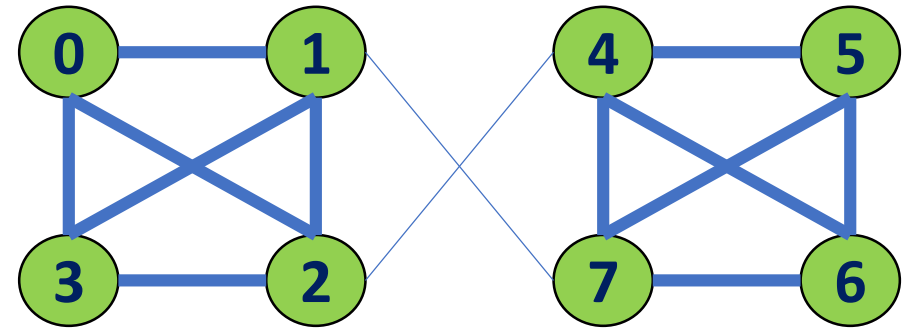
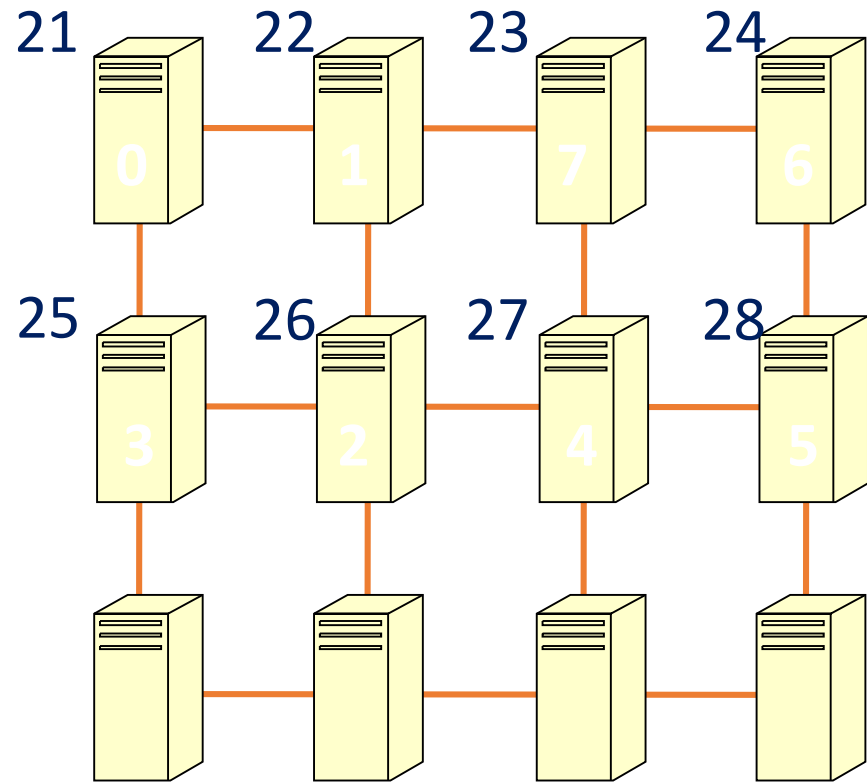
- Communication pattern of MPI processes
  - Graphical representation of communications
    - Nearest neighbor in a mesh
    - All-to-all
    - ...
  - Convenient way to represent communications
- Note: Virtual topology set up before execution

# Physical Topology



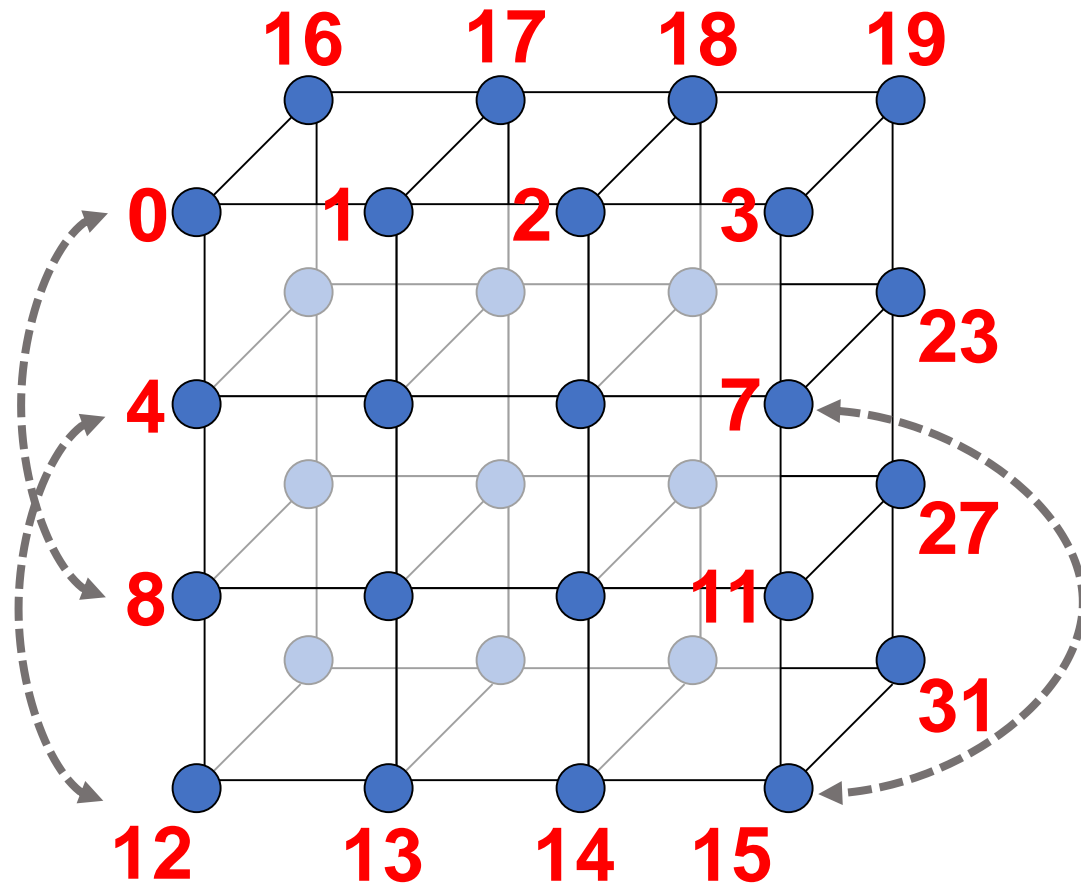
- Connections between allocated cores
- Default placement of ranks based on node IDs
- Mapping: Placement of ranks onto cores
- **Topology-aware mapping**: Mapping that minimizes all communication times taking into account the physical topology

# Rank placement

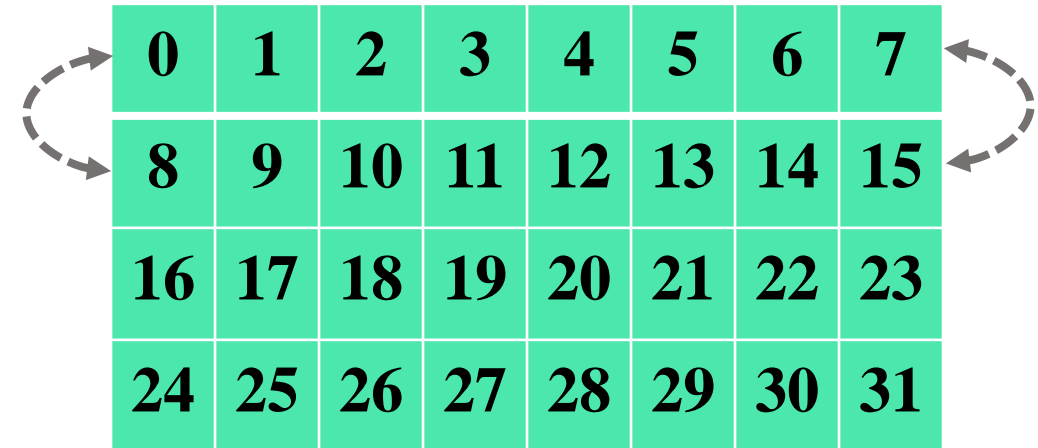


- May place ranks anywhere in the allocated nodes based on the communication pattern

# Process-to-processor Mapping



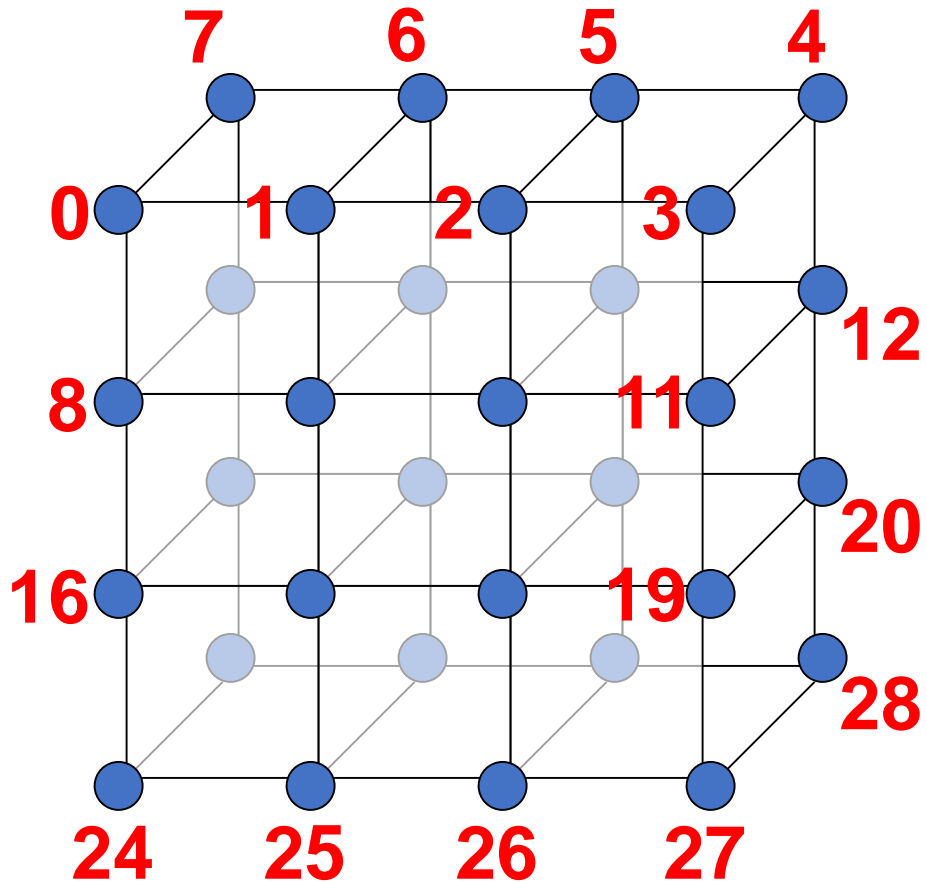
4 x 4 x 2 3D torus



8 x 4 2D virtual process topology



# Topology-aware Mapping

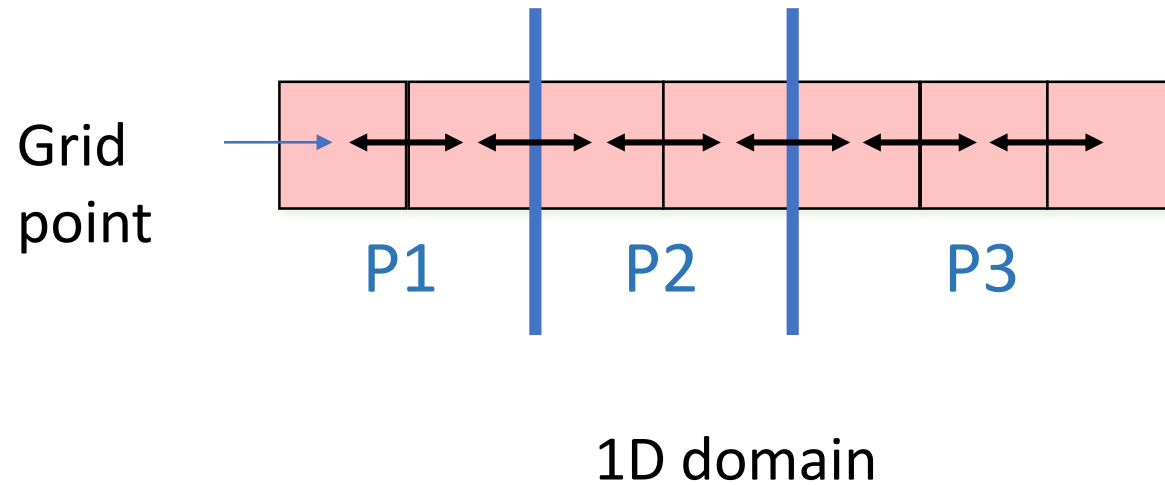


4 x 4 x 2 3D torus

0	1	2	3	4	5	6	7
8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23
24	25	26	27	28	29	30	31

8 x 4 2D virtual process topology

# 1D Domain decomposition

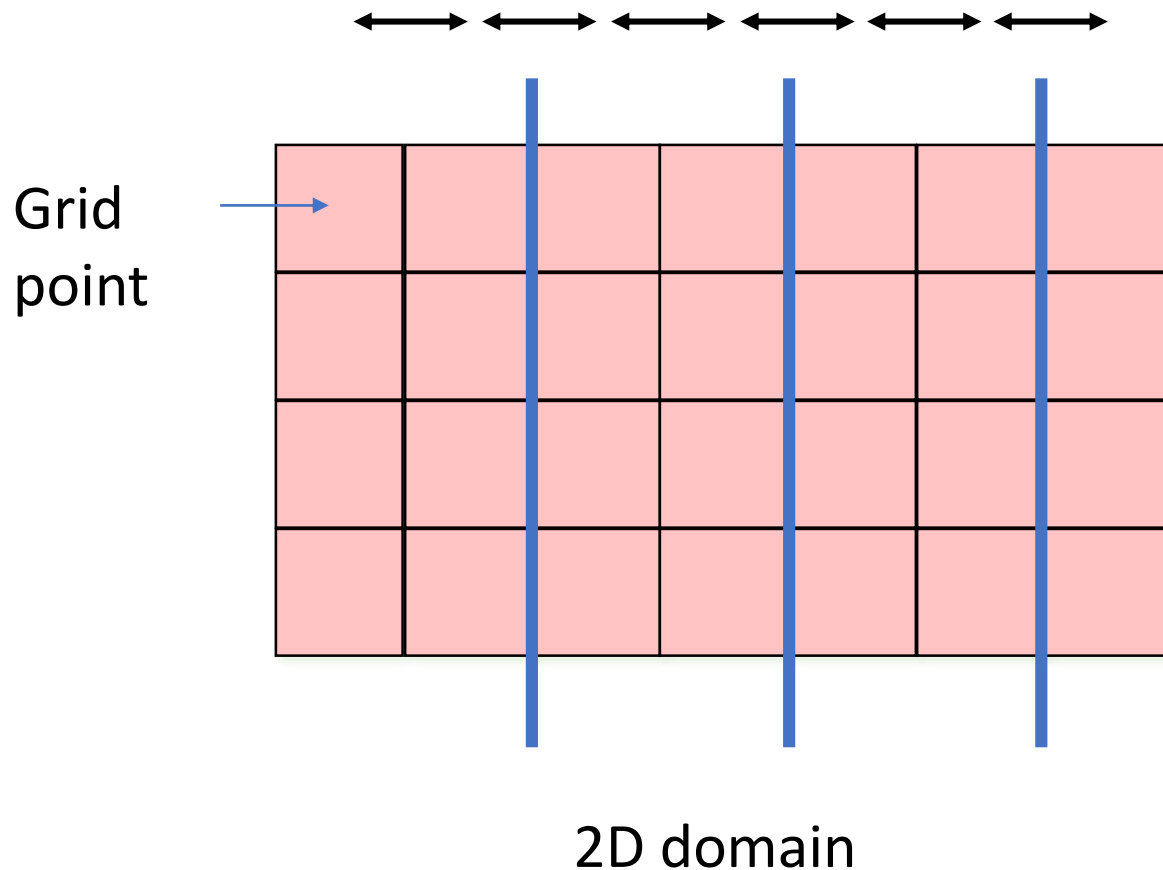


N grid points  
P processes  
 $N/P$  points per process

Communications?

2 sends()  
2 recvs()

# 1D Domain decomposition

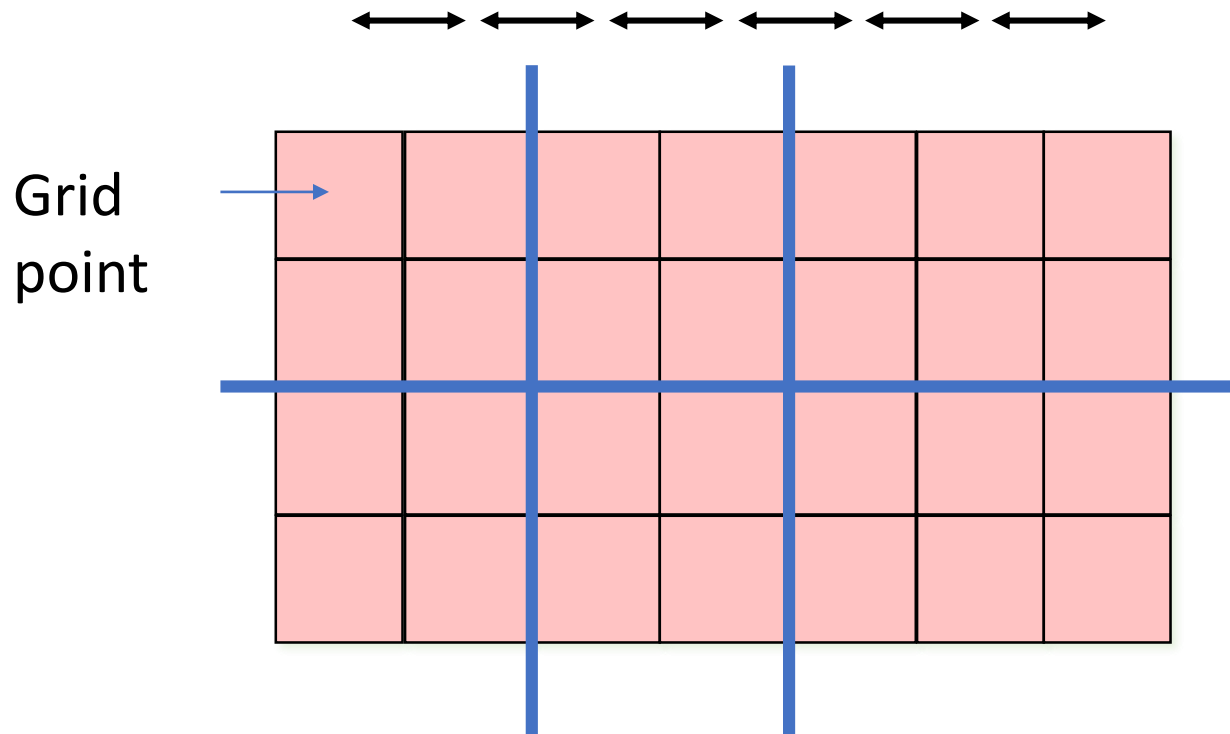


N grid points  
P processes  
 $N/P$  points per process

Decomposition  
Assignment  
Orchestration  
Mapping

Q1: Tunable parameters?  
Q2: One drawback?

# 2D Domain decomposition



N grid points  
P processes  
 $N/P$  points per process

- + Several parallel communications
- + Lower communication volume/process

# Send / Recv Options

0	1	2	3				
4	5	6	7				
8	9	10	11				

MPI\_Pack (buf)

MPI\_Recv (buf)

MPI\_Pack (buf)

MPI\_Unpack (buf)

MPI\_Pack (buf)

MPI\_Unpack (buf)

MPI\_Send (buf)

MPI\_Unpack (buf)

Create a new datatype

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

# MPI Derived Datatypes

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

0'	1'	2'	3'
----	----	----	----

MPI\_Type\_contiguous

Count = ?

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

0	1	2	3	4	5	6	7
---	---	---	---	---	---	---	---

MPI\_Type\_vector

count = 2, blocklength = 1, stride = 4

MPI\_Type\_vector (count, blocklength, stride, MPI\_INT, newtype)

# MPI Derived Datatypes

MPI\_Datatype newtype

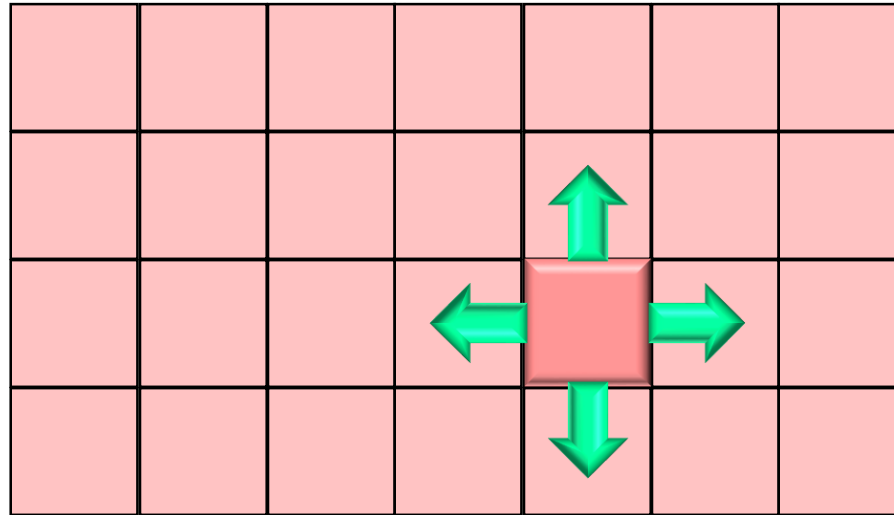
- MPI\_Type\_contiguous (count, oldtype, newtype)
- MPI\_Type\_vector (count, blocklength, stride, oldtype, newtype)
- MPI\_Type\_create\_subarray (ndims, array\_of\_sizes, array\_of\_subsizes, array\_of\_starts, order, oldtype, newtype)
- MPI\_Type\_create\_struct (count, array\_of\_blocklengths, array\_of\_displacements, array\_of\_types, newtype)

MPI\_Type\_commit (newtype)

.....

MPI\_Type\_free (newtype)

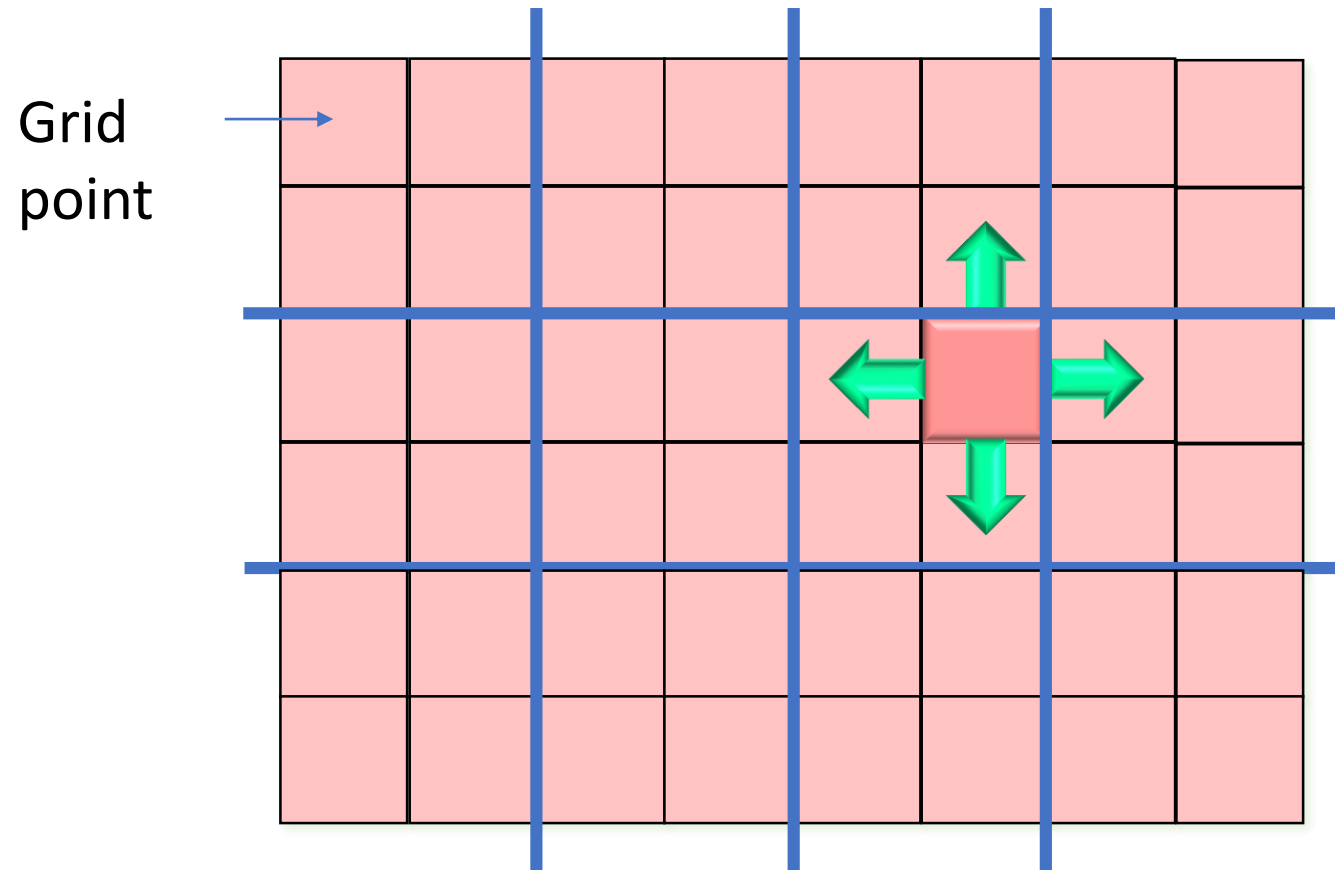
# Stencils



Five-point stencil: Each grid point new value is the average of itself and its four neighbors'



# 2D Domain decomposition



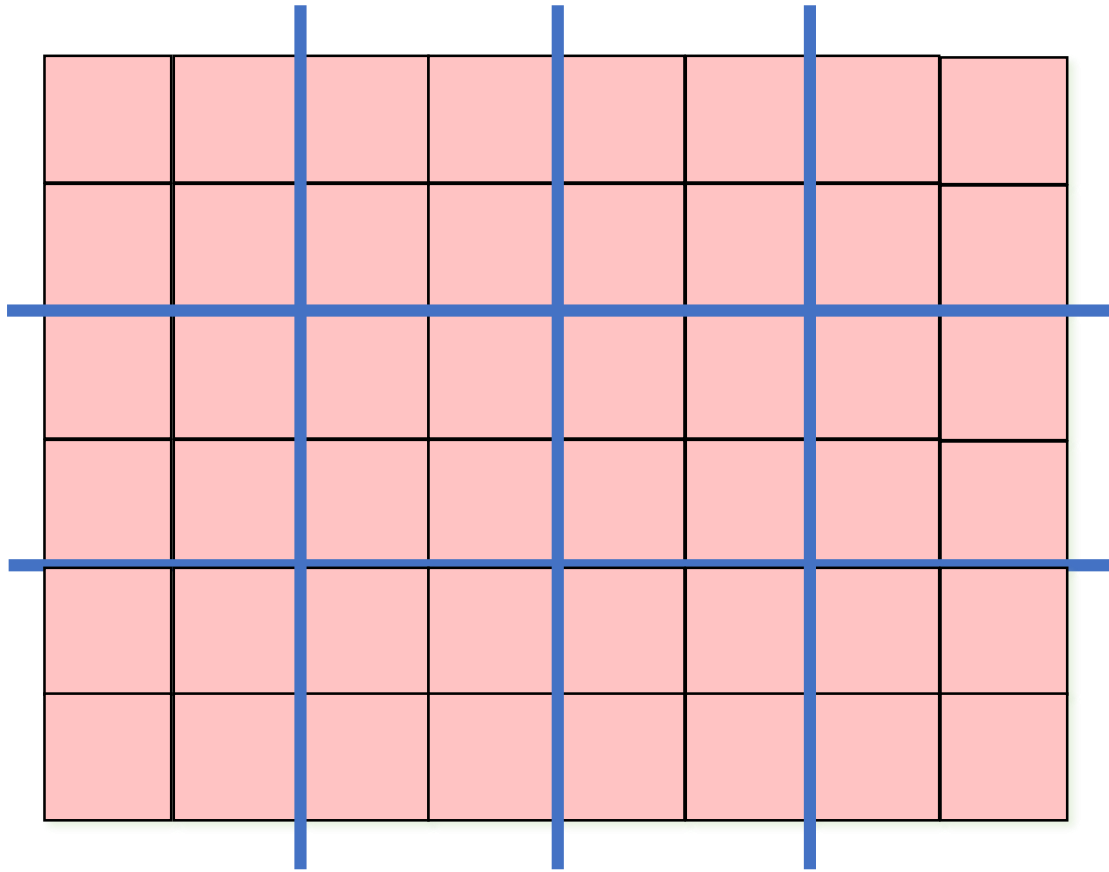
Communications?

4 Isends()  
4 Irecv()

Design considerations  
for multiple variables?

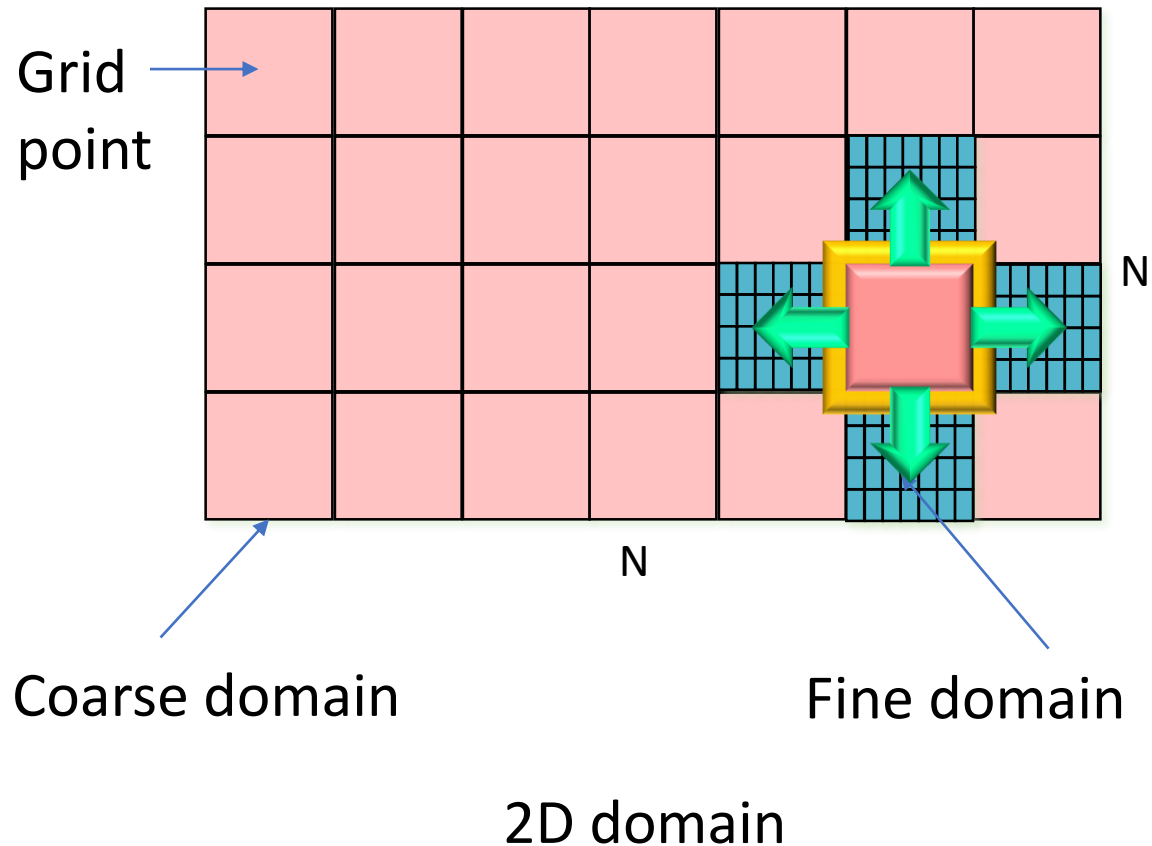
Pack  
Unpack

# 2D Domain decomposition



Does the order of processes matter in 2D?

# Domain refinement



Halo exchange

- Each cell has some ghost regions
- Communication with neighbors

*#Computations and communication volume of each cell?*