



# Parameterized Algorithms for Non-uniform All-to-all

Presenter: Jens Domke

Authors: Ke Fan, Jens Domke, Seydou Ba, and Sidharth Kumar



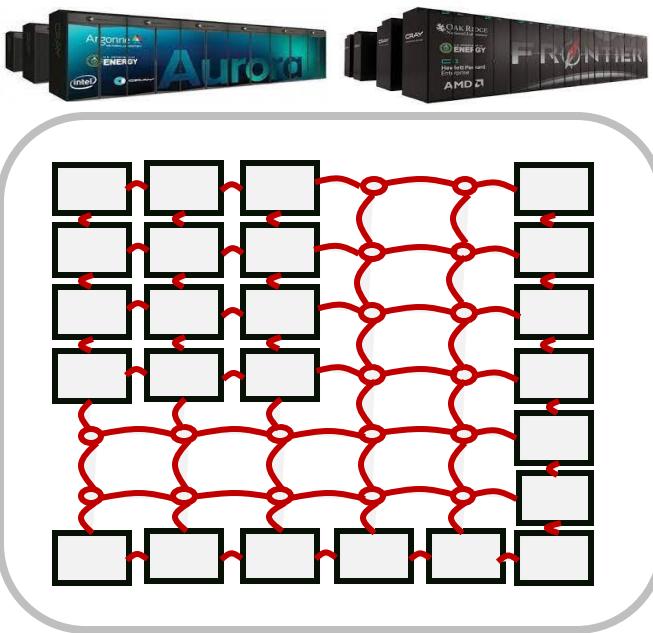
# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- Parameterized Logarithmic Algorithm
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- Parameterized Logarithmic Algorithm
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

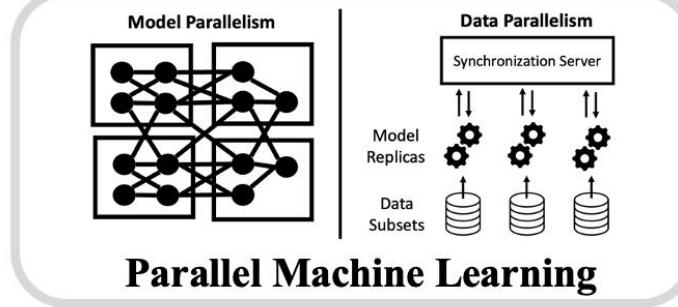
# Exascale Computing Necessaries Effective Parallelism



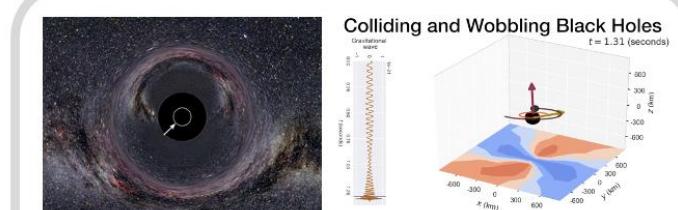
Exascale Supercomputers



Earth and Climate Modeling



Parallel Machine Learning



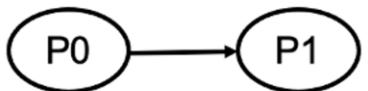
Black Holes Simulation



Inter-process Data  
Movement

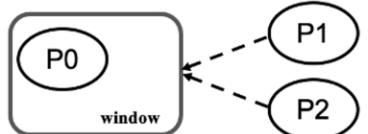
# Three Fundamental Interfaces of Data Movement

## Point-to-point



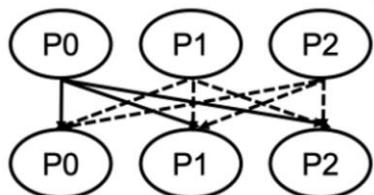
MPI\_Send  
MPI\_Recv

## One-sided

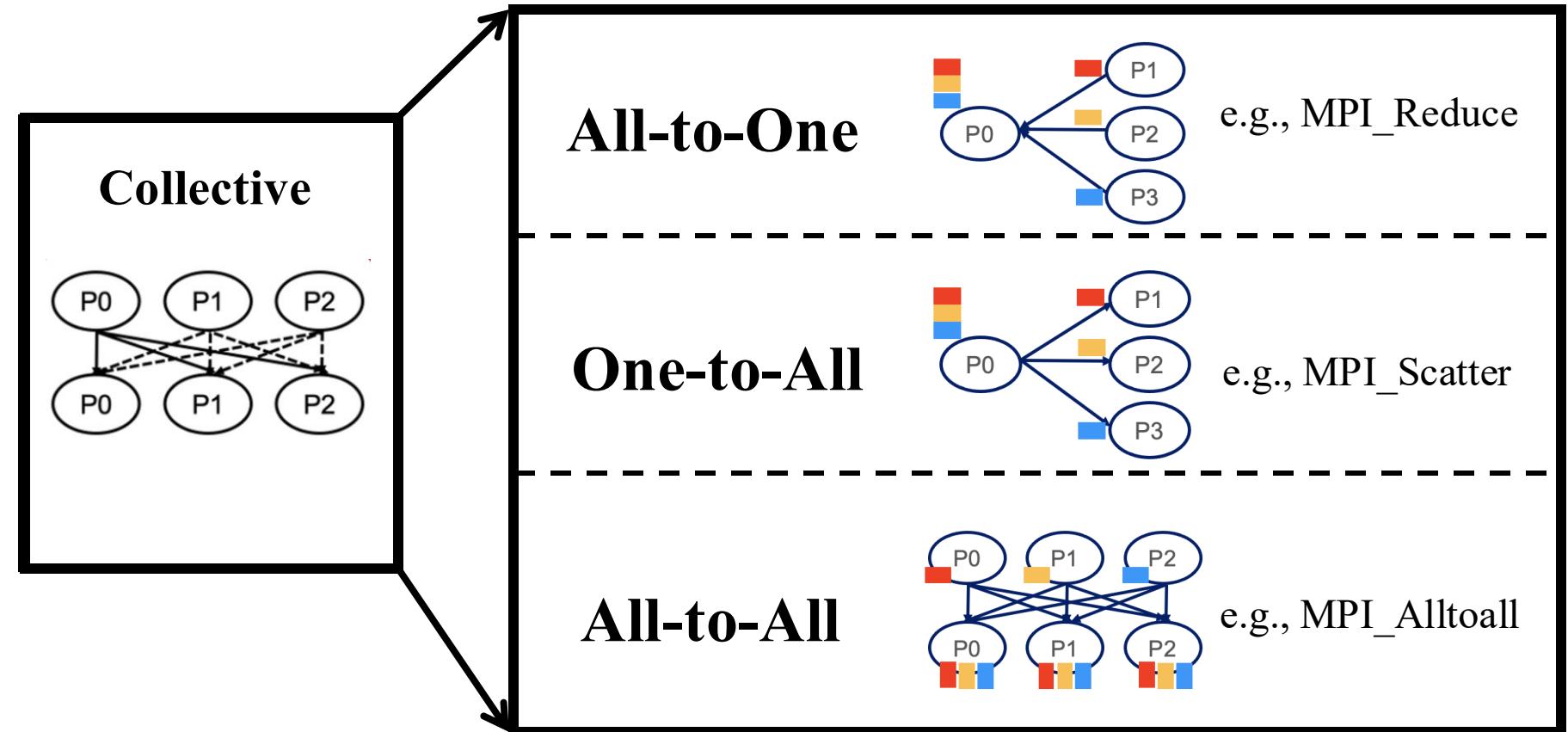
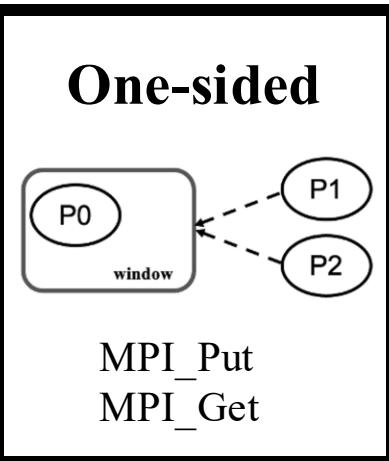
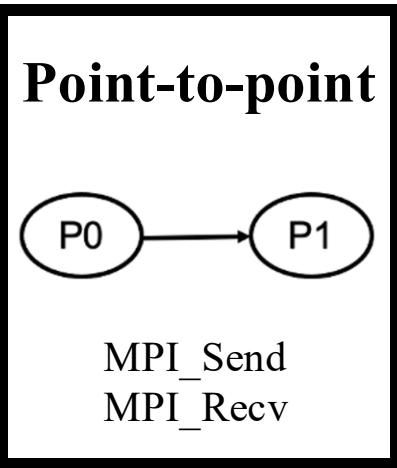


MPI\_Put  
MPI\_Get

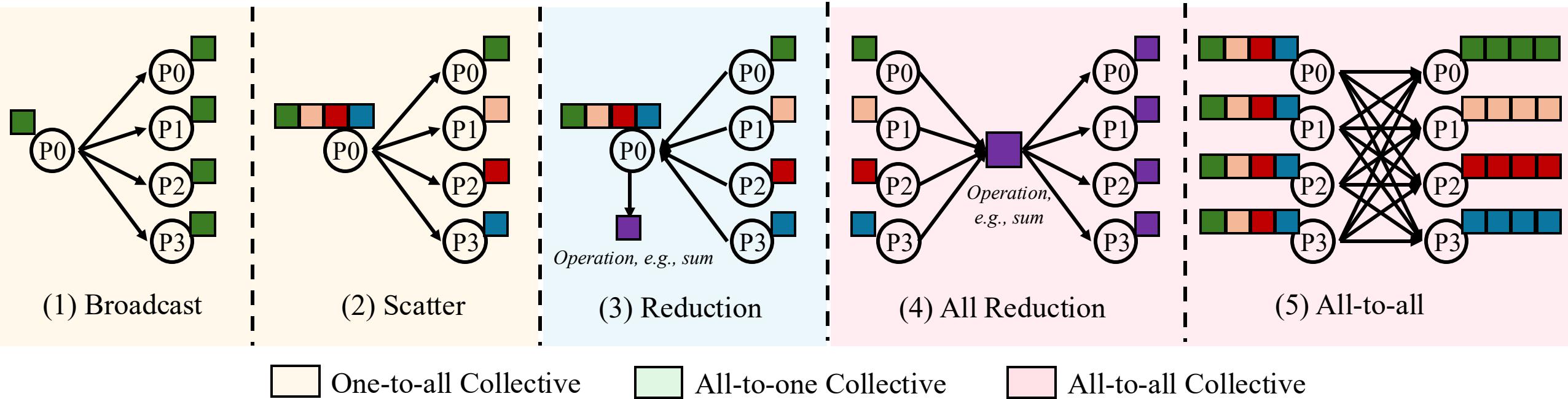
## Collective



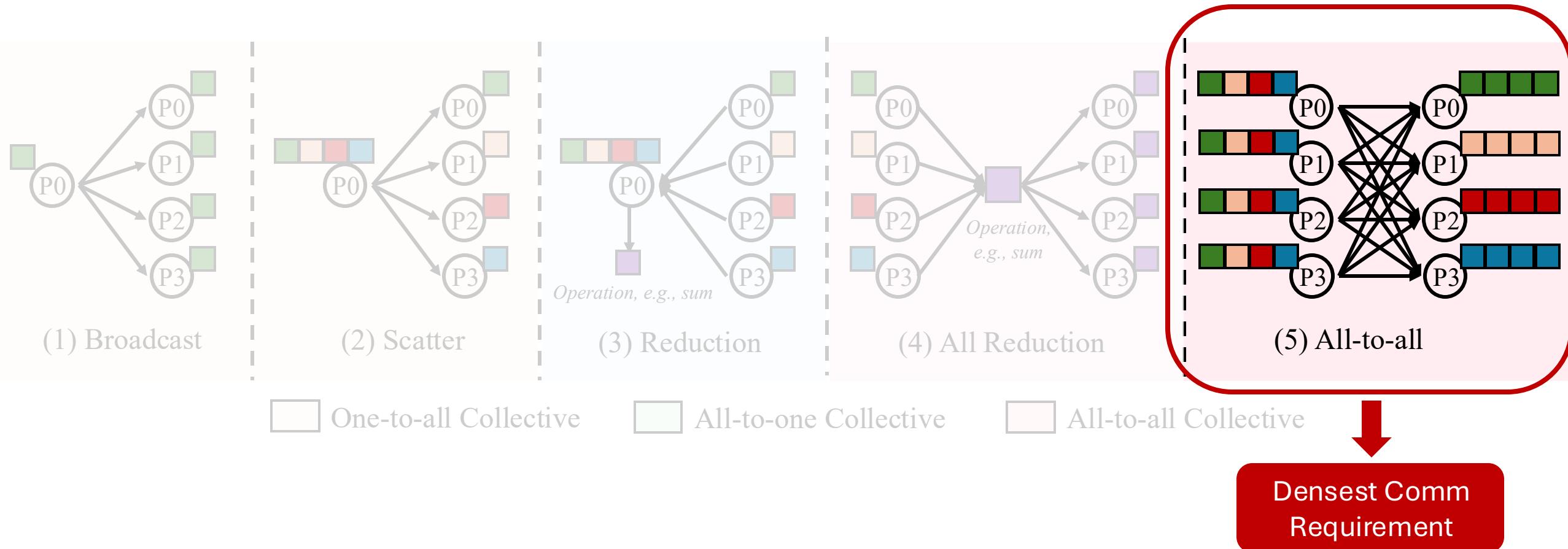
# Collective Involves Communication Among All Processes



# Some Widely Used Collective Primitives



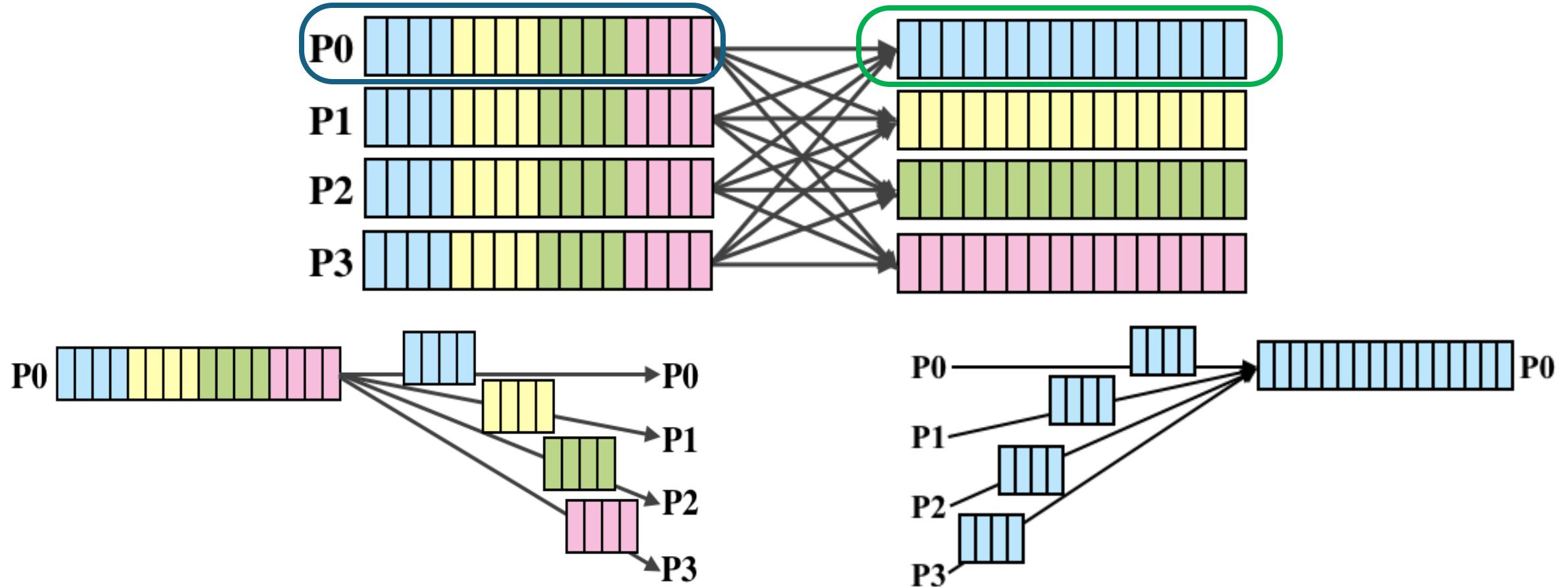
# All-to-all Data Exchange is the Most Challenging to Scale and Optimize



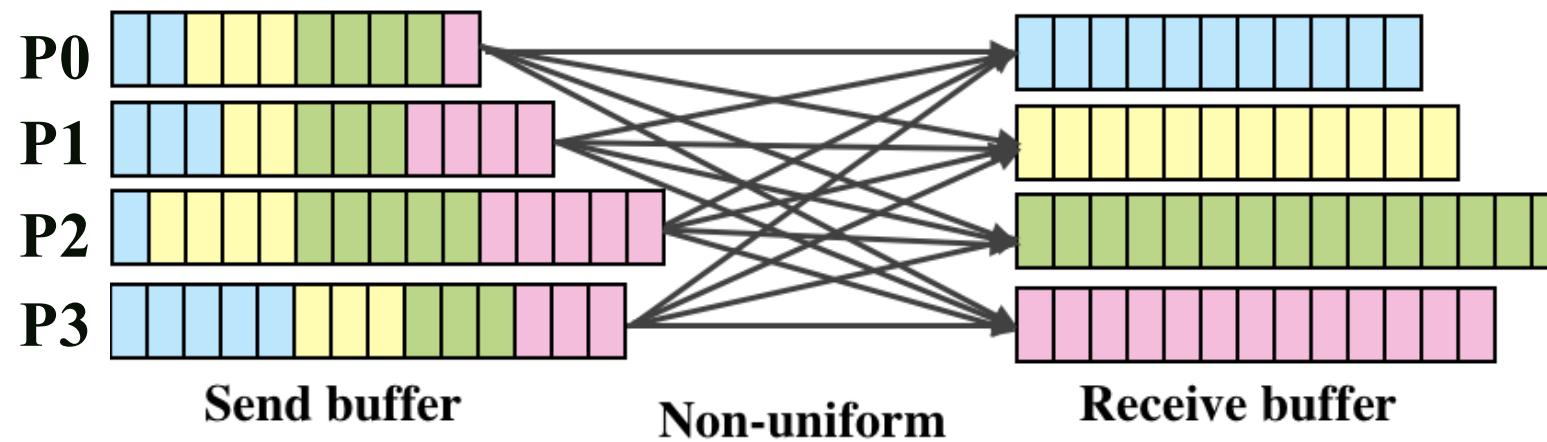
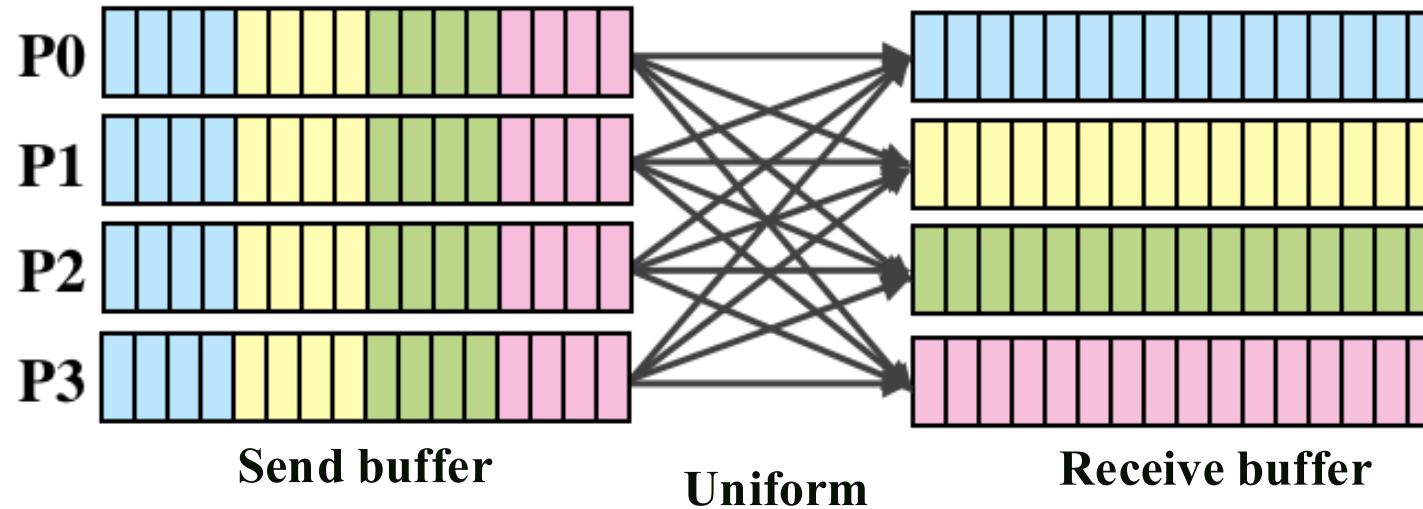
# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- Parameterized Logarithmic Algorithm
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

# All-to-all – Every process Sends and receives data from every other process



# Uniform and Non-uniform All-to-all



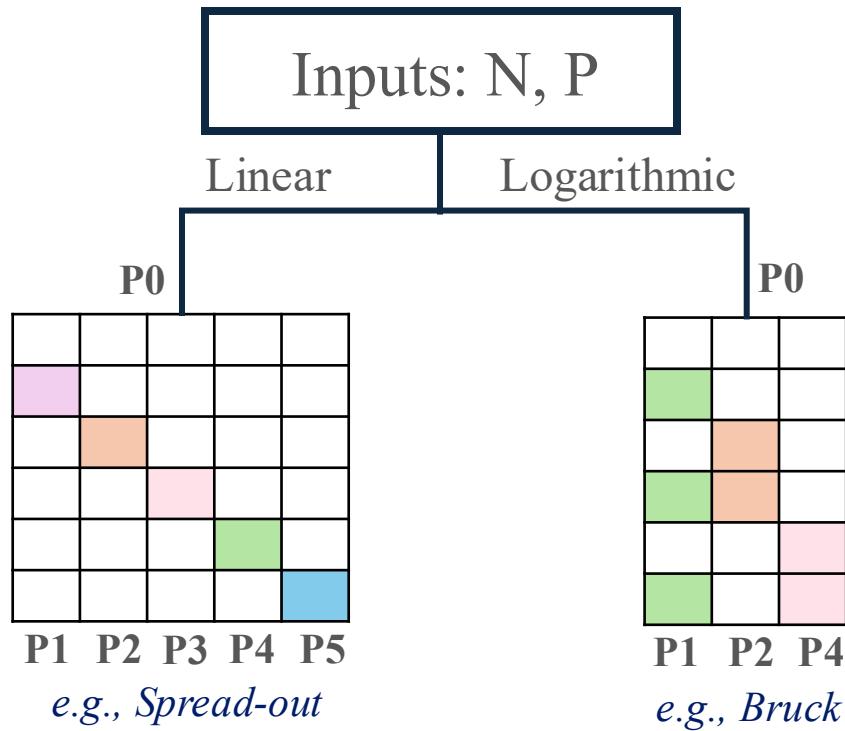
# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- Parameterized Logarithmic Algorithm
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

# Standard MPI All-to-all Implementations

*N: Size of data-block*

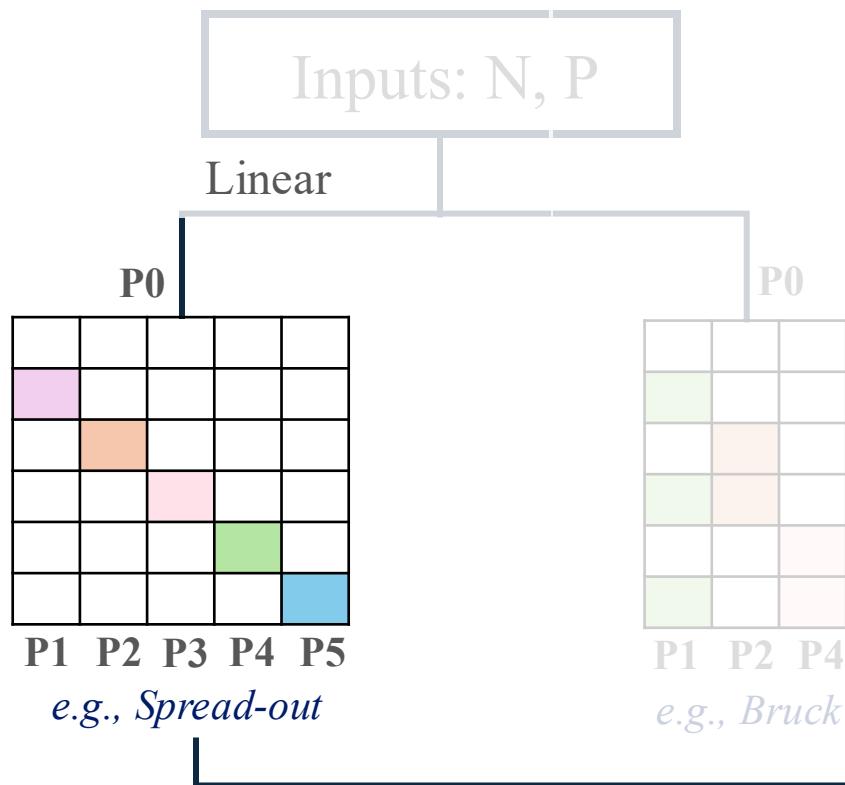
*P: Process count*



# Standard MPI All-to-all Implementations: Spread-out

*N: Size of data-block*

*P: Process count*



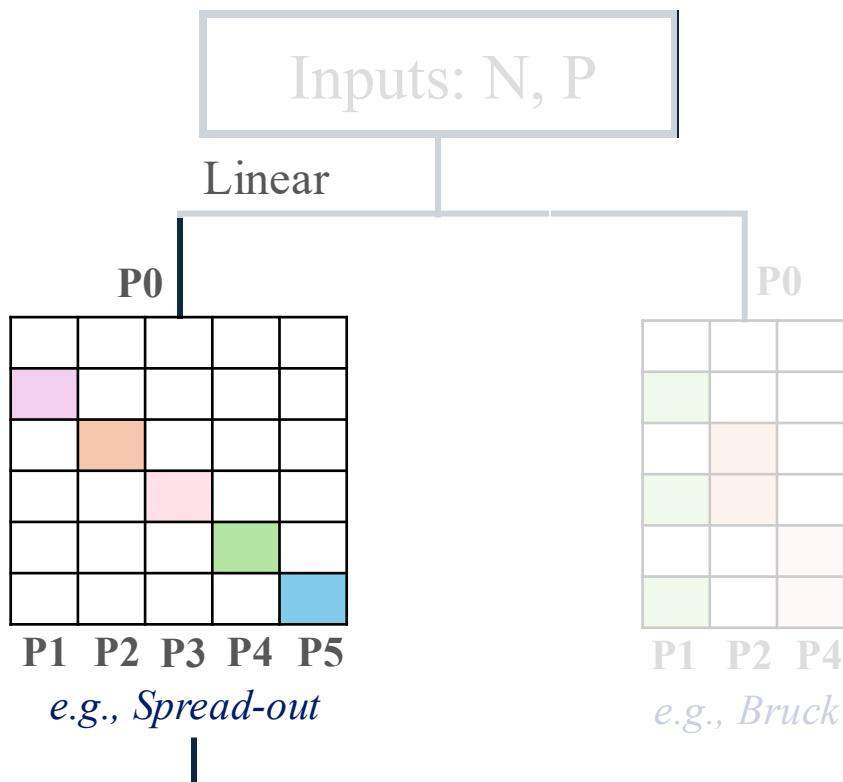
Spread-out algorithm



# Standard MPI All-to-all Implementations: Spread-out

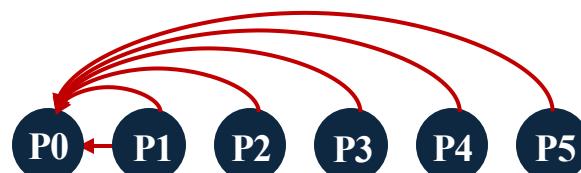
*N: Size of data-block*

*P: Process count*



Spread-out algorithm

**Pose receive requests  
using `MPI_Irecv`**



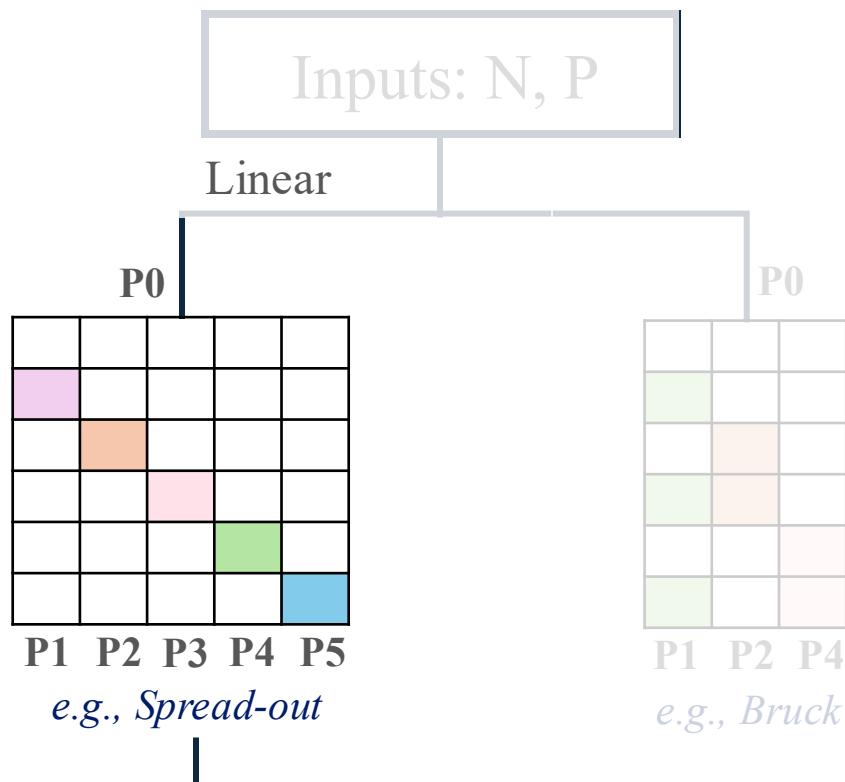
$(rank + i) \% P$

Taking  $P_0$  as example

# Standard MPI All-to-all Implementations: Spread-out

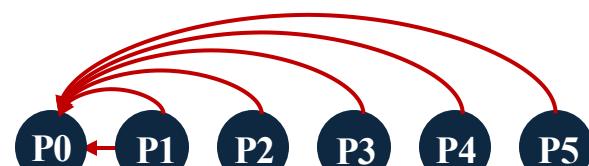
*N: Size of data-block*

*P: Process count*

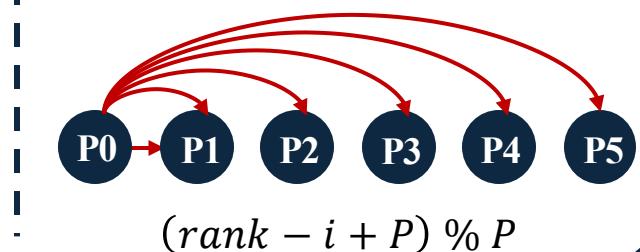


Spread-out algorithm

**Pose receive requests  
using MPI\_Irecv**



**Pose send requests  
using MPI\_Isend**

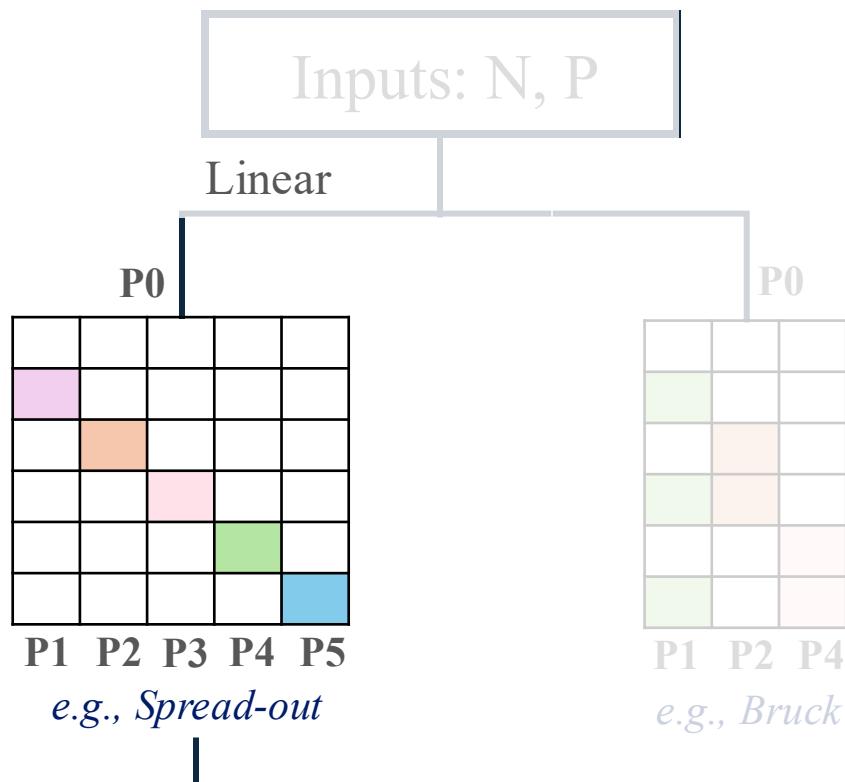


Taking P0 as example

# Standard MPI All-to-all Implementations: Spread-out

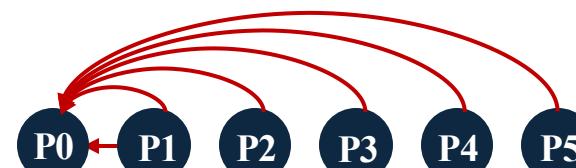
$N$ : Size of data-block

$P$ : Process count

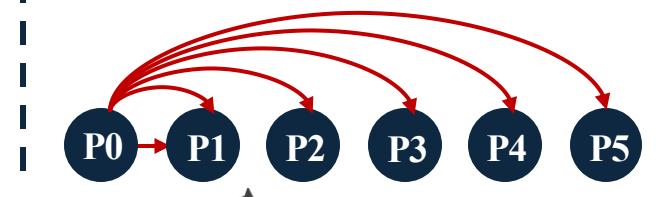


Spread-out algorithm

Pose receive requests  
using **`MPI_Irecv`**



Pose send requests  
using **`MPI_Isend`**

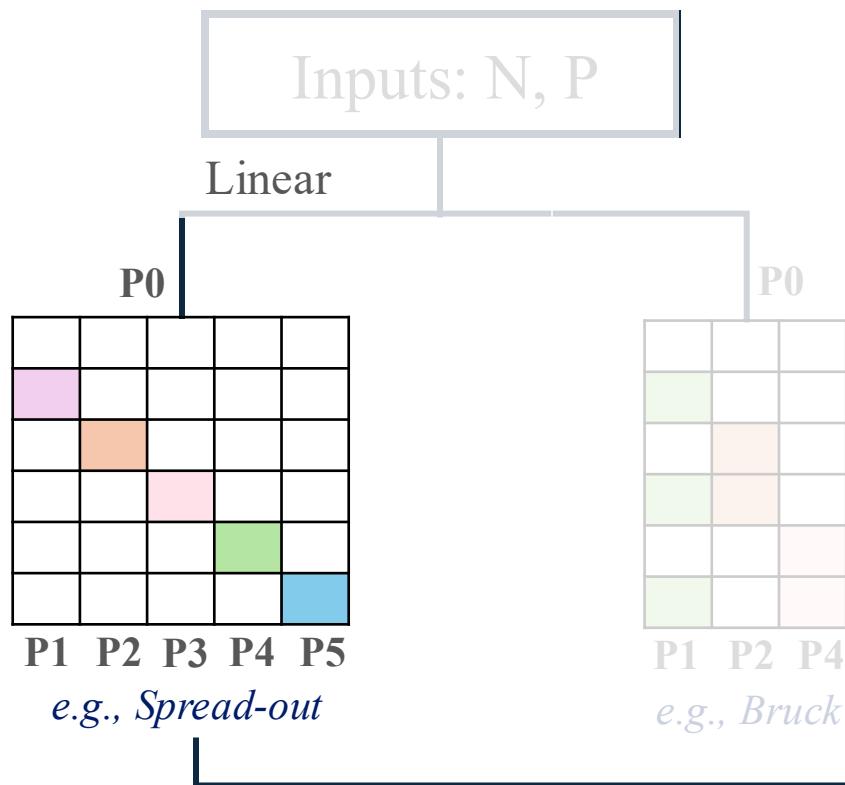


Taking  $P0$  as example

# Standard MPI All-to-all Implementations: Spread-out

*N: Size of data-block*

*P: Process count*



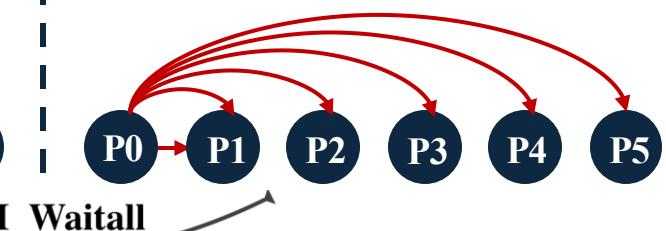
Spread-out algorithm

**Pose receive requests  
using MPI\_Irecv**



**MPI\_Waitall**

**Pose send requests  
using MPI\_Isend**

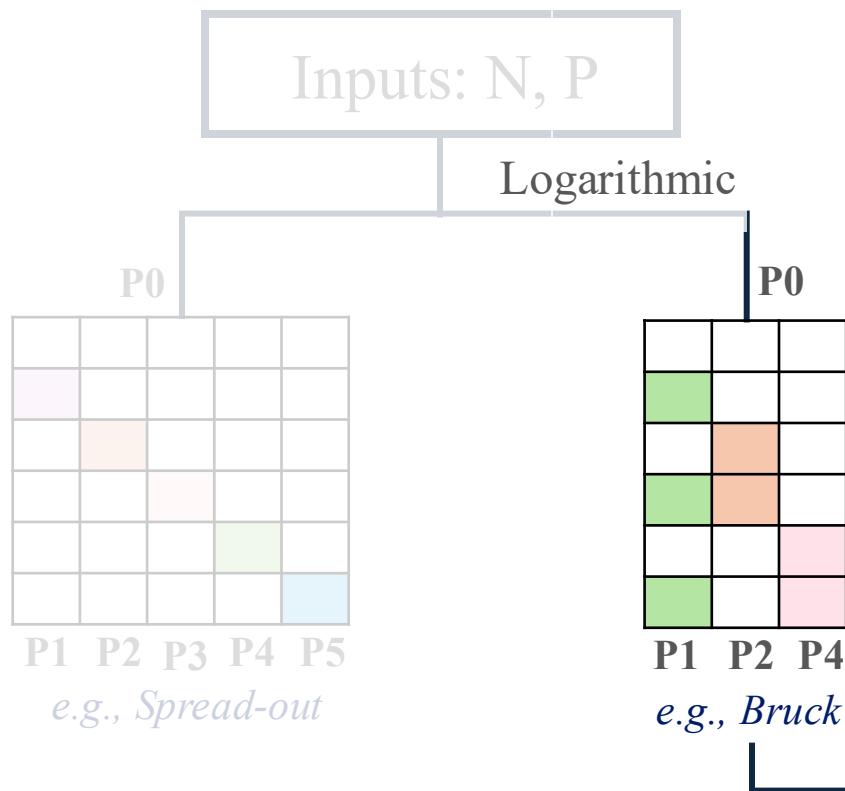


( $P - 1$ ) communication steps.

# Standard MPI All-to-all Implementations: Bruck

$N$ : Size of data-block

$P$ : Process count



Bruck algorithm

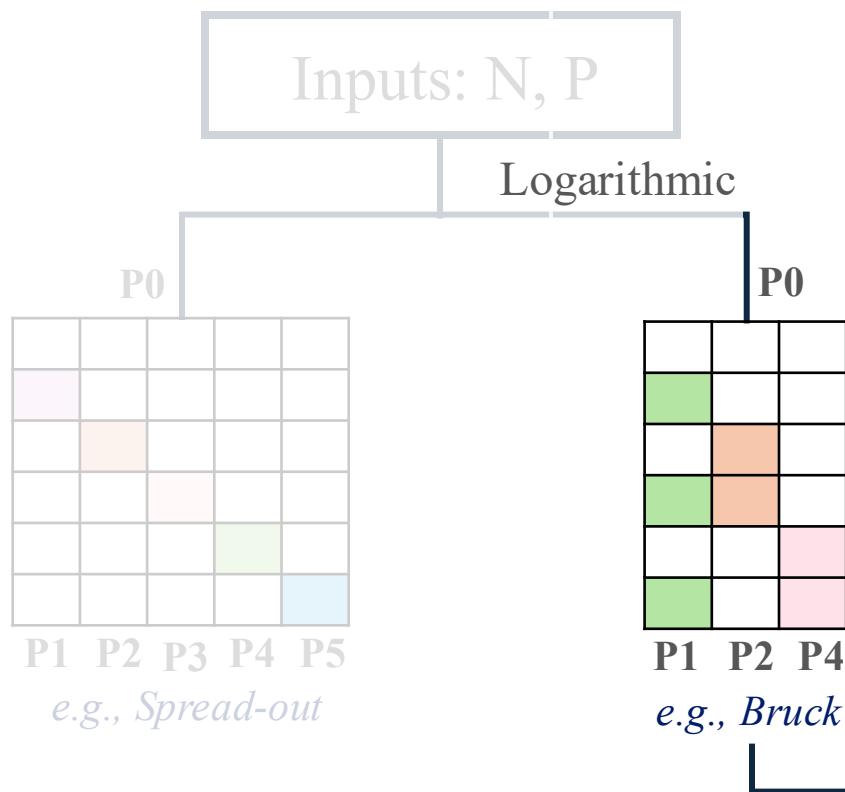


$\log(P)$  communication steps

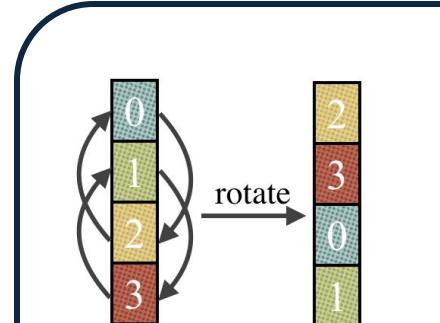
# Standard MPI All-to-all Implementations: Bruck

$N$ : Size of data-block

$P$ : Process count



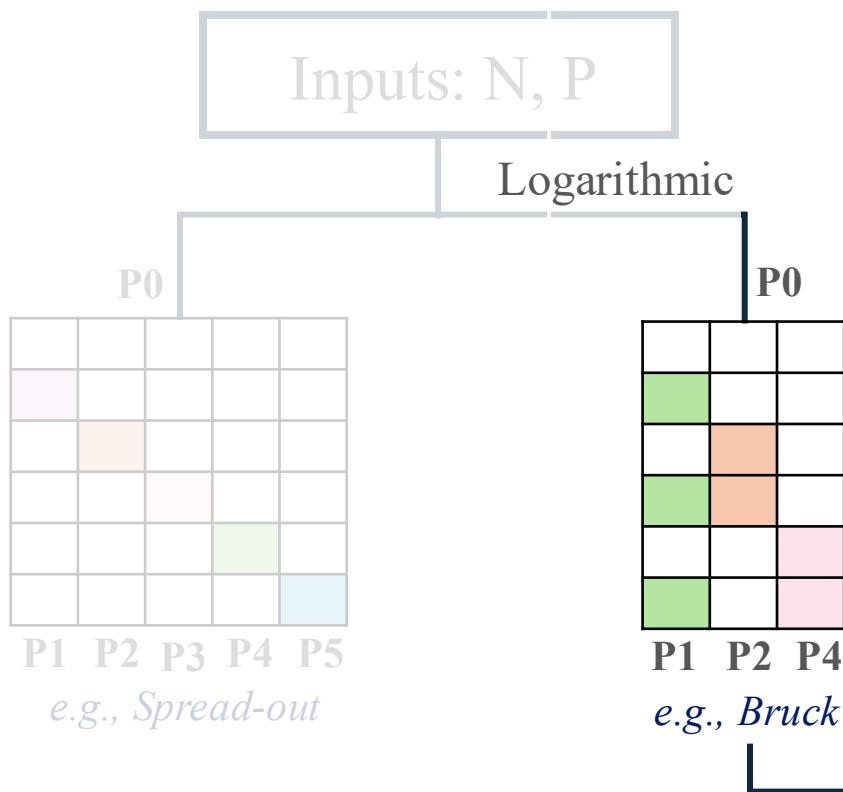
Bruck algorithm



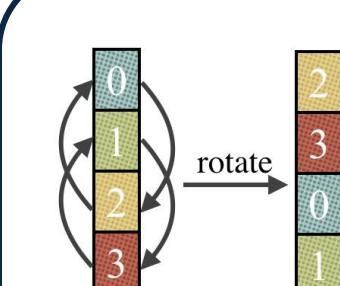
# Standard MPI All-to-all Implementations: Bruck

*N: Size of data-block*

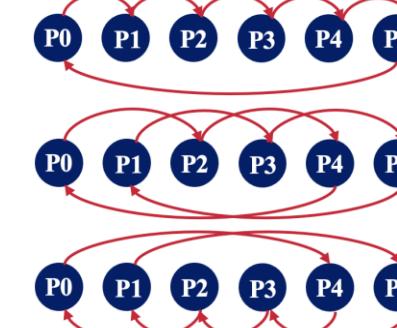
*P: Process count*



Bruck algorithm



2.  $\log(P)$  communication steps.

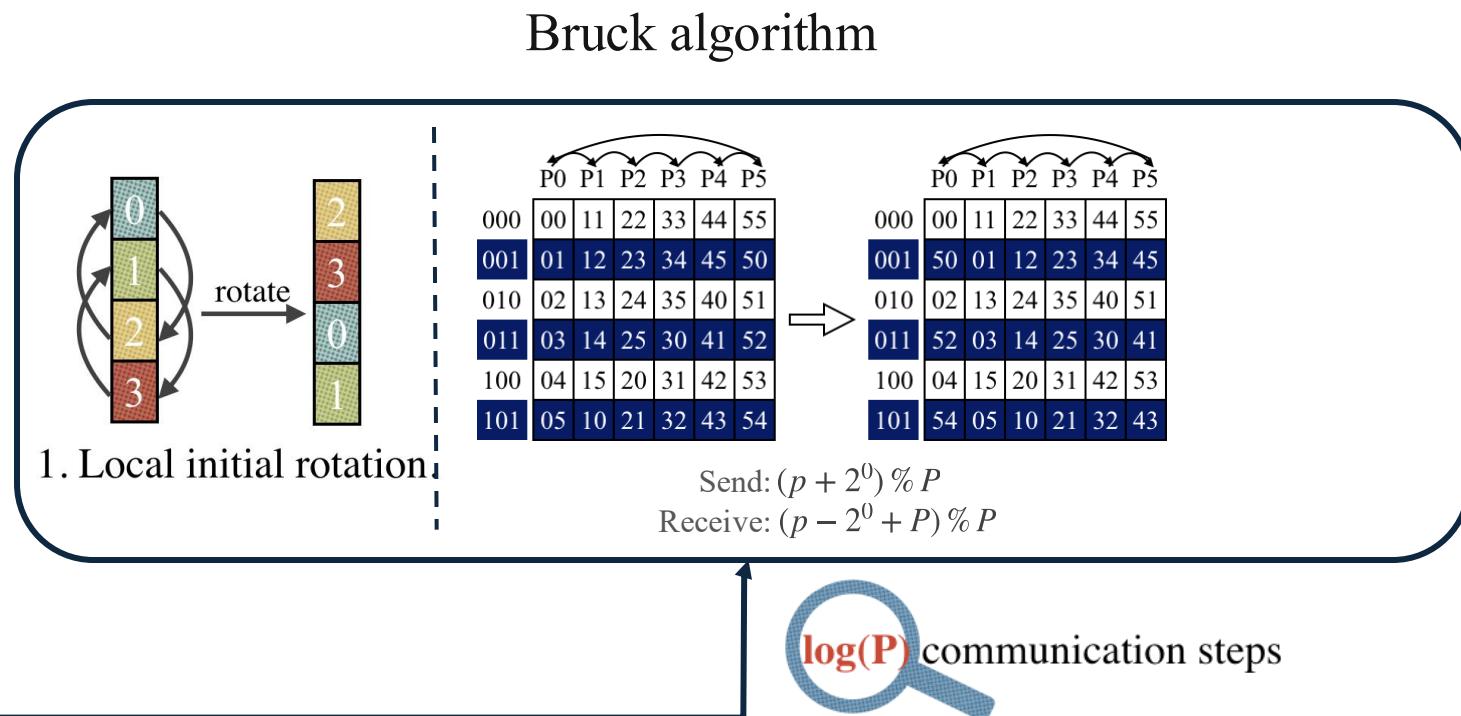
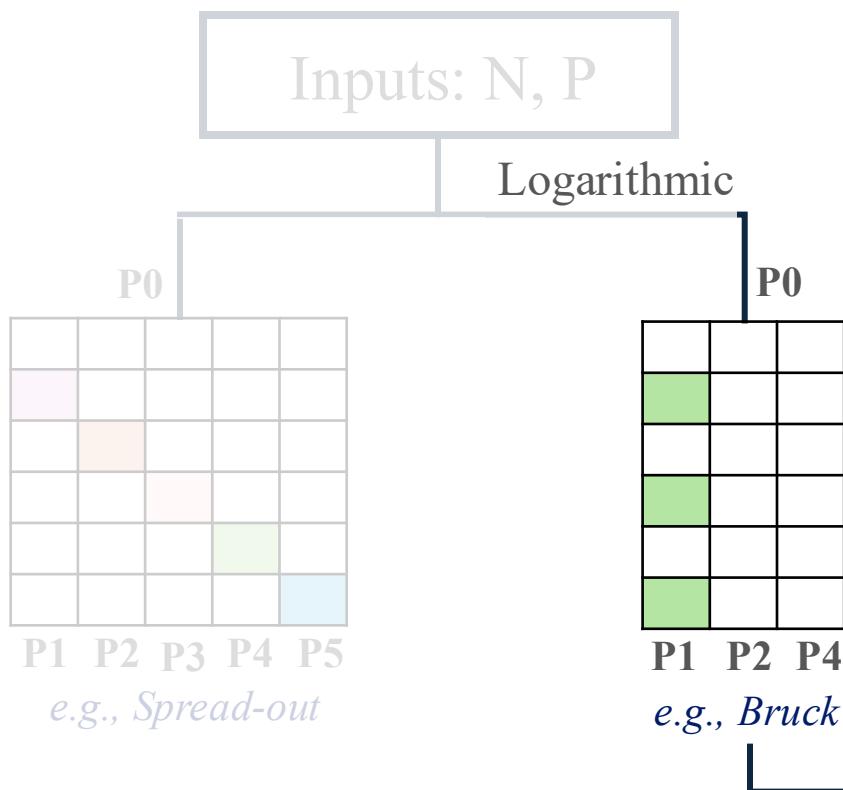


$\log(P)$  communication steps

# Standard MPI All-to-all Implementations: Comm-0

*N: Size of data-block*

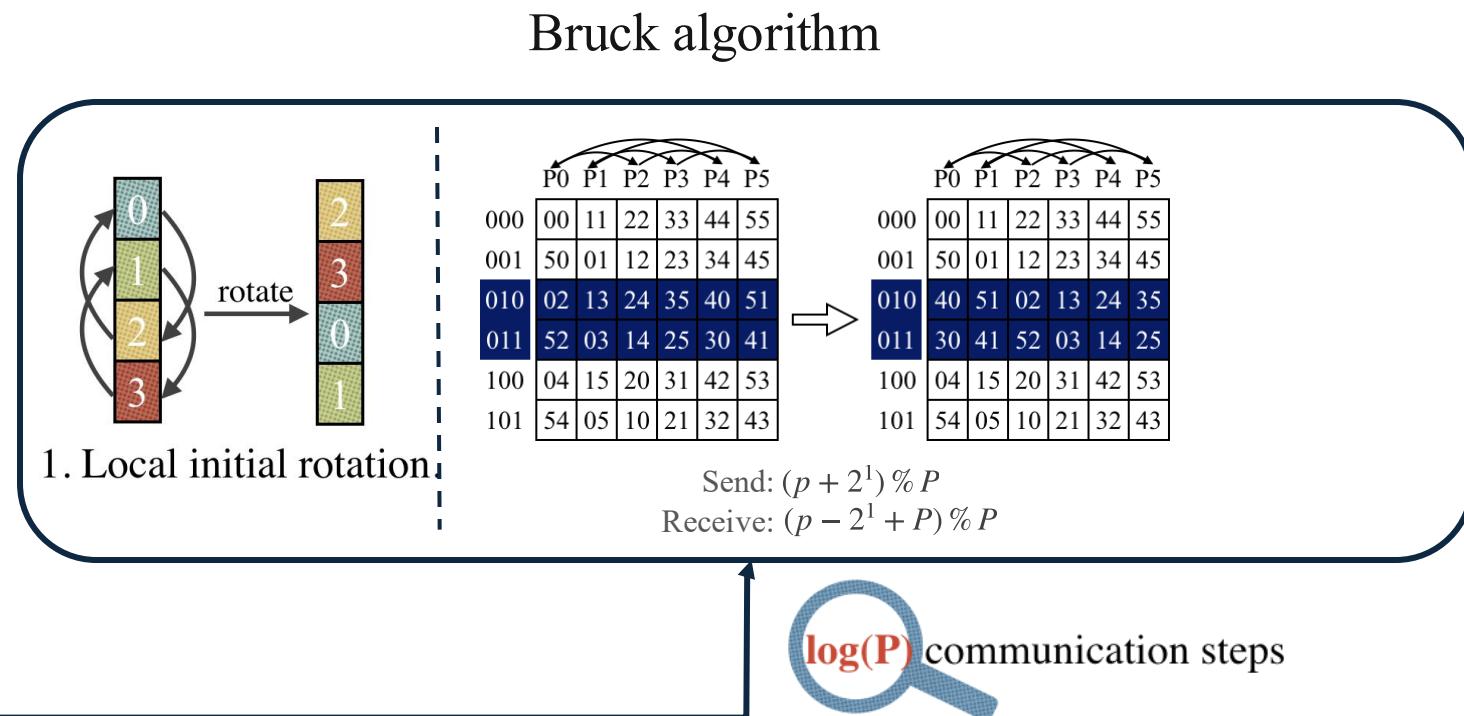
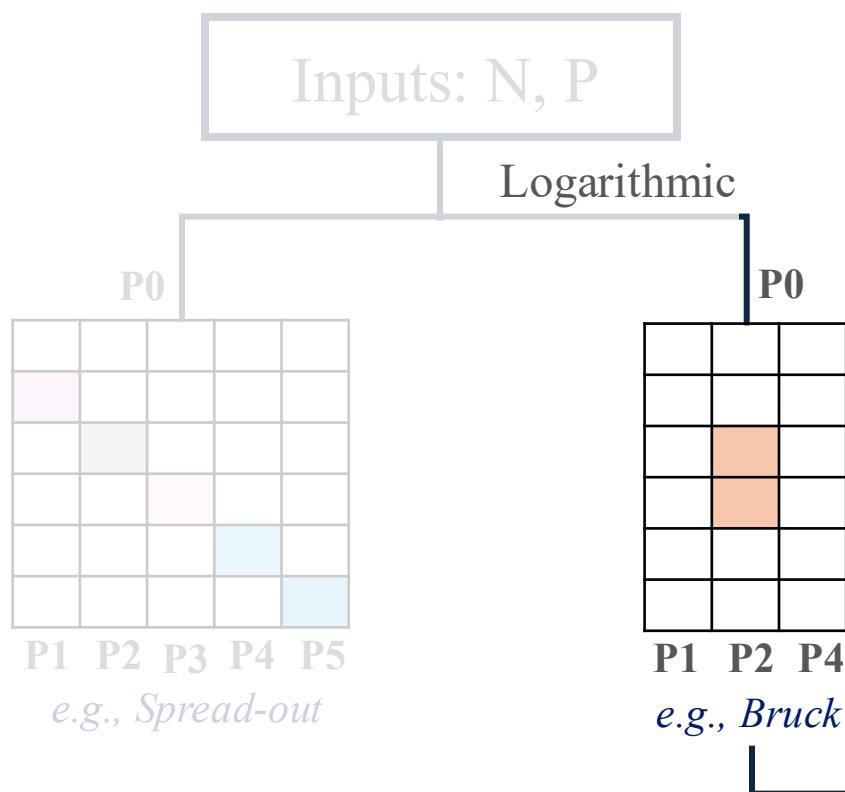
*P: Process count*



# Standard MPI All-to-all Implementations: Comm-1

*N: Size of data-block*

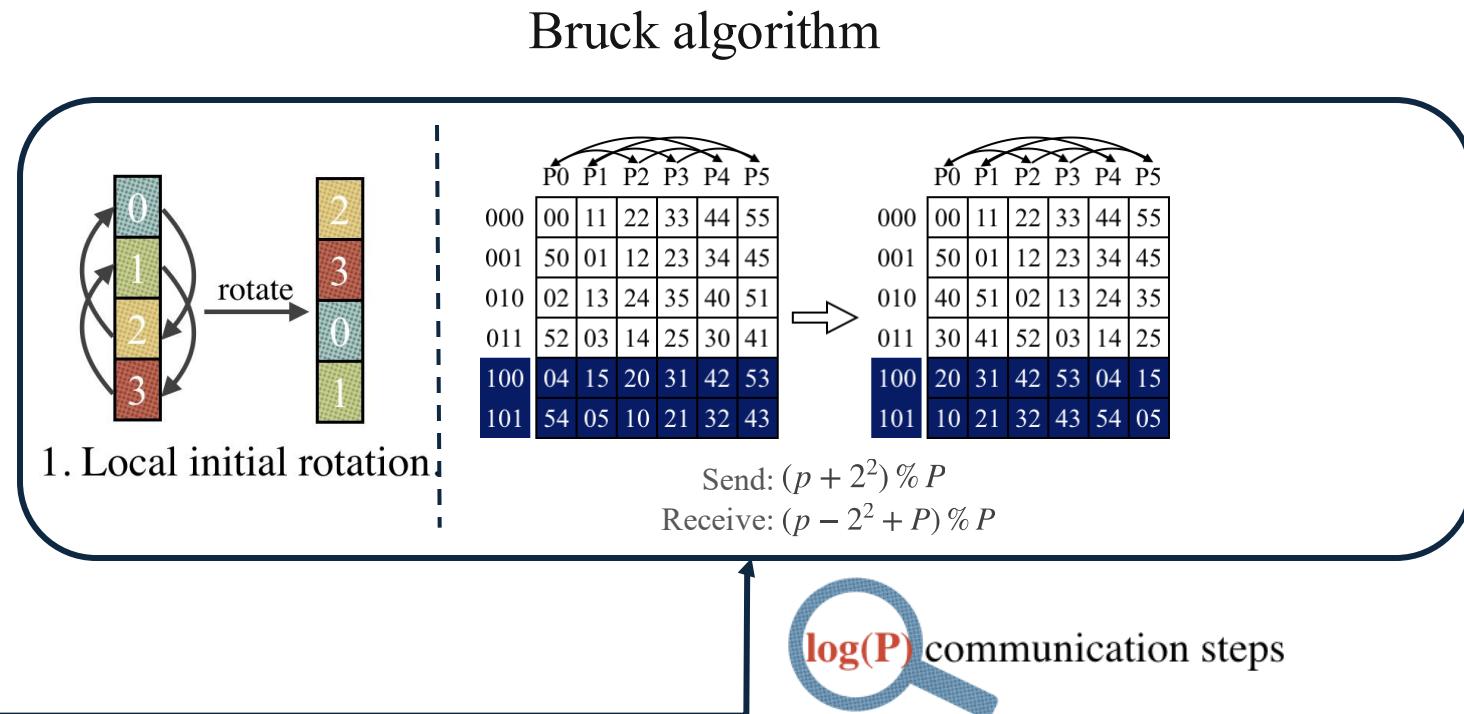
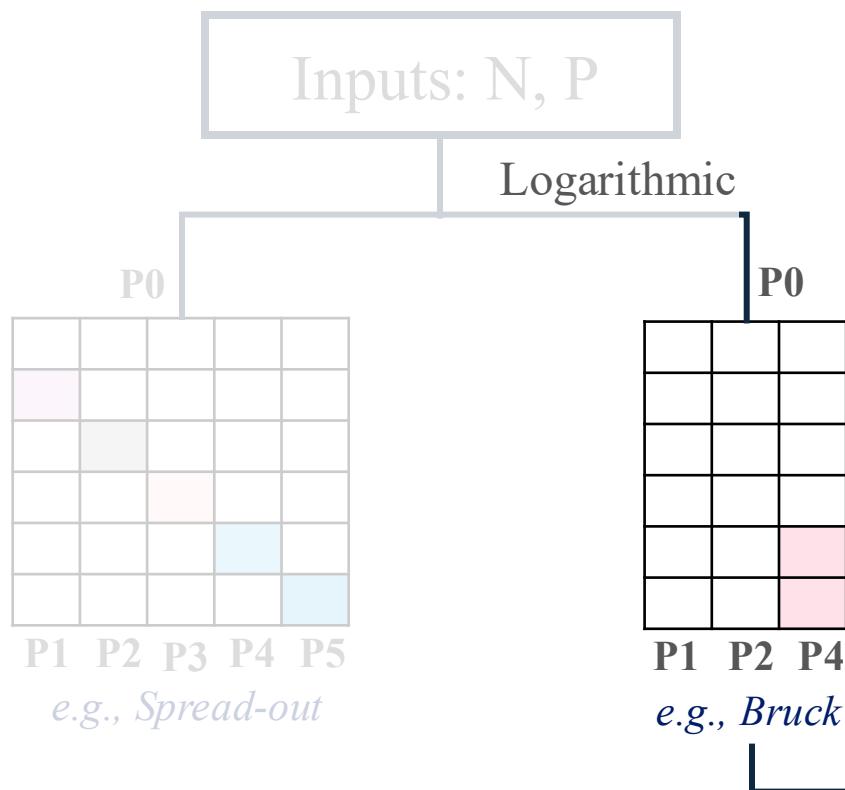
*P: Process count*



# Standard MPI All-to-all Implementations: Comm-2

*N: Size of data-block*

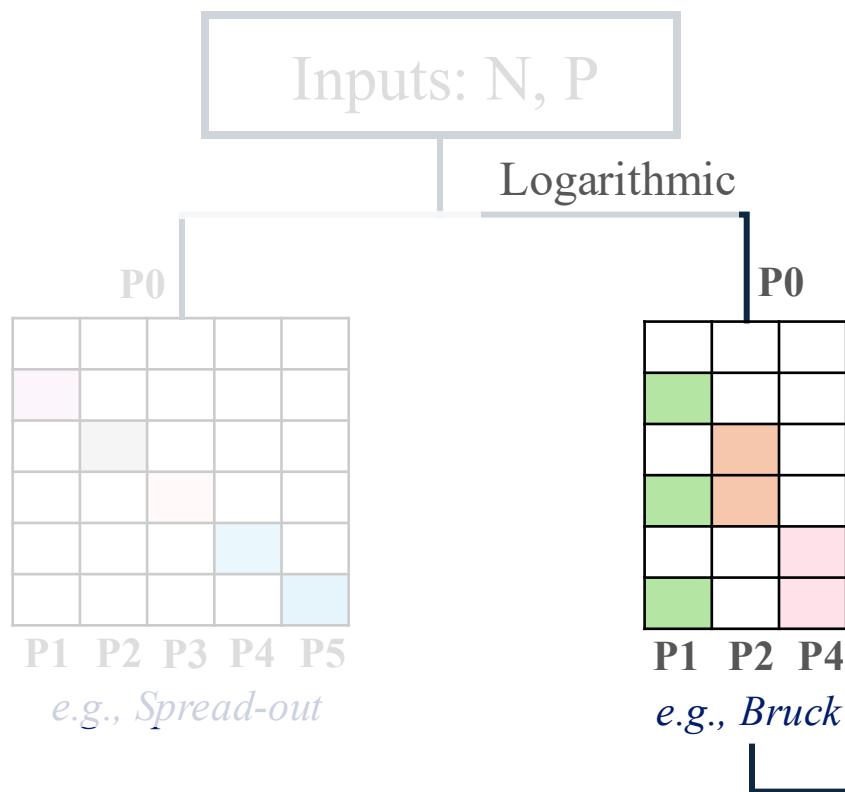
*P: Process count*



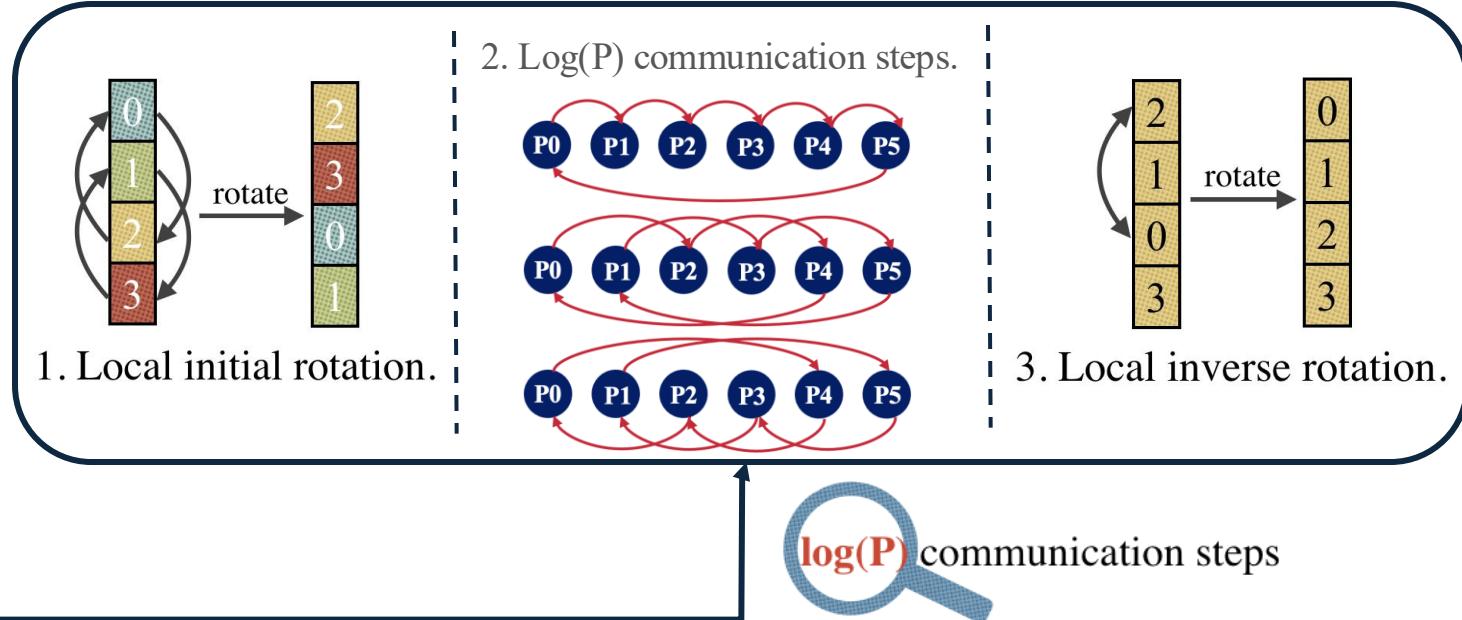
# Standard MPI All-to-all Implementations: Bruck

$N$ : Size of data-block

$P$ : Process count



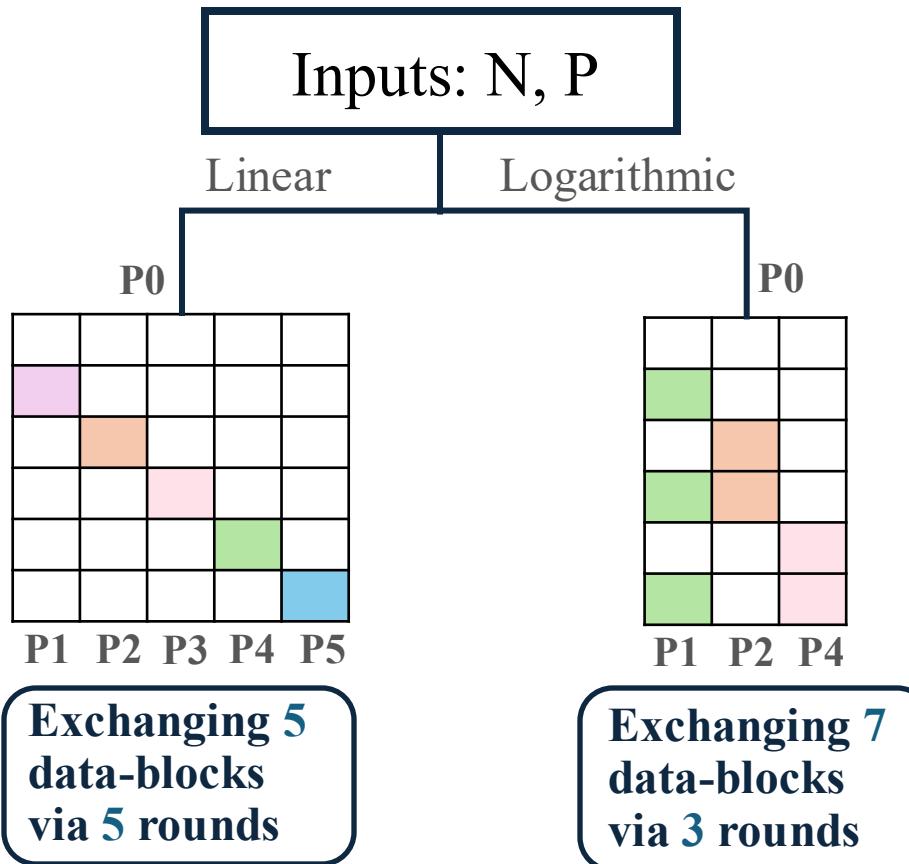
Bruck algorithm



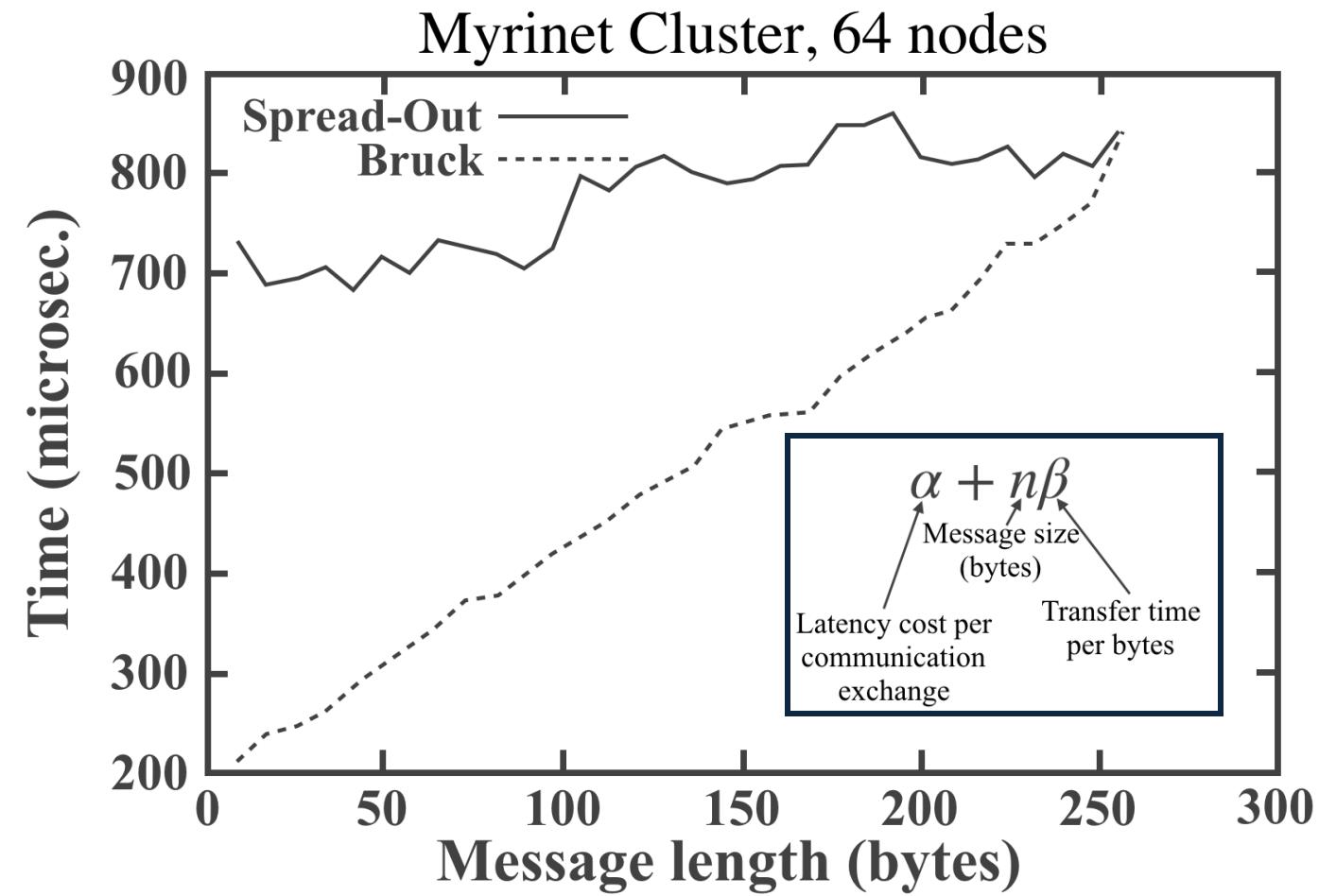
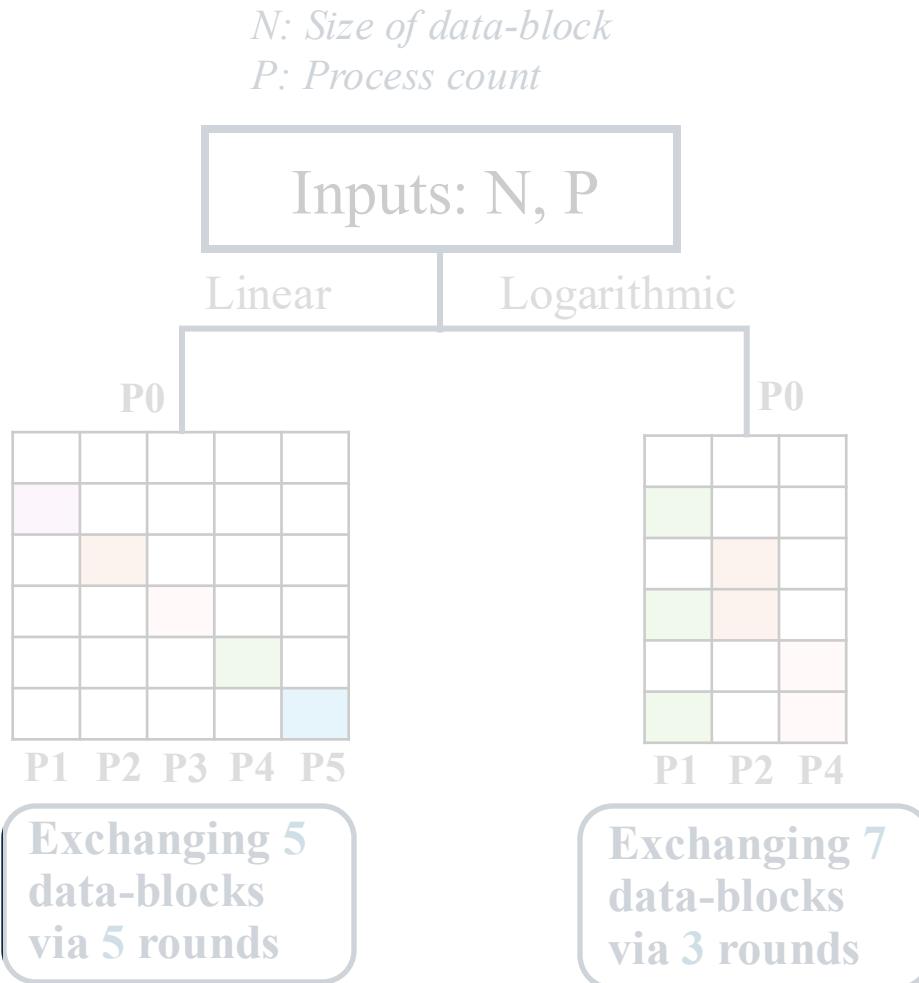
# Comparison of These Two Algorithms

$N$ : Size of data-block

$P$ : Process count



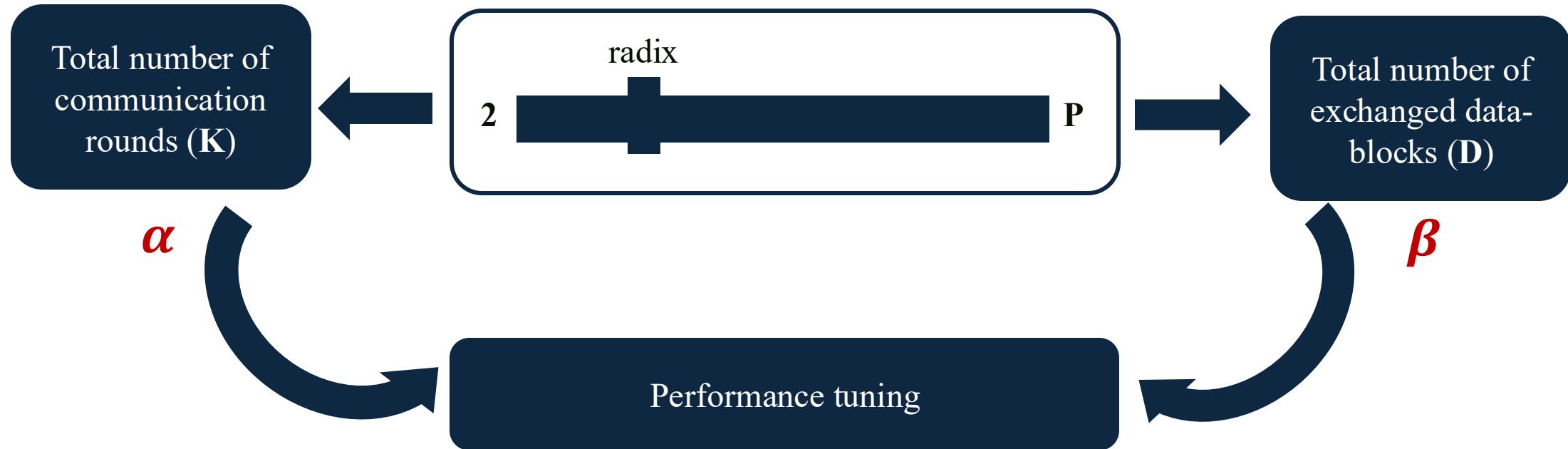
# Comparison of These Two Algorithms



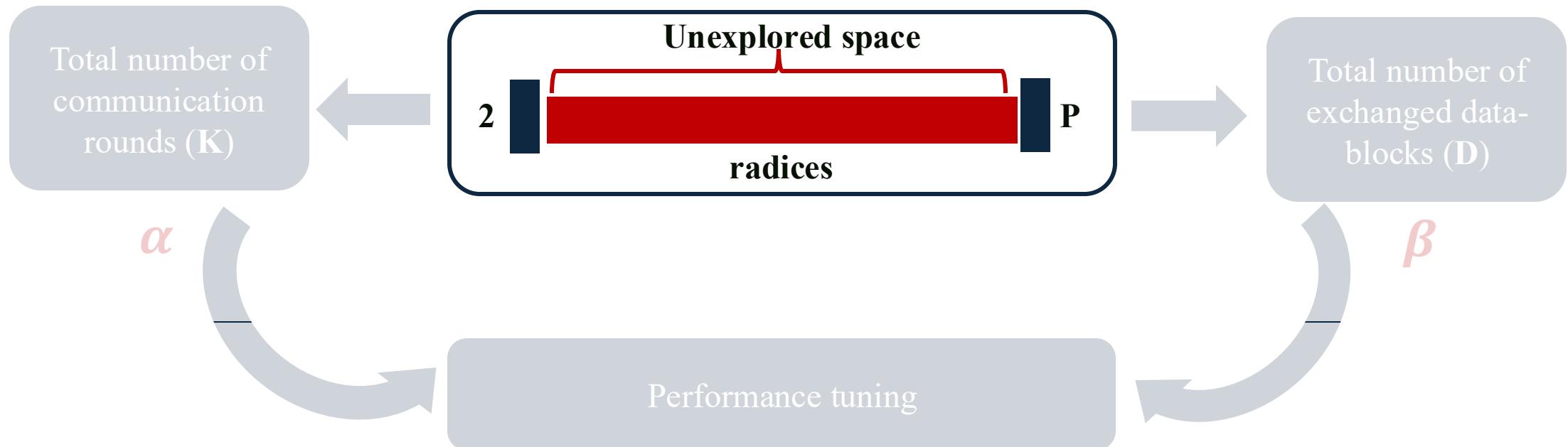
# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- Parameterized Logarithmic Algorithm
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

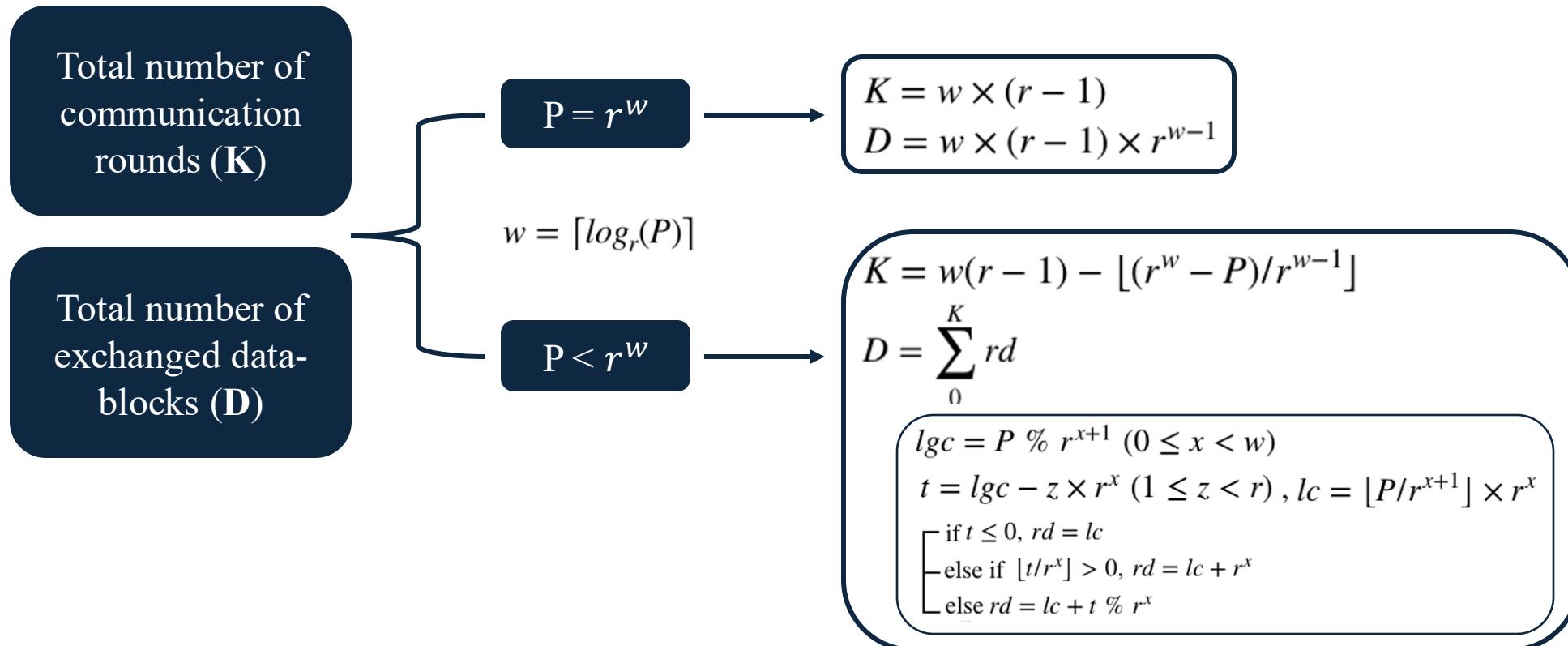
# Tuning Radices Enables Performance Tuning



# Standard MPI All-to-all Implementation Leaves a Large Unexplored Space



# Mathematics of Tunable Radix All-to-all



# Tunable Radix All-to-all Example

$$P = 8, r = 3$$

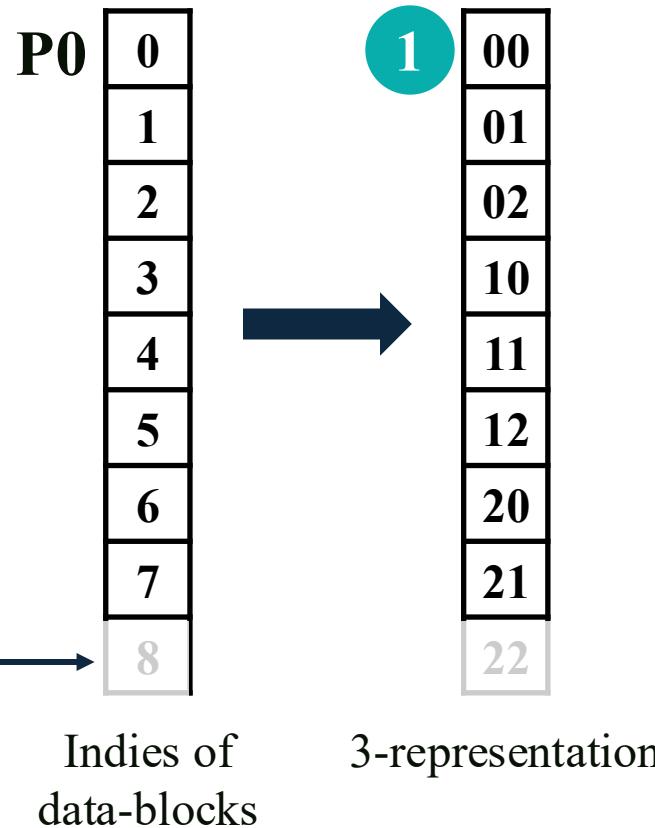
P0
0
1
2
3
4
5
6
7
8

if  $\log_r(P)$   
is integer

Indies of  
data-blocks

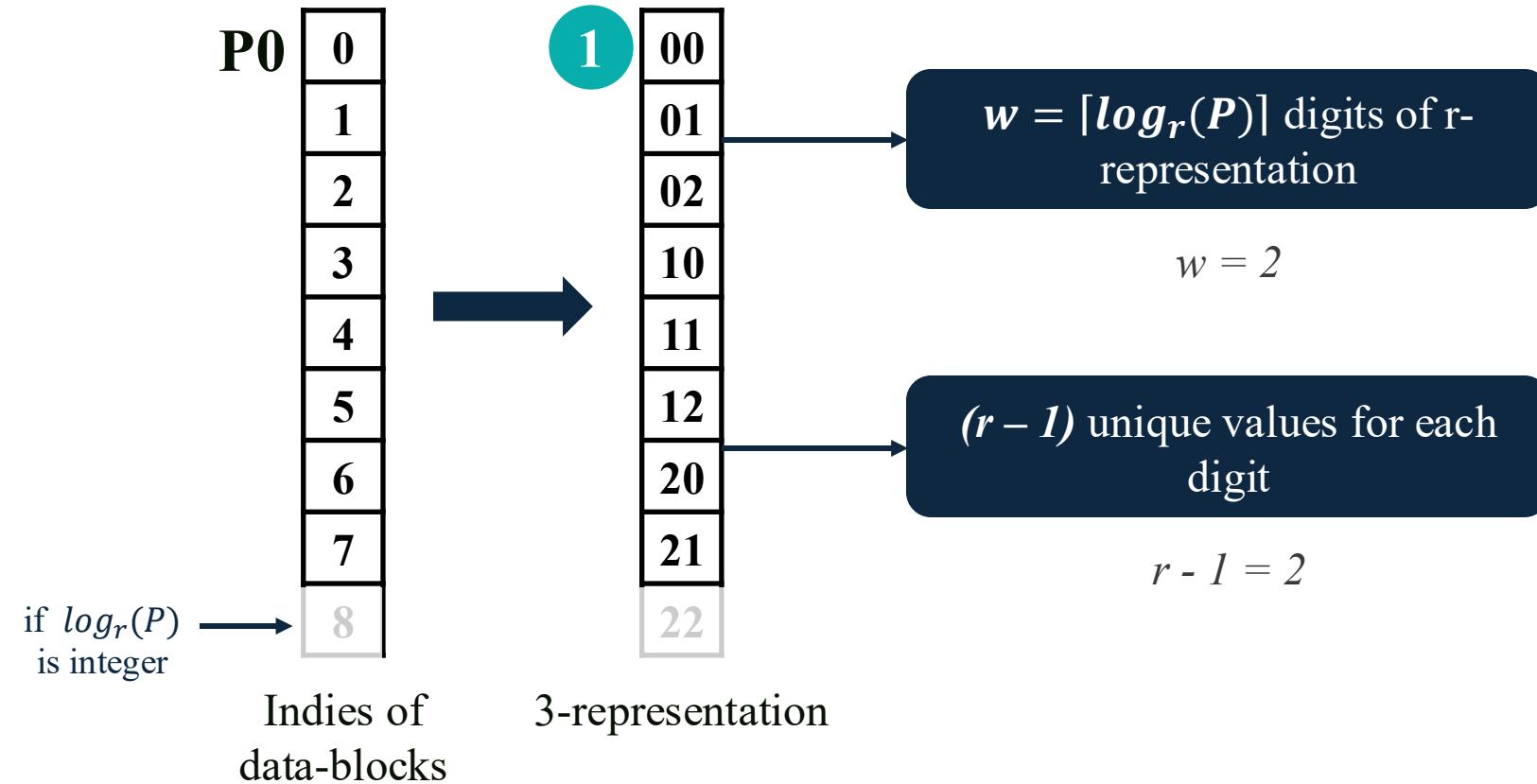
# Tunable Radix All-to-all Example: r-representation

$$P = 8, r = 3$$



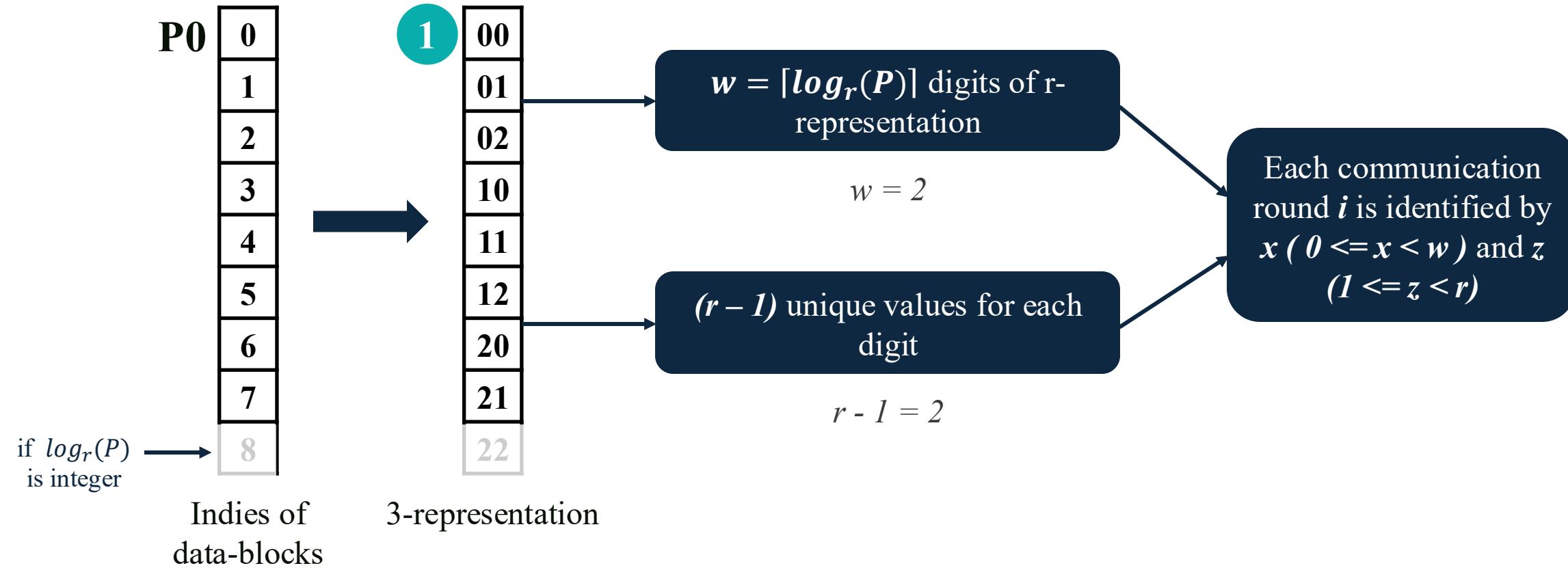
# Tunable Radix All-to-all Example: r-representation

$$P = 8, r = 3$$



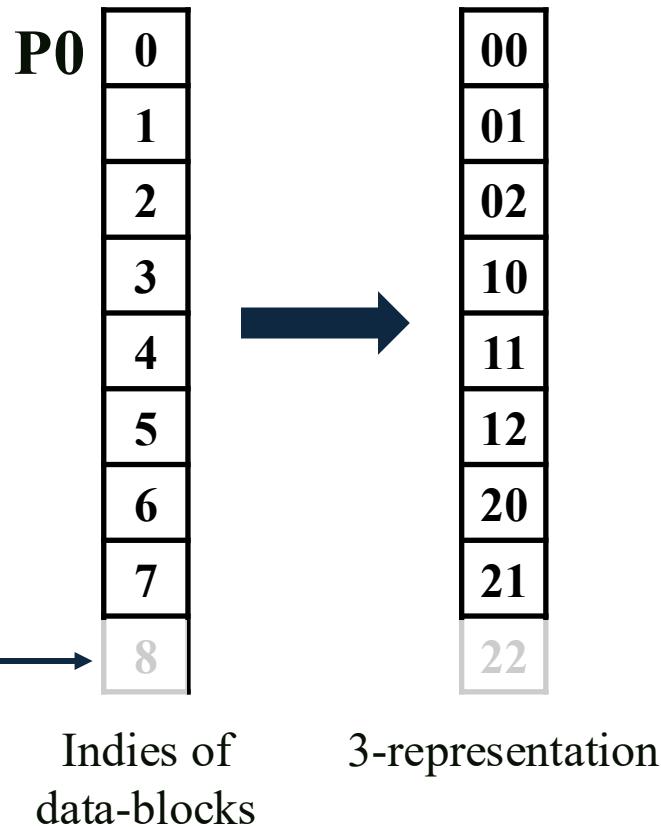
# Tunable Radix All-to-all Example: r-representation

$$P = 8, r = 3$$



# Tunable Radix All-to-all Example: Communication

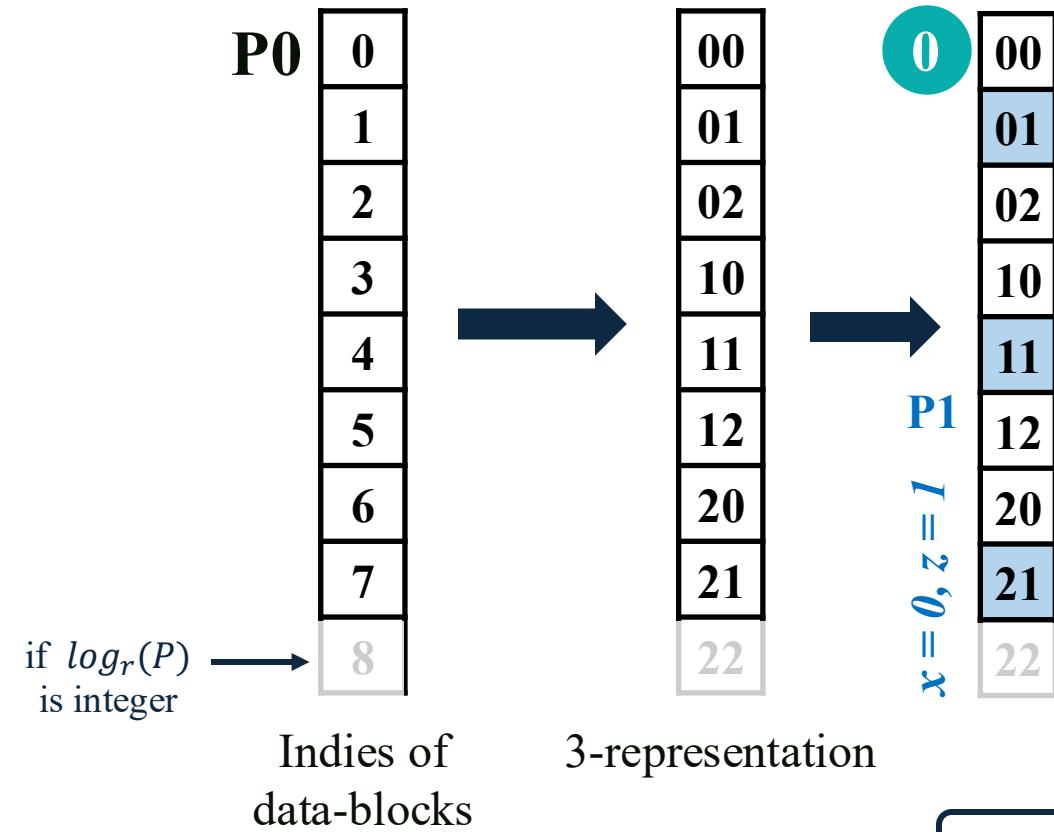
$$P = 8, r = 3$$



P0 sends data-blocks whose  
x digit equals to z

# Tunable Radix All-to-all Example: Comm-round-0

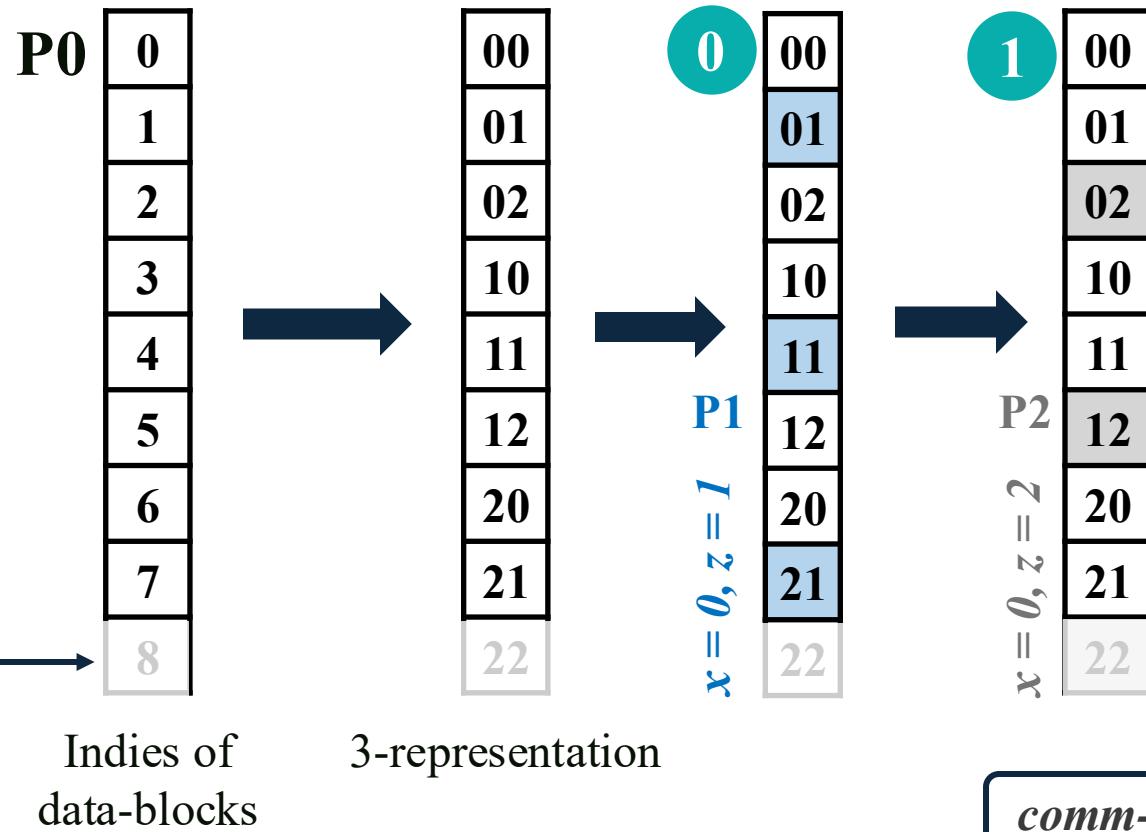
$$P = 8, r = 3$$



*comm-round-0:  $x = 0, z = 1, P_0$  sends to  $P_1$*

# Tunable Radix All-to-all Example: Comm-round-1

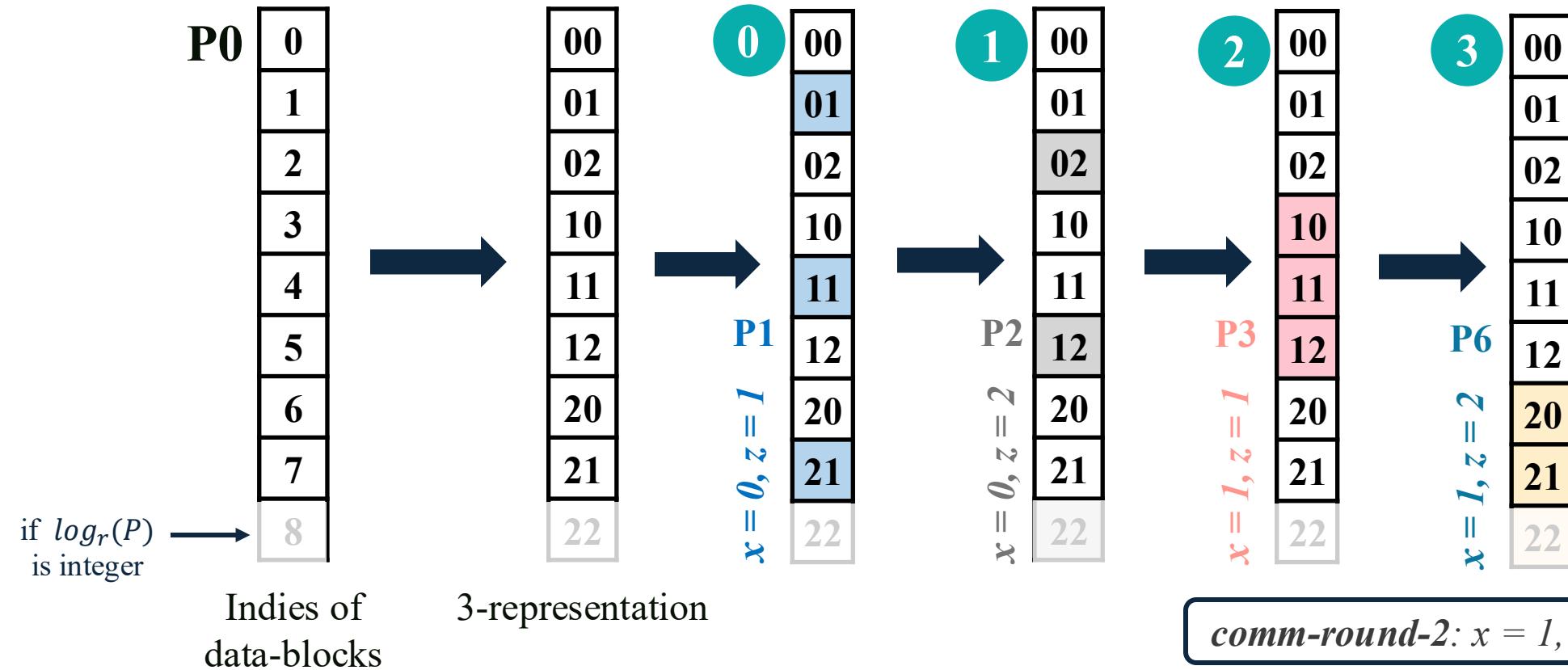
$$P = 8, r = 3$$



**comm-round-1:**  $x = 0, z = 2, P_0$  sends to  $P_2$

# Tunable Radix All-to-all Example: Comm-round-2,3

$$P = 8, r = 3$$



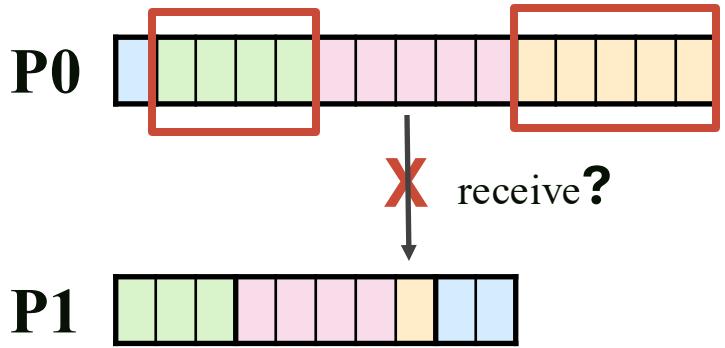
**comm-round-2:  $x = 1, z = 1, P0$  sends to  $P3$**

**comm-round-3:  $x = 1, z = 2, P0$  sends to  $P6$**

# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - **Logarithmic Non-uniform All-to-all**
- Parameterized Logarithmic Algorithm
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

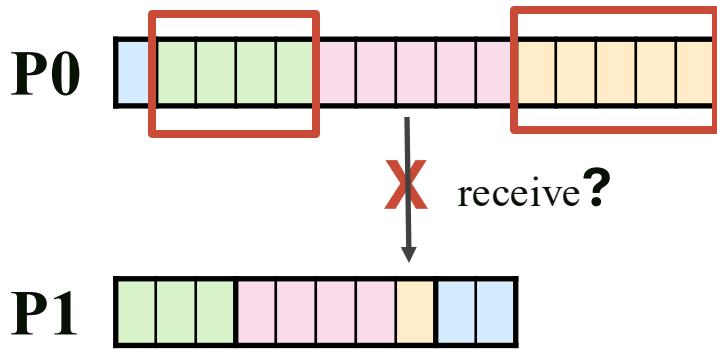
# Challenges of Applying the Logarithmic Algorithm to Non-uniform All-to-all Data Exchange



Each process is unaware of how much data to expect during each intermediate communication round.

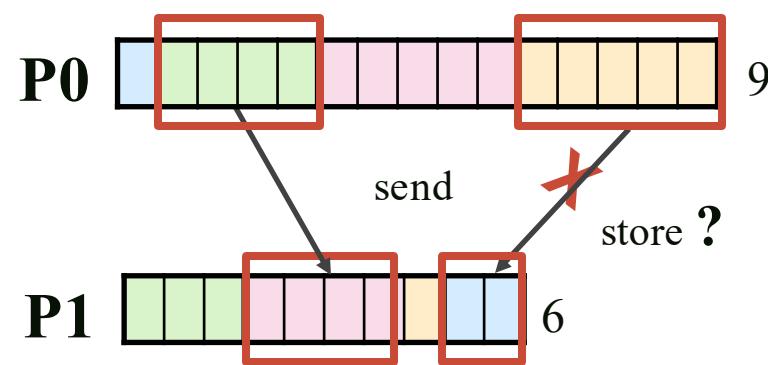
Non-uniform  
workloads

# Challenges of Applying the Logarithmic Algorithm to Non-uniform All-to-all Data Exchange



Each process is unaware of how much data to expect during each intermediate communication round.

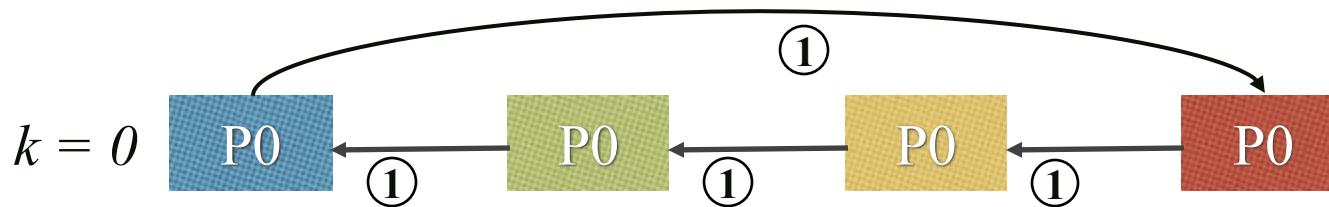
Non-uniform workloads



The received data elements could be too large to fit into the segment in the send buffer.

Store-and-forward

# Two-phase Communication Scheme

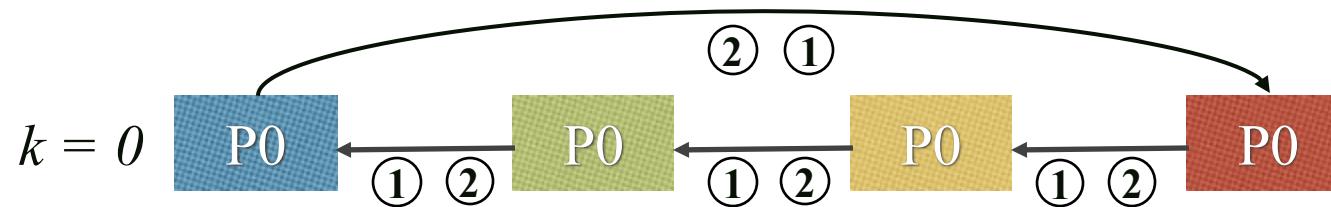


① *Metadata Exchange*

a	b
---	---

 Size of each block

# Two-phase Communication Scheme



- ① *Metadata Exchange*

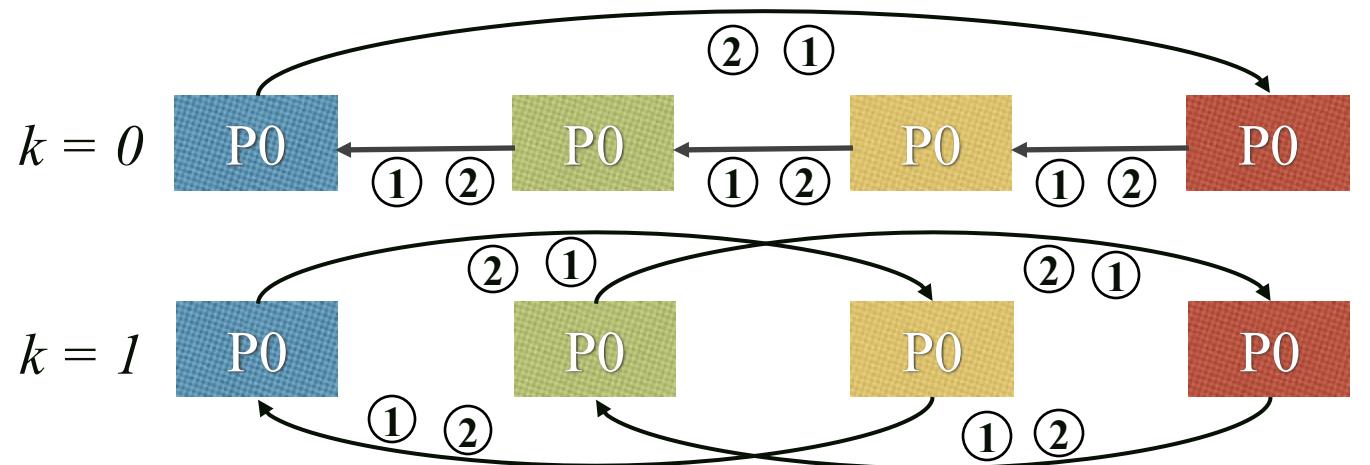
a	b
---	---

 Size of each block
- ② *Data Exchange*

Block 1	Block 3
---------	---------

 Actual data-blocks

# Two-phase Communication Scheme



① *Metadata Exchange*

a	b
---	---

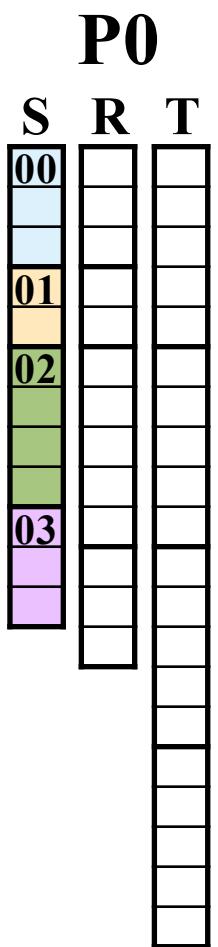
 Size of each block

② *Data Exchange*

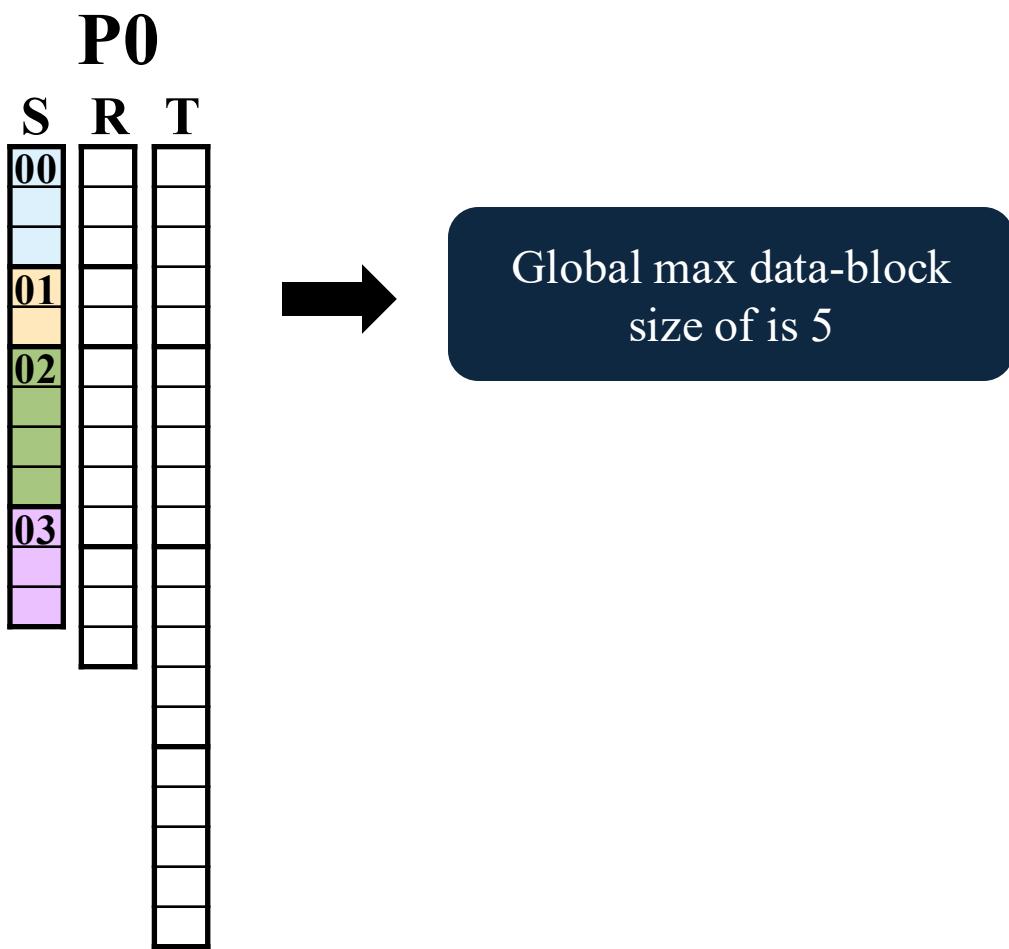
Block 1	Block 3
---------	---------

 Actual data-blocks

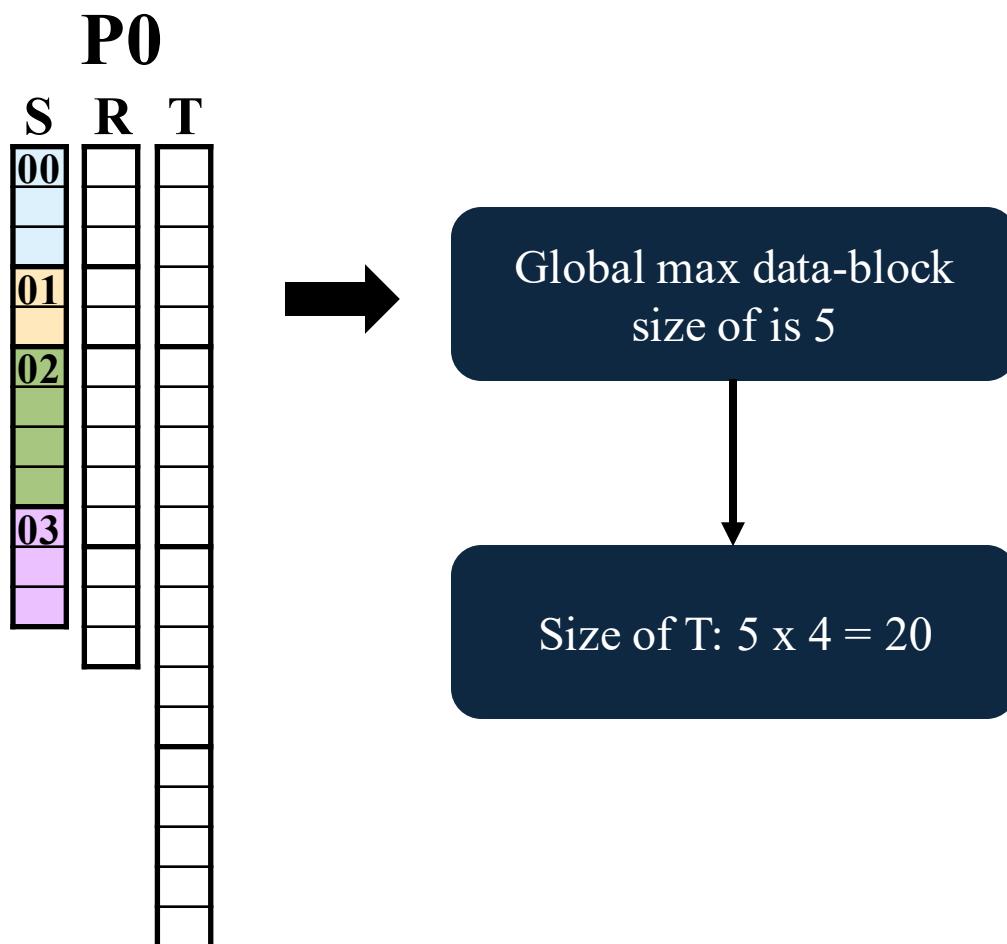
# Using Temporary Buffer to Solve Challenge 2



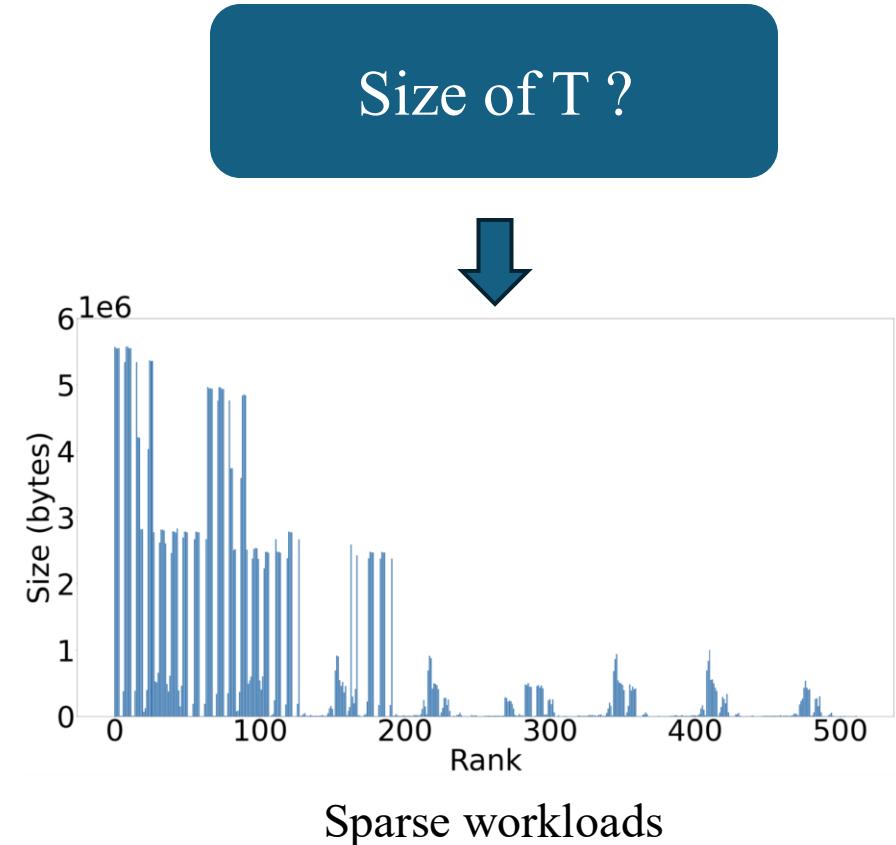
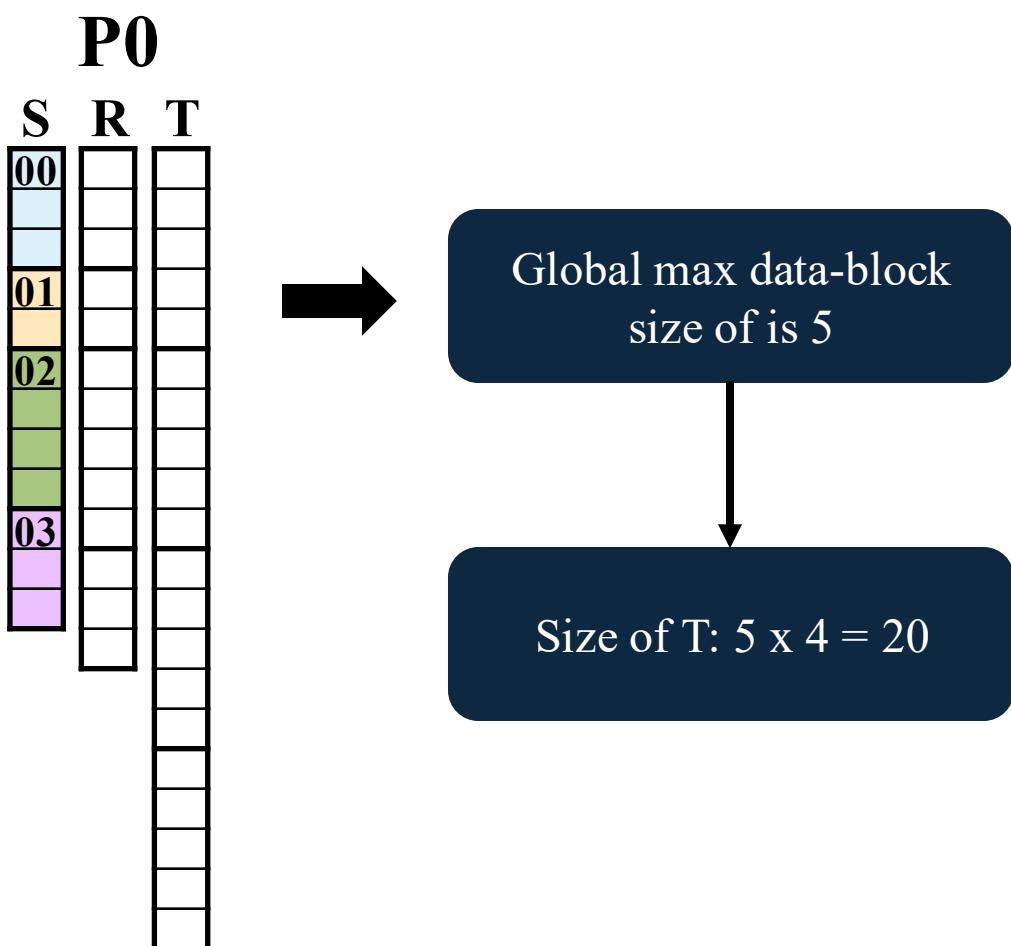
# Estimating Temporary Buffer Size



# Estimating Temporary Buffer Size

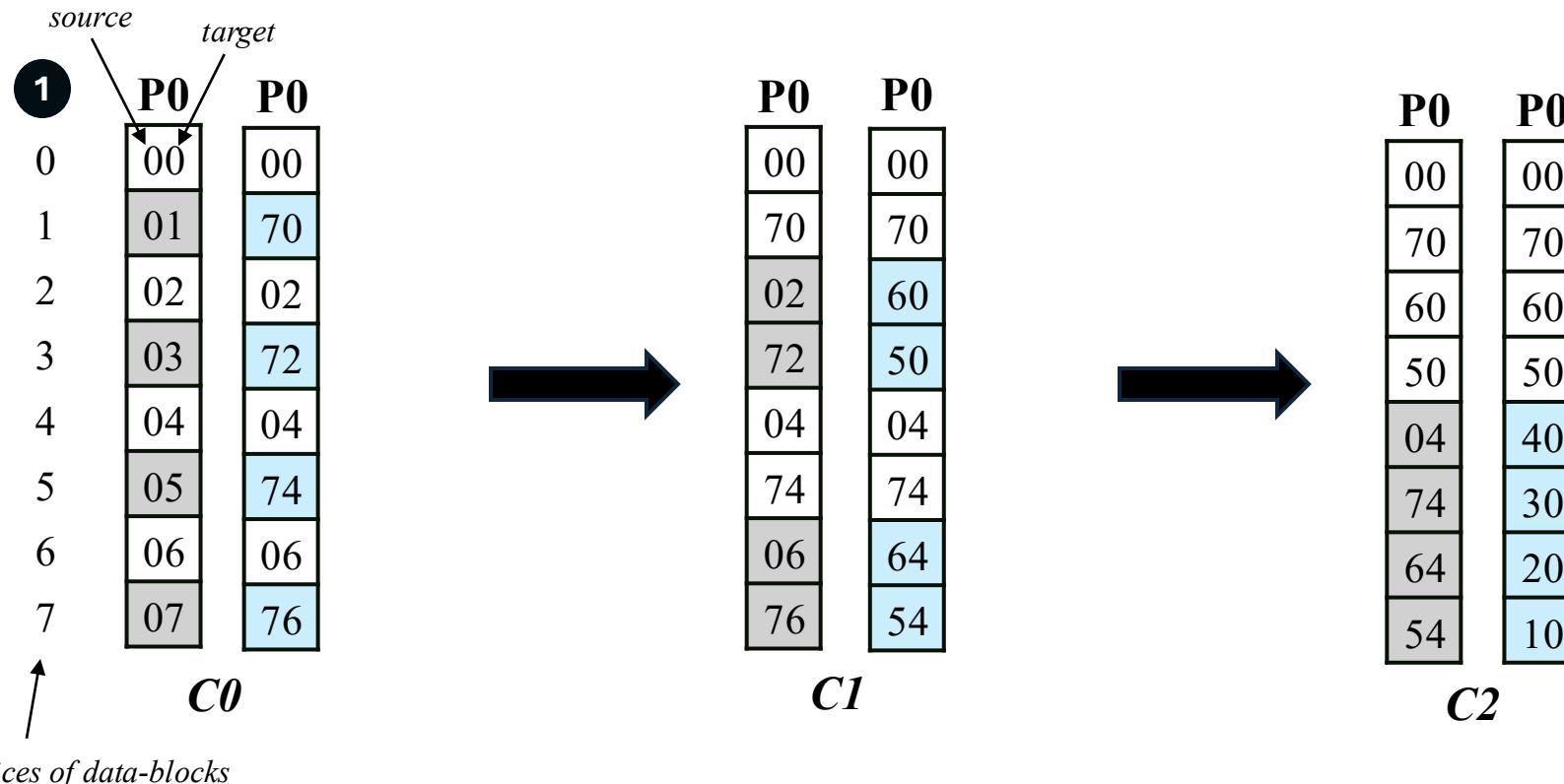


# Temporary Buffer Size Can Be Very Large



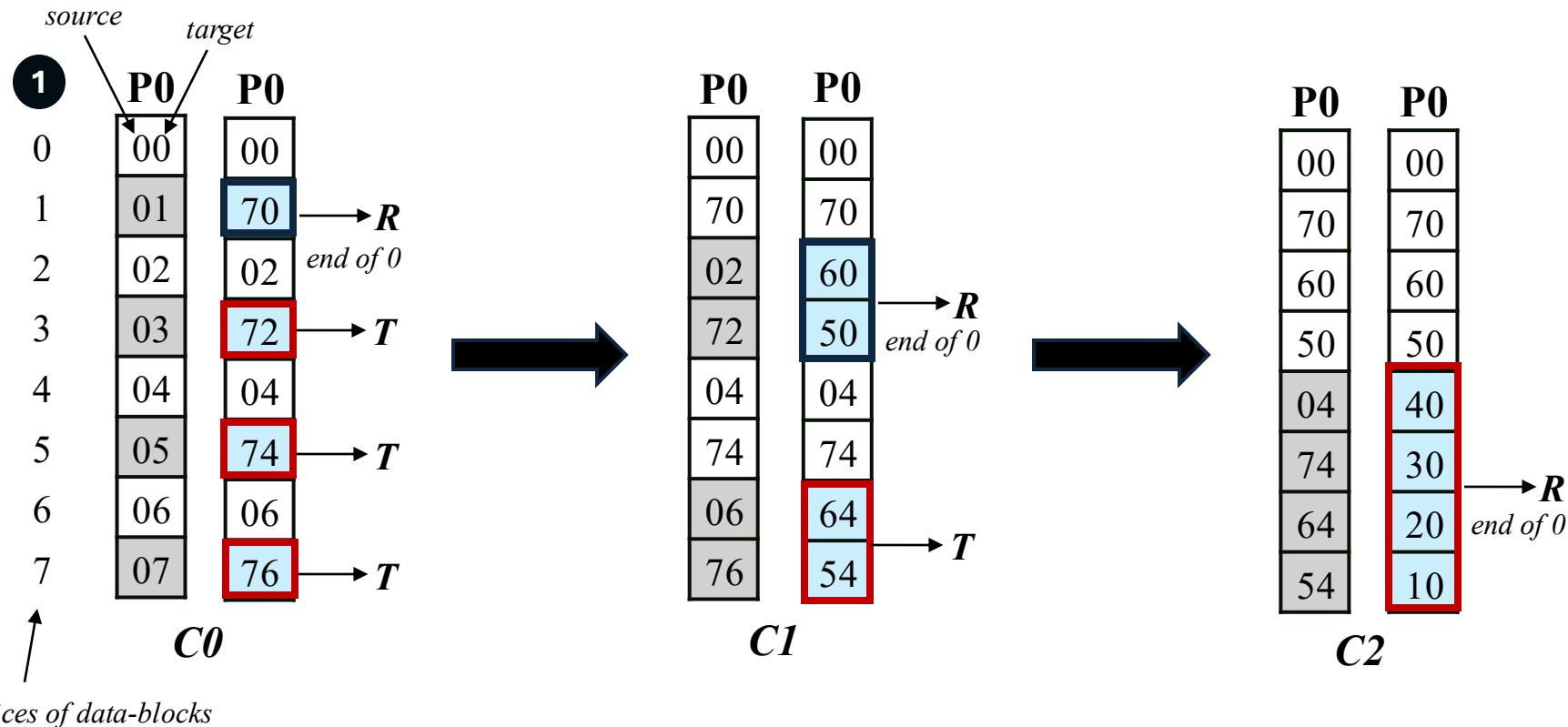
# Temporary Buffer Size Analysis

$$P = 8, r = 2$$



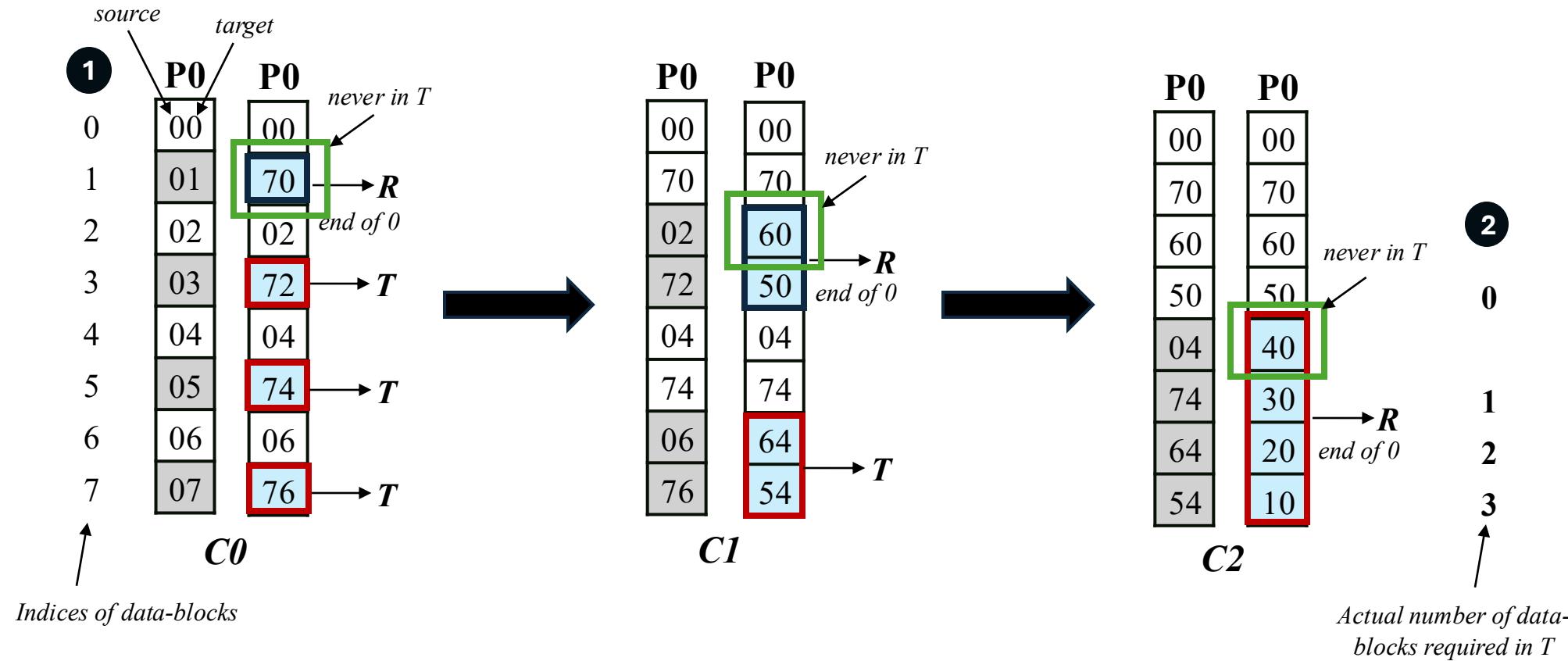
# Temporary Buffer Size Analysis

$$P=8, r=2$$

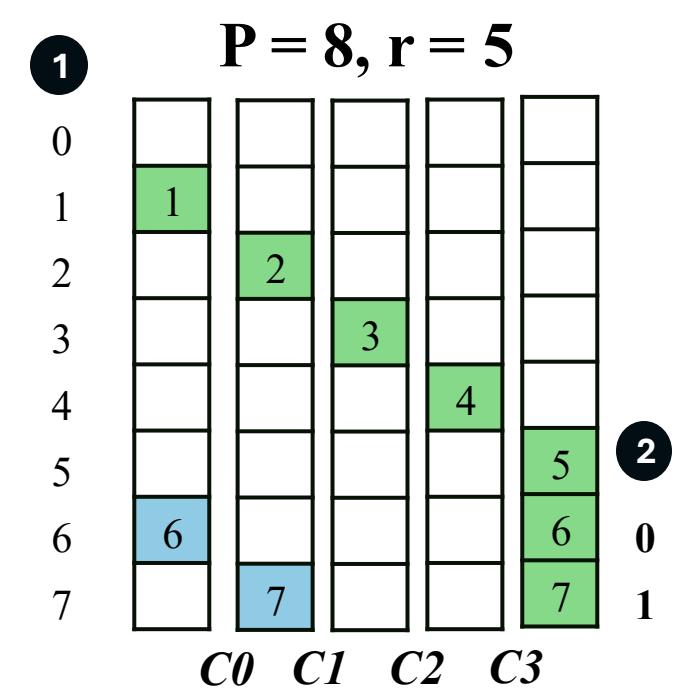
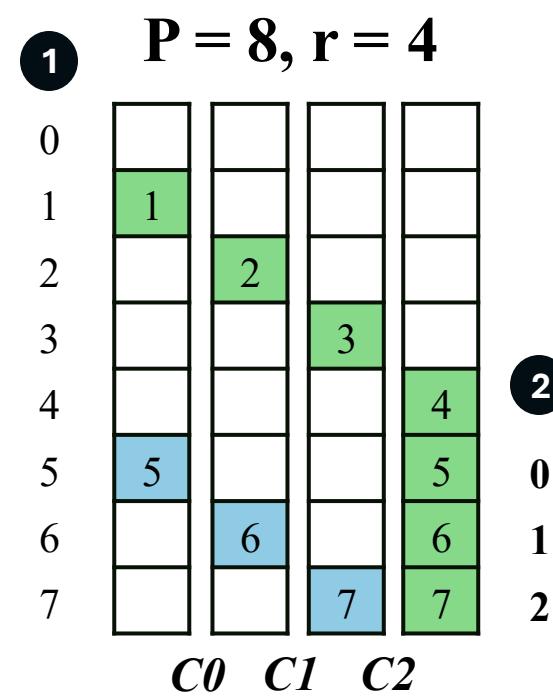
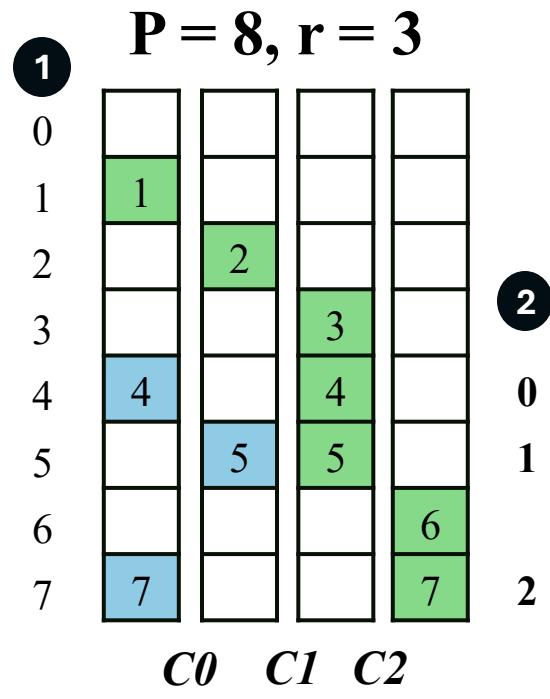


# Temporary Buffer Size Analysis

$$P = 8, r = 2$$



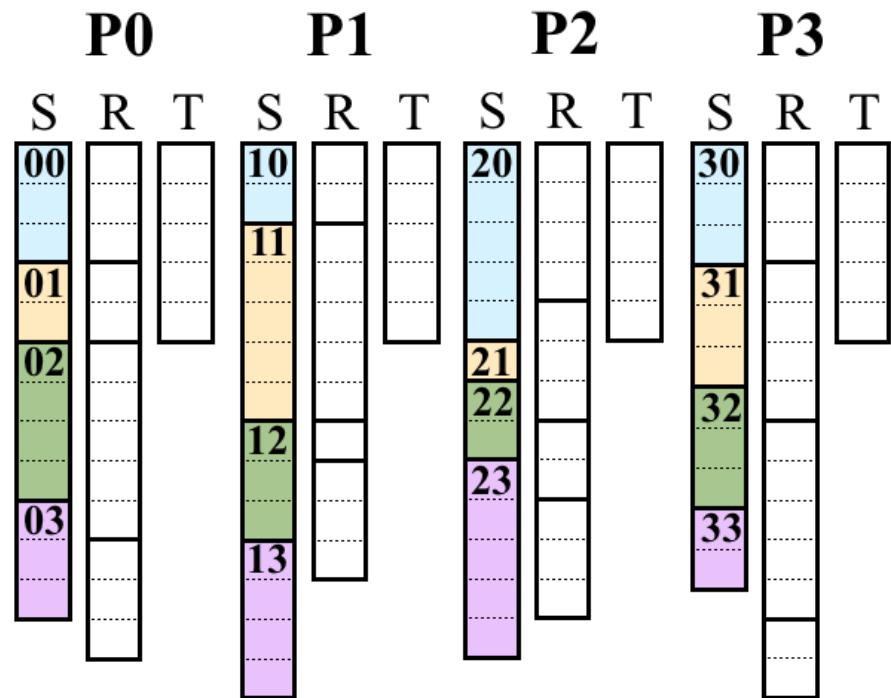
# Temporary Buffer Size Analysis



# Agenda

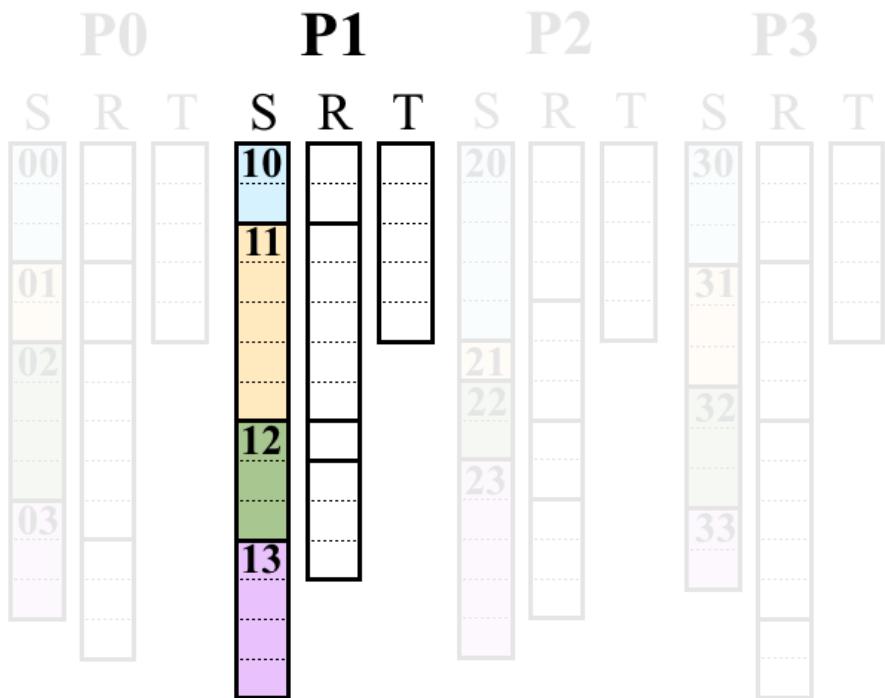
- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- **Parameterized Logarithmic Algorithm**
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

# Example of the ParLogNa with $P = 4$ and $r = 2$ .



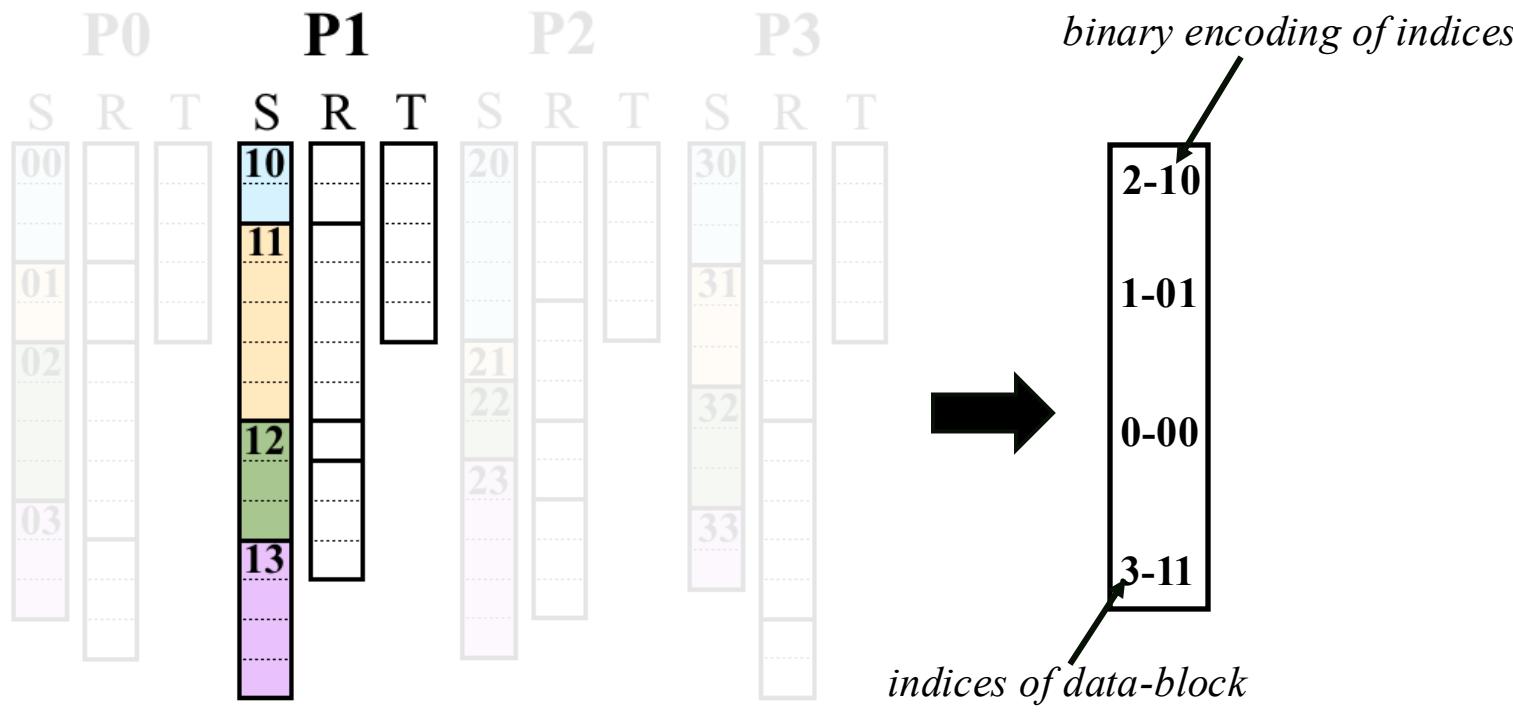
Initial data for 4 process

# Example of the ParLogNa with $P = 4$ and $r = 2$ .

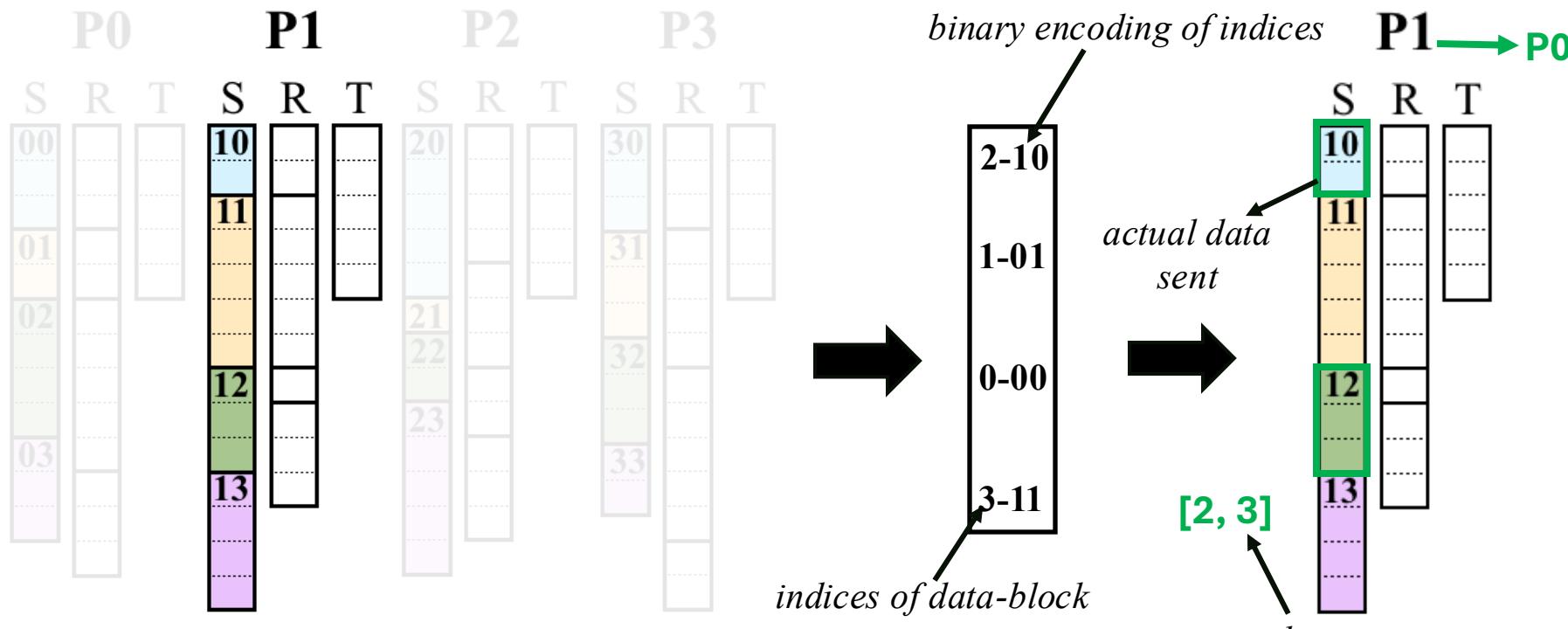


Taking P1 for example

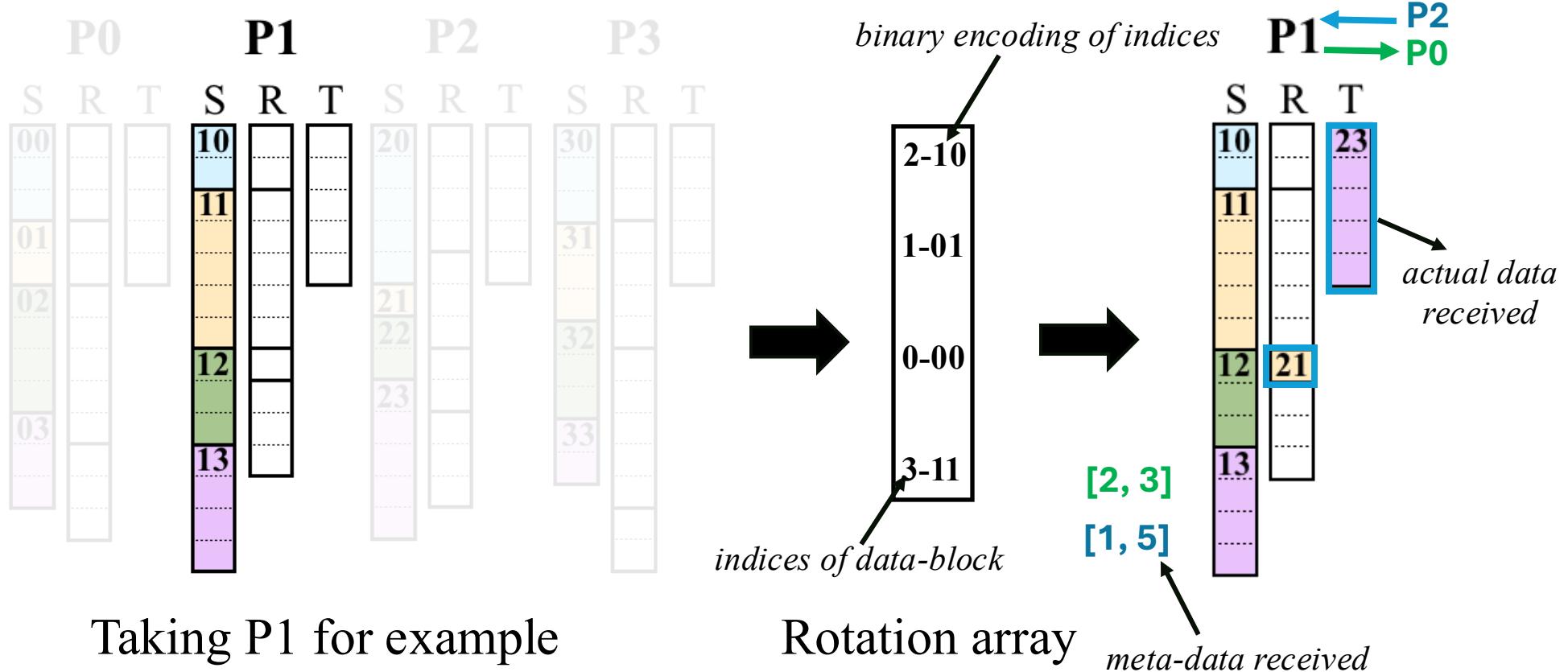
# Example of the ParLogNa: Rotation Array and Binary Encoding



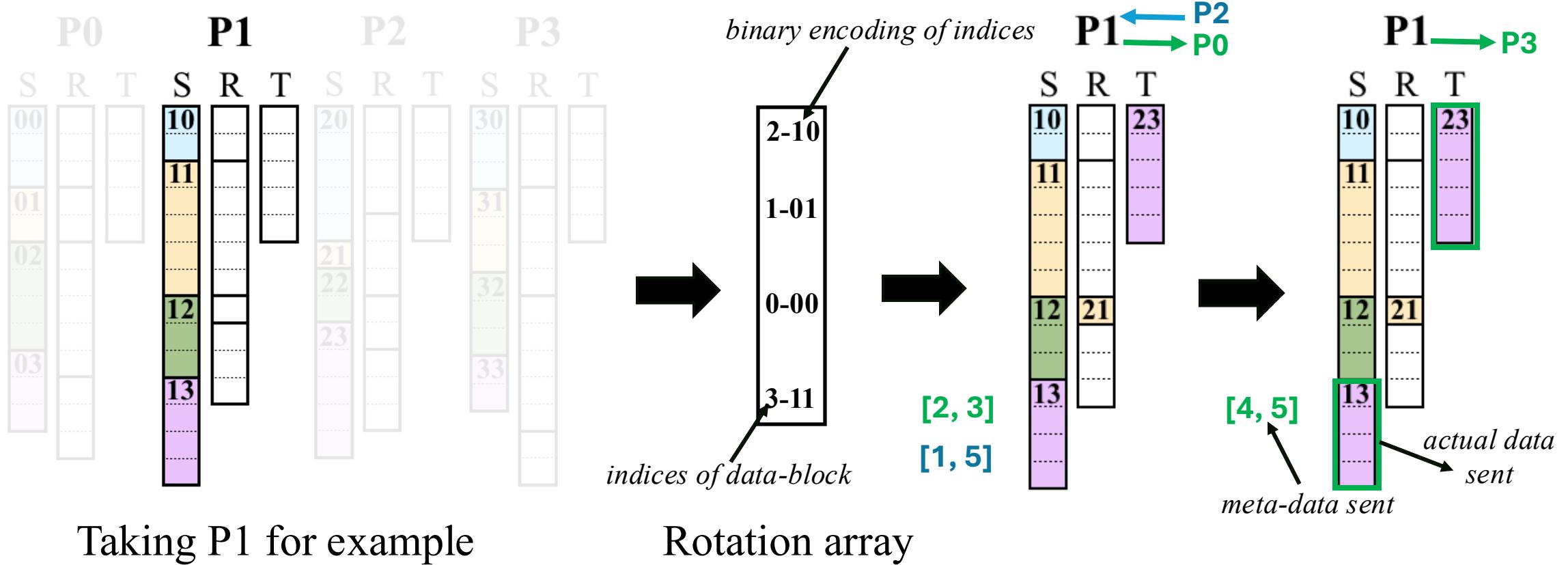
# Example of the ParLogNa: Comm Round 0



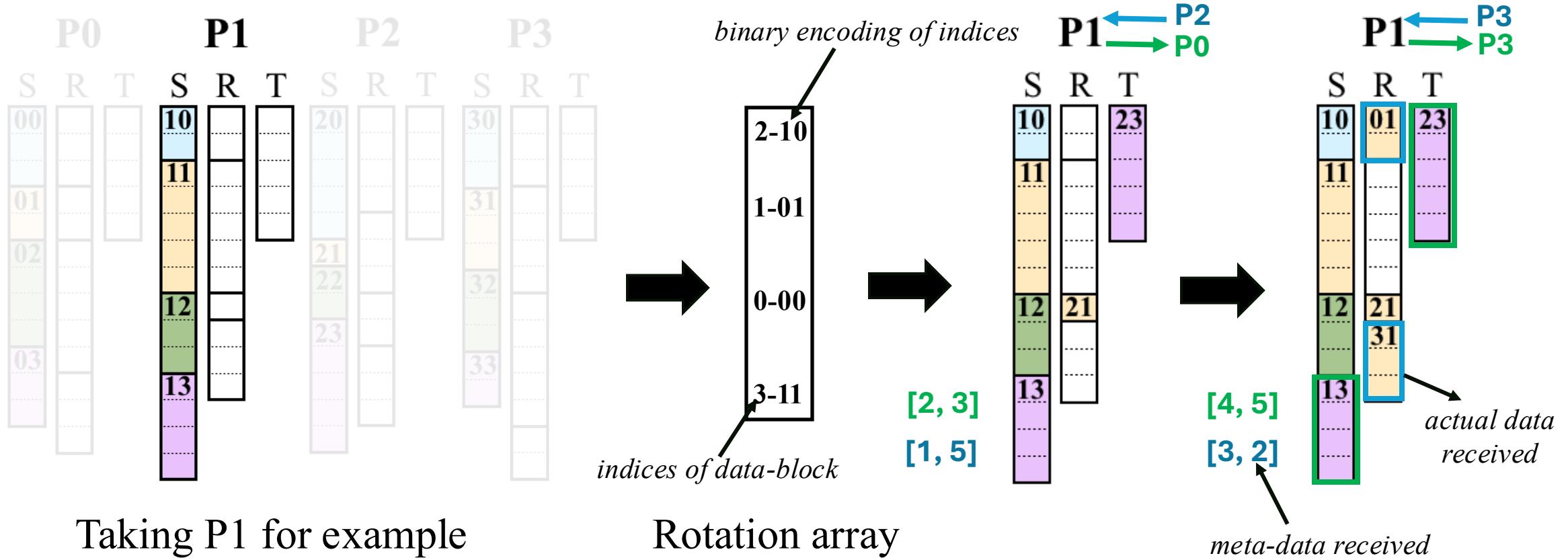
# Example of the ParLogNa: Comm Round 0



# Example of the ParLogNa: Comm Round 1



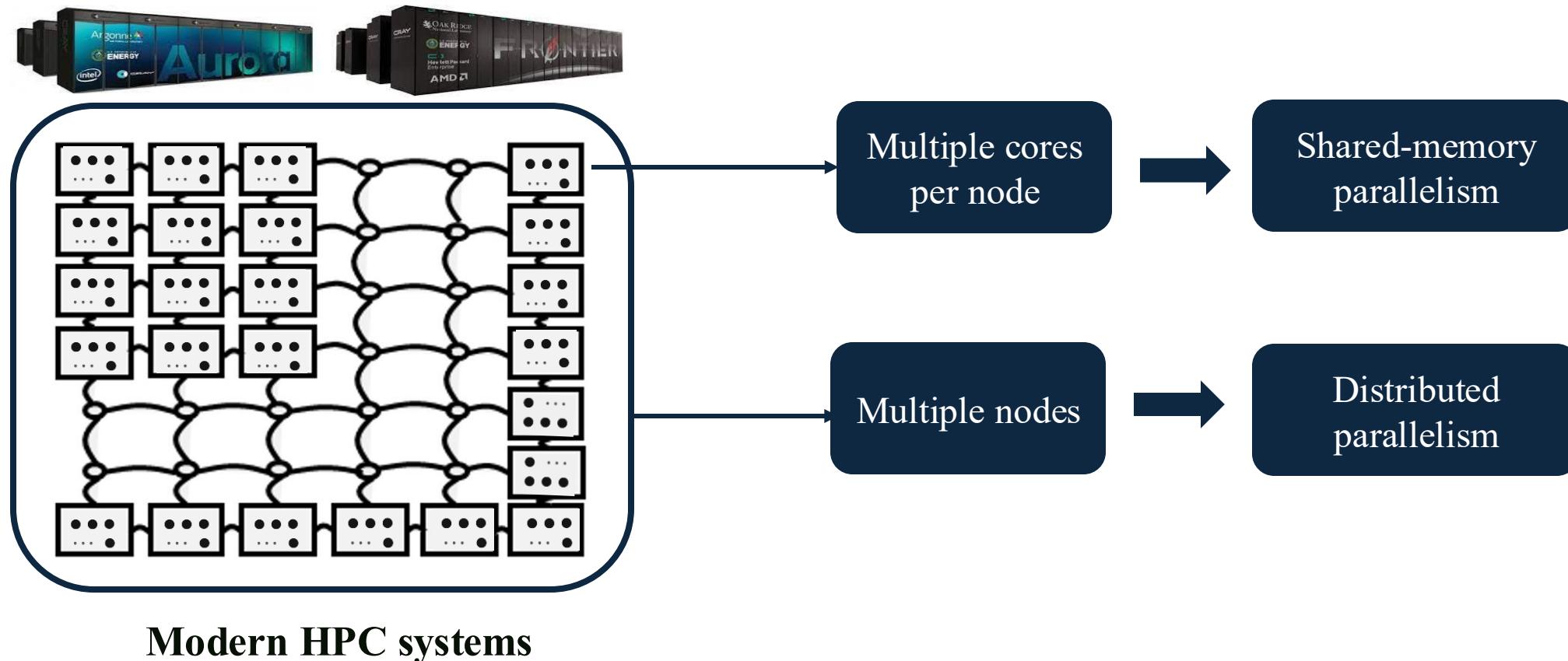
# Example of the ParLogNa: Comm Round 1



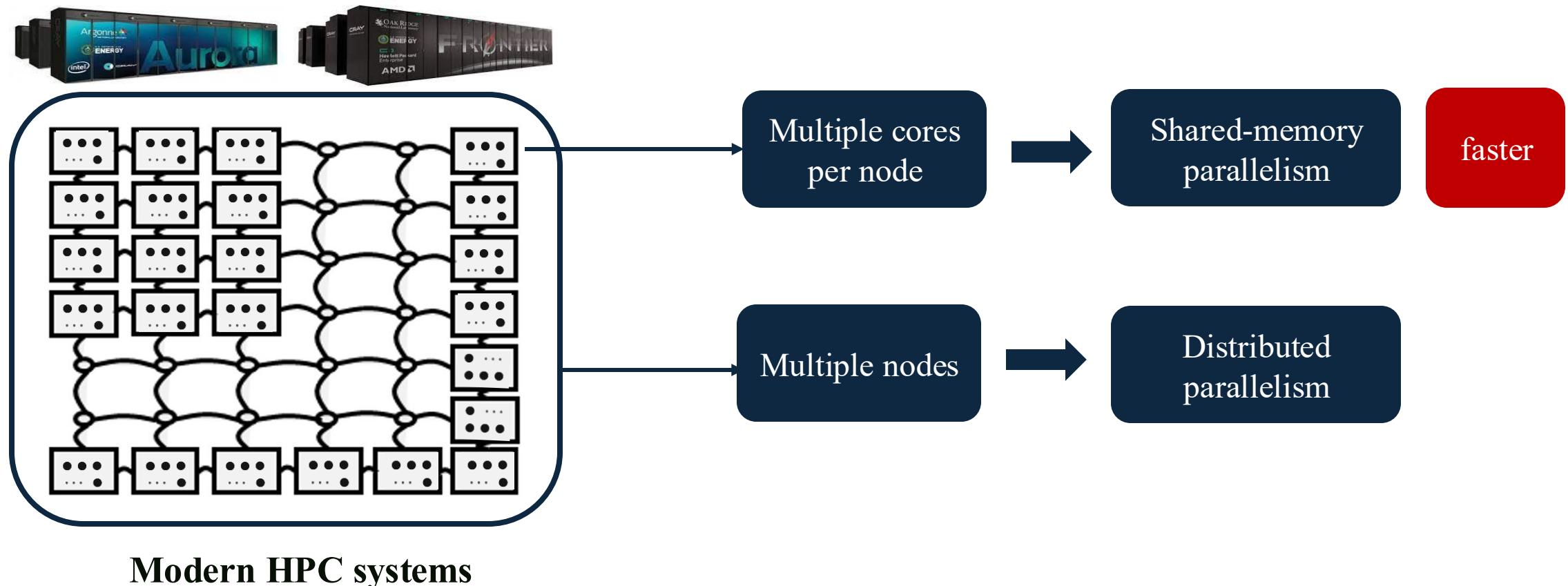
# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- Parameterized Logarithmic Algorithm
- **Parameterized Linear Algorithm**
- Evaluation
- Application
- Conclusion

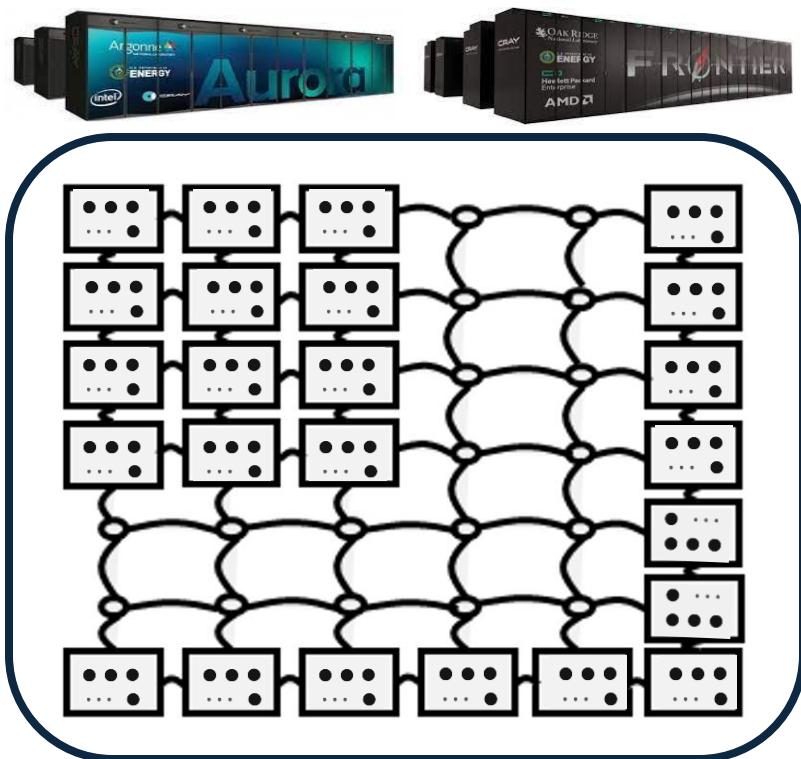
# Architecture of Modern HPC Systems



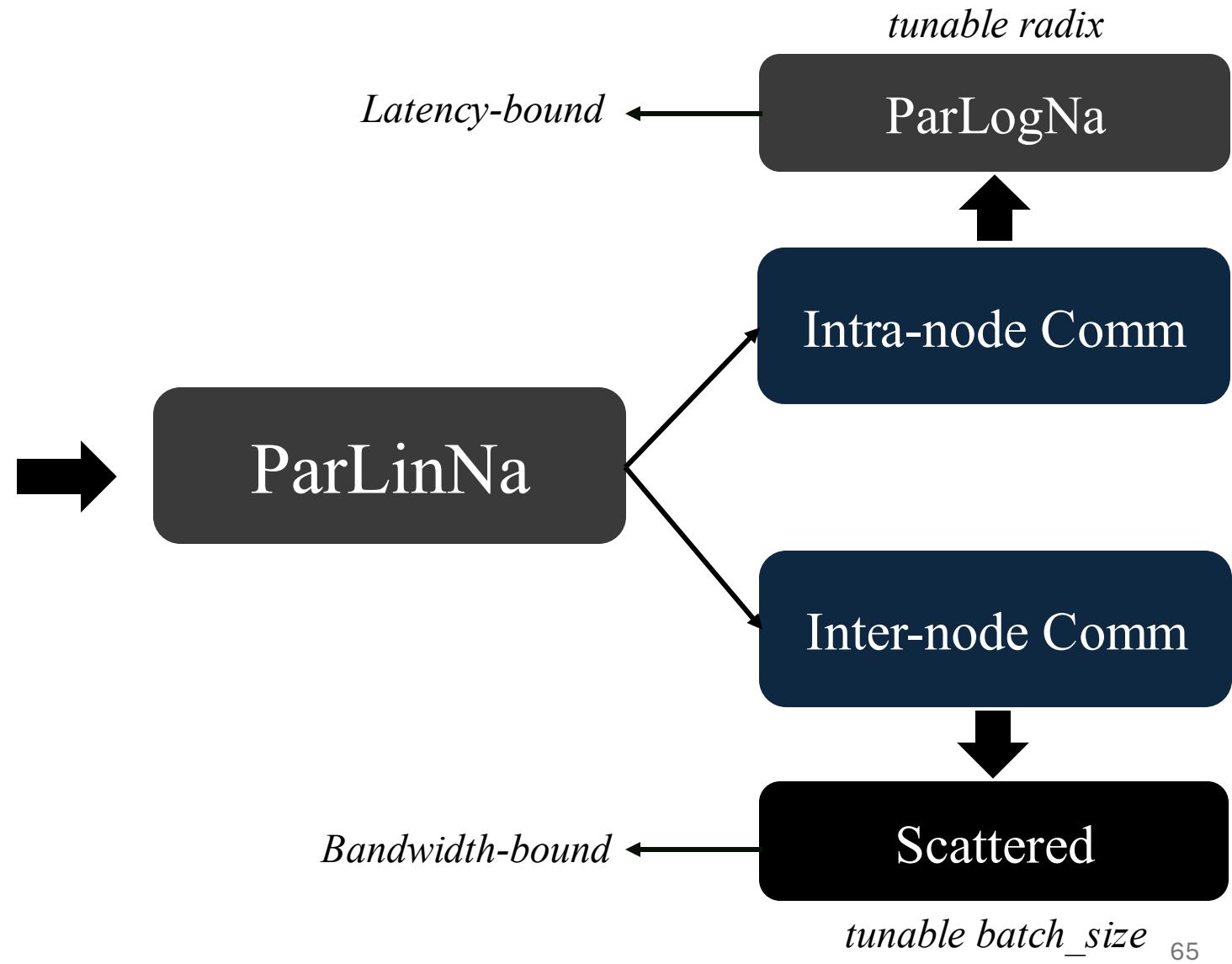
# Architecture of Modern HPC Systems



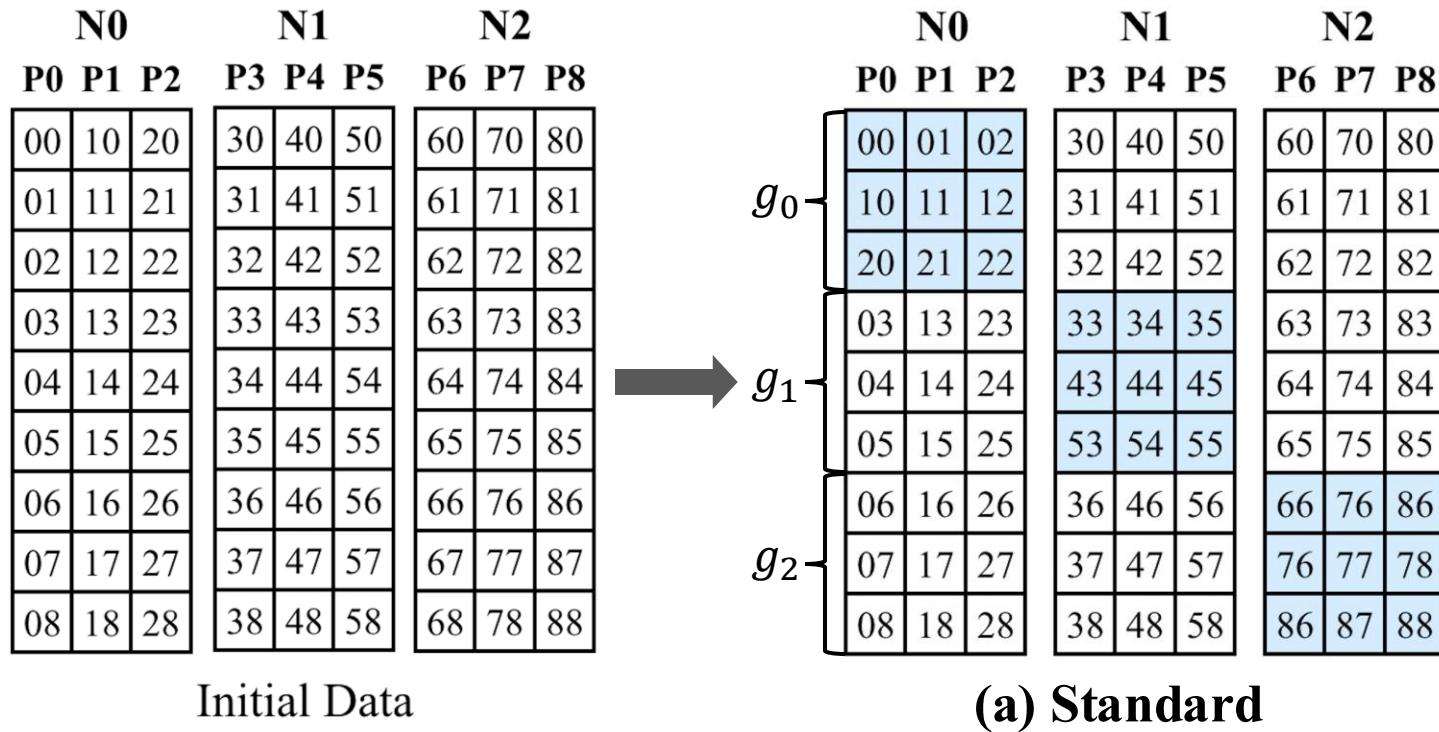
# Parameterized Linear Non-uniform All-to-all (ParLinNa)



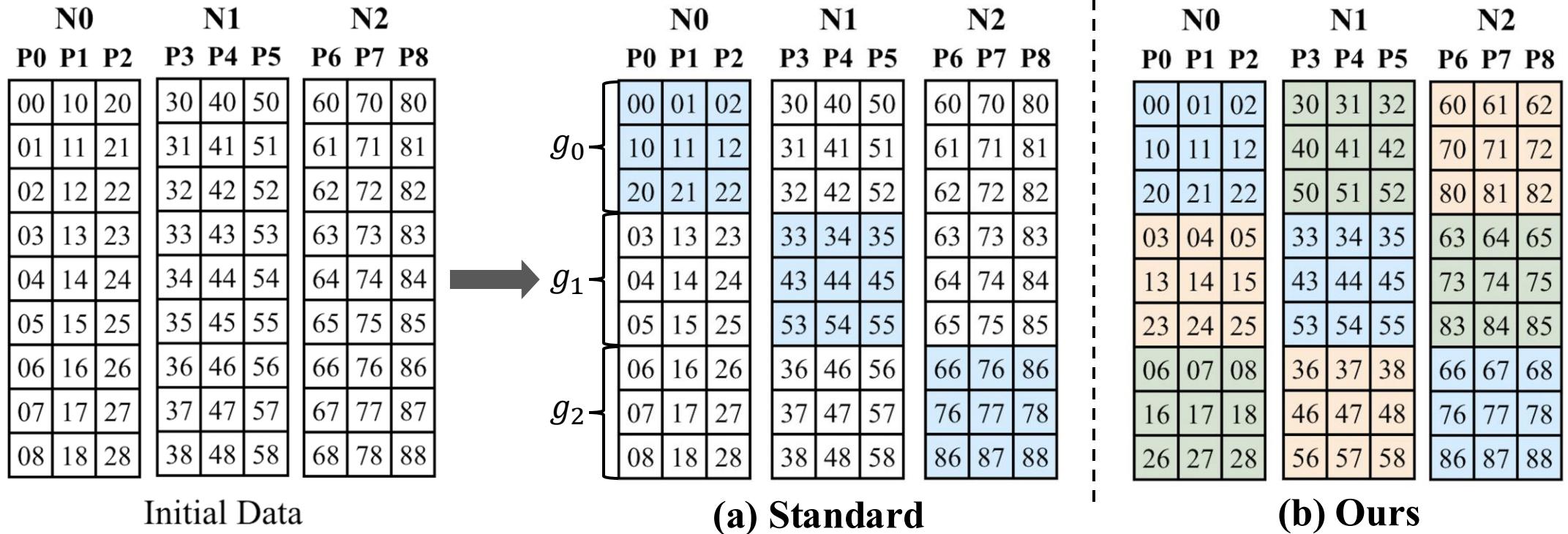
Modern HPC systems



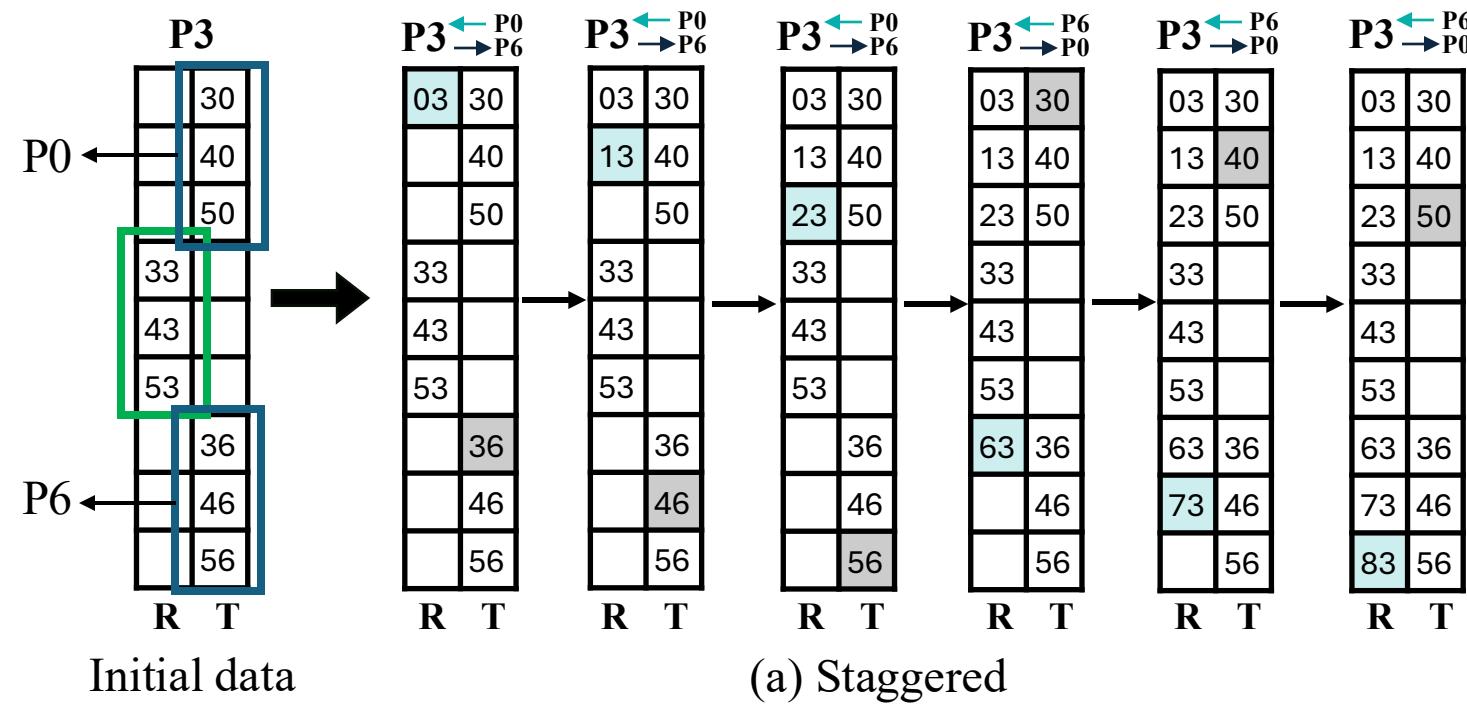
# Intra-node Communication: Two Strategies



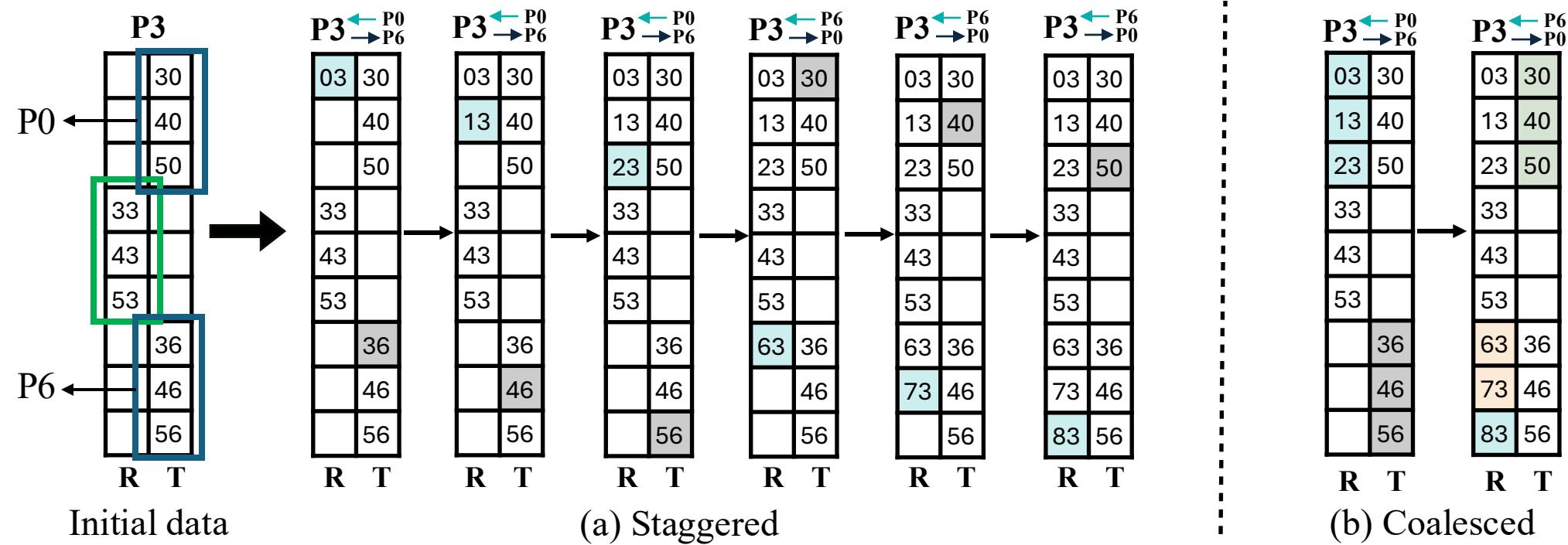
# Intra-node Communication: Two Strategies



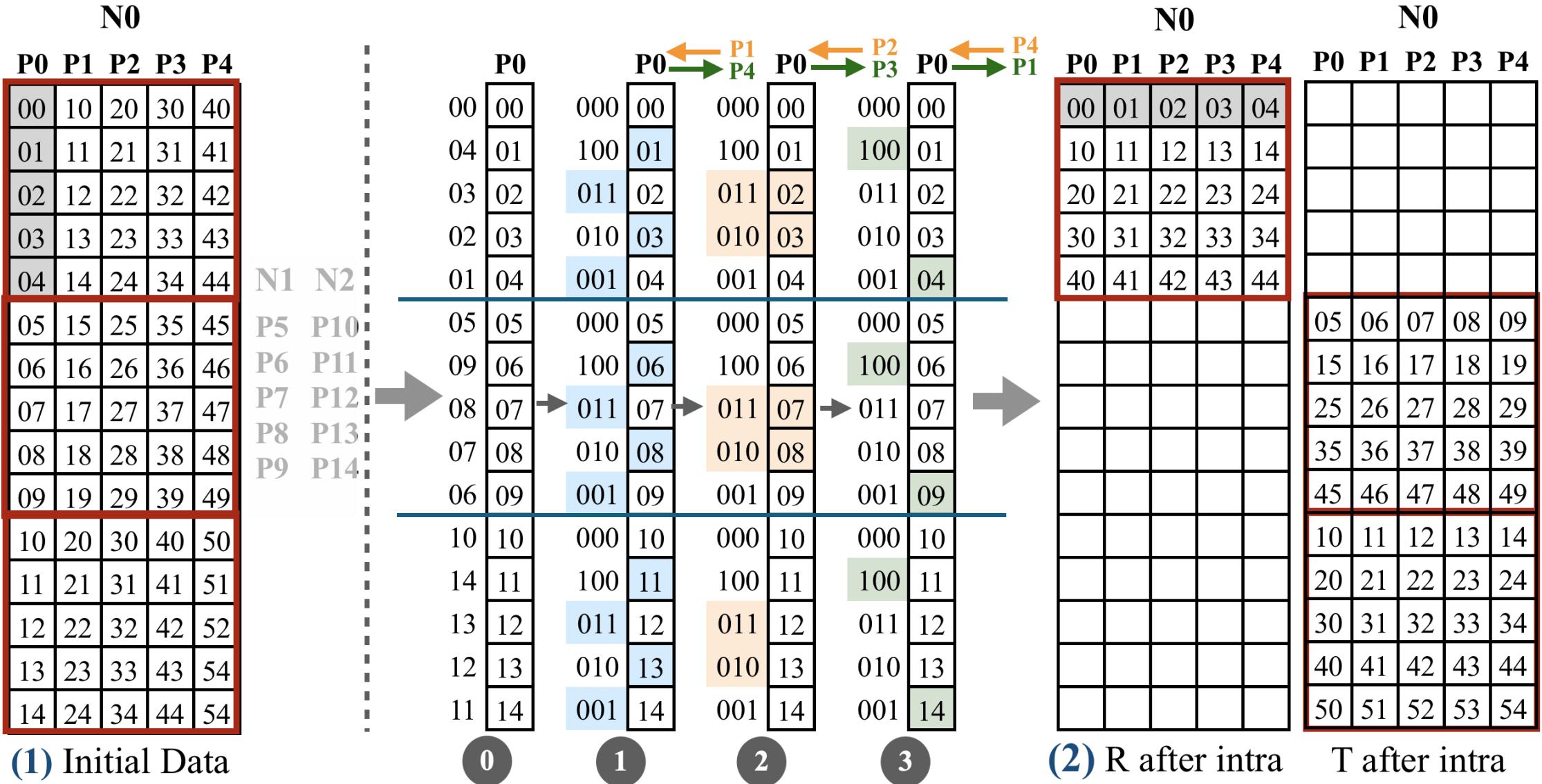
# Inter-node Communication: Two Patterns



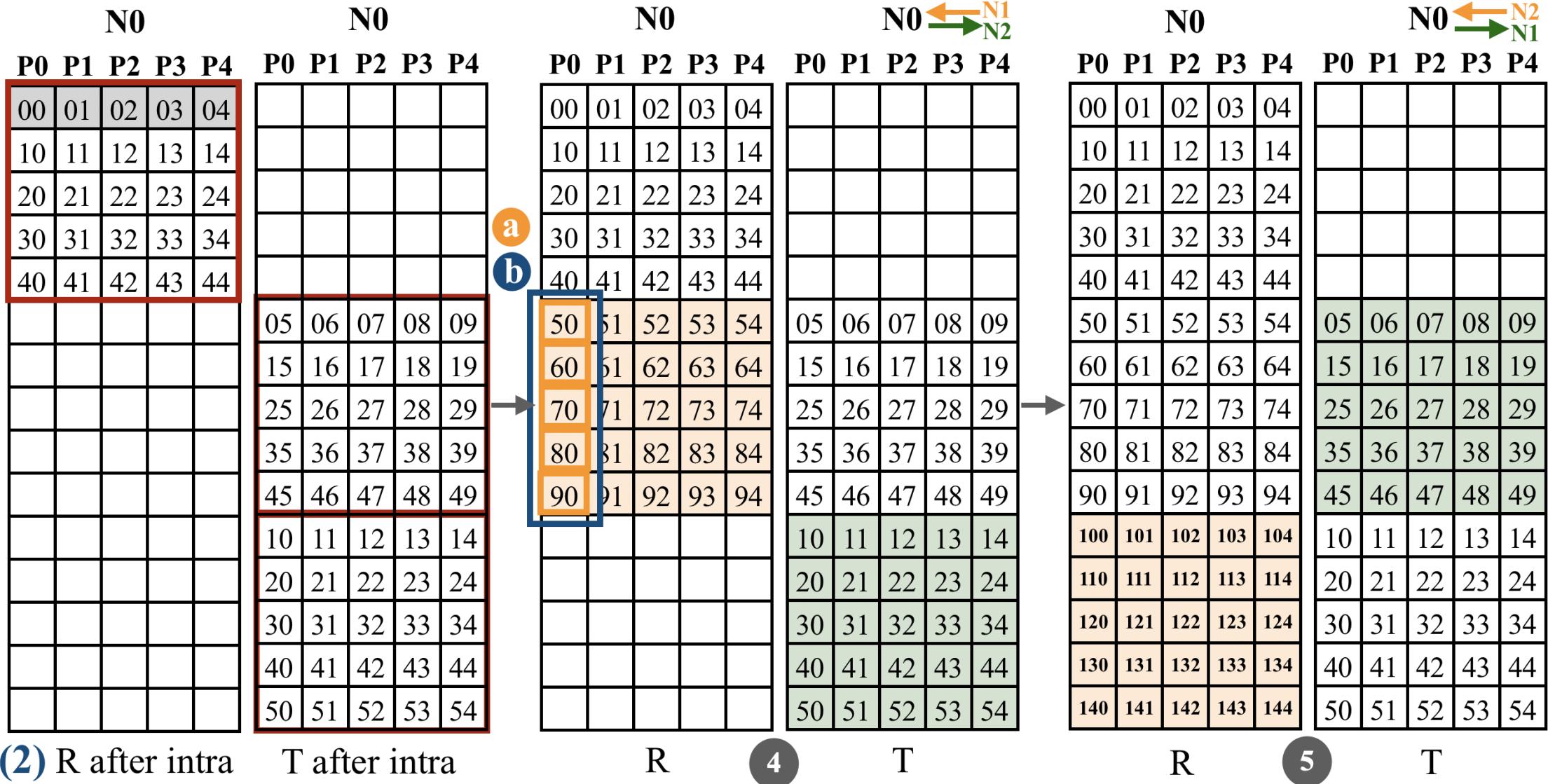
# Inter-node Communication: Two Patterns



# Intra-node Communication Example



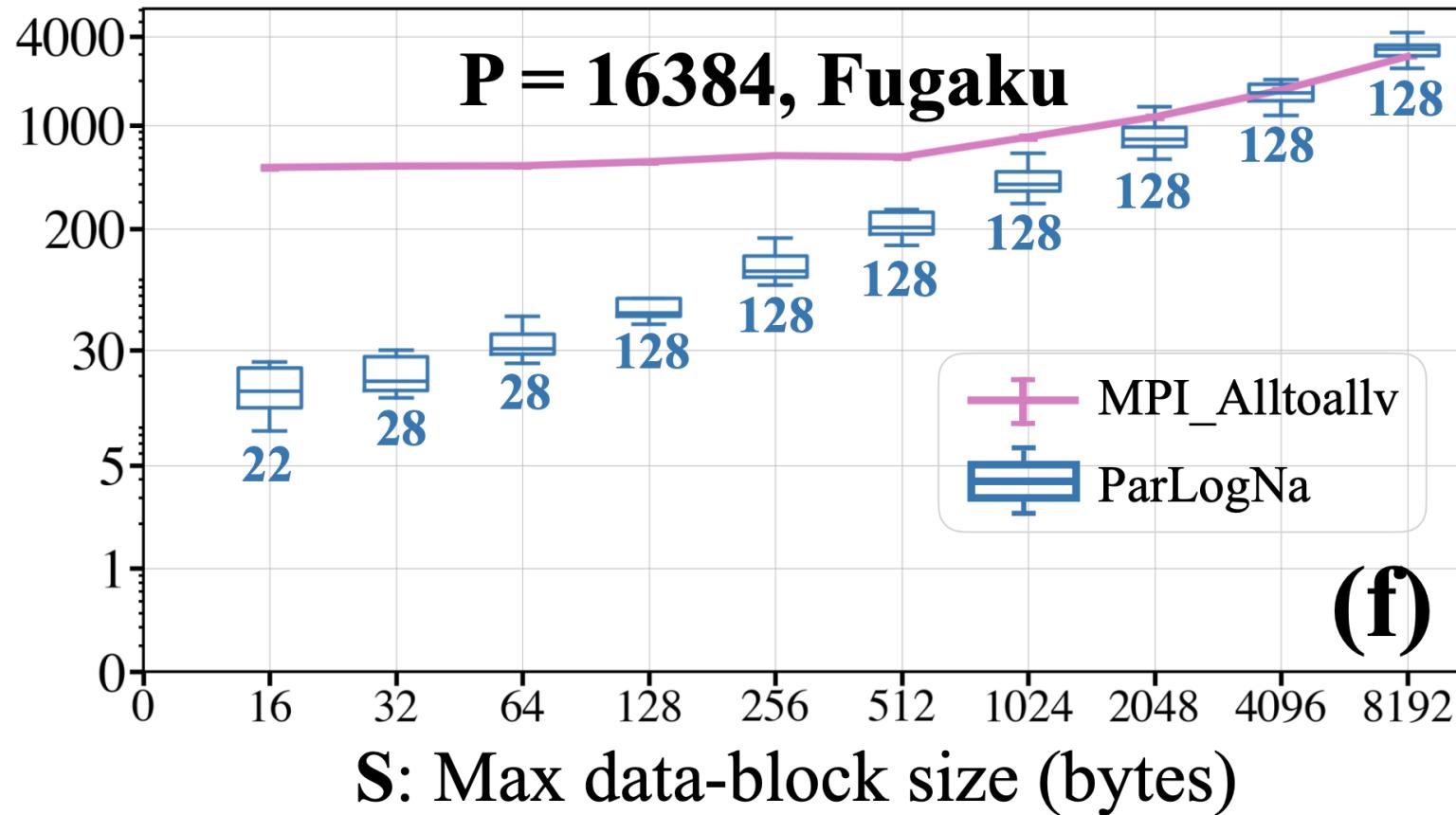
# Inter-node Communication Example



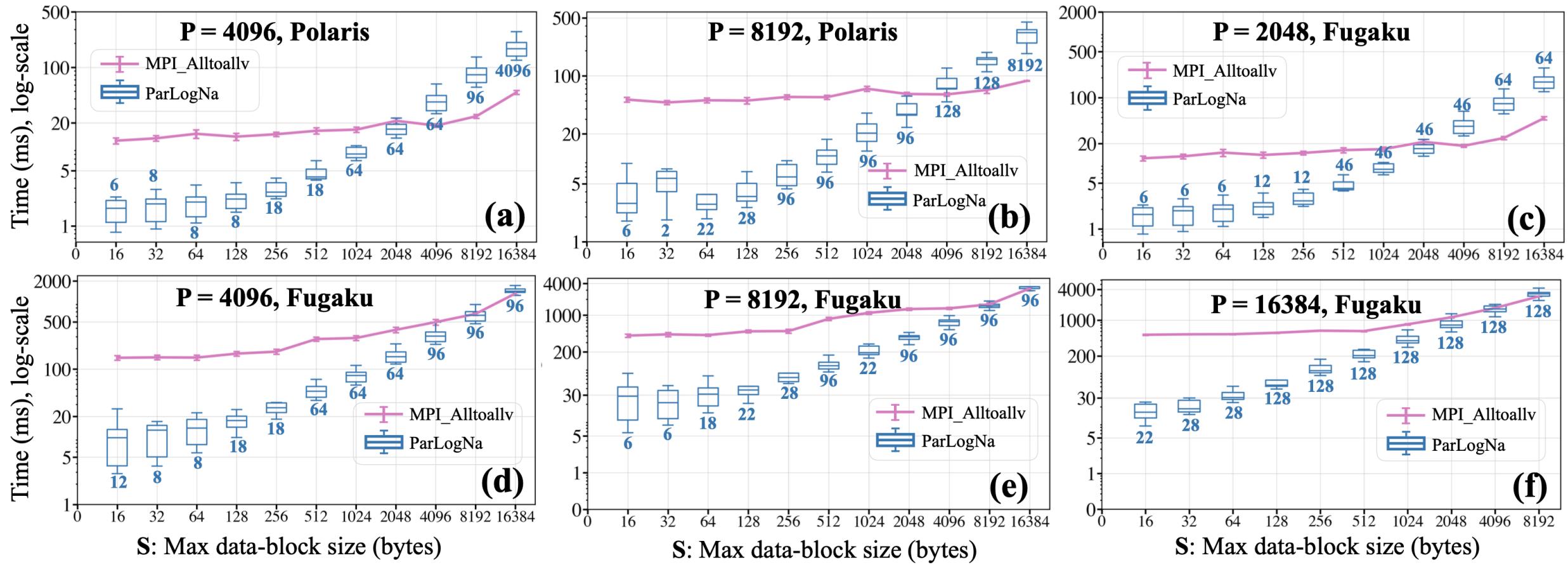
# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- Parameterized Logarithmic Algorithm
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

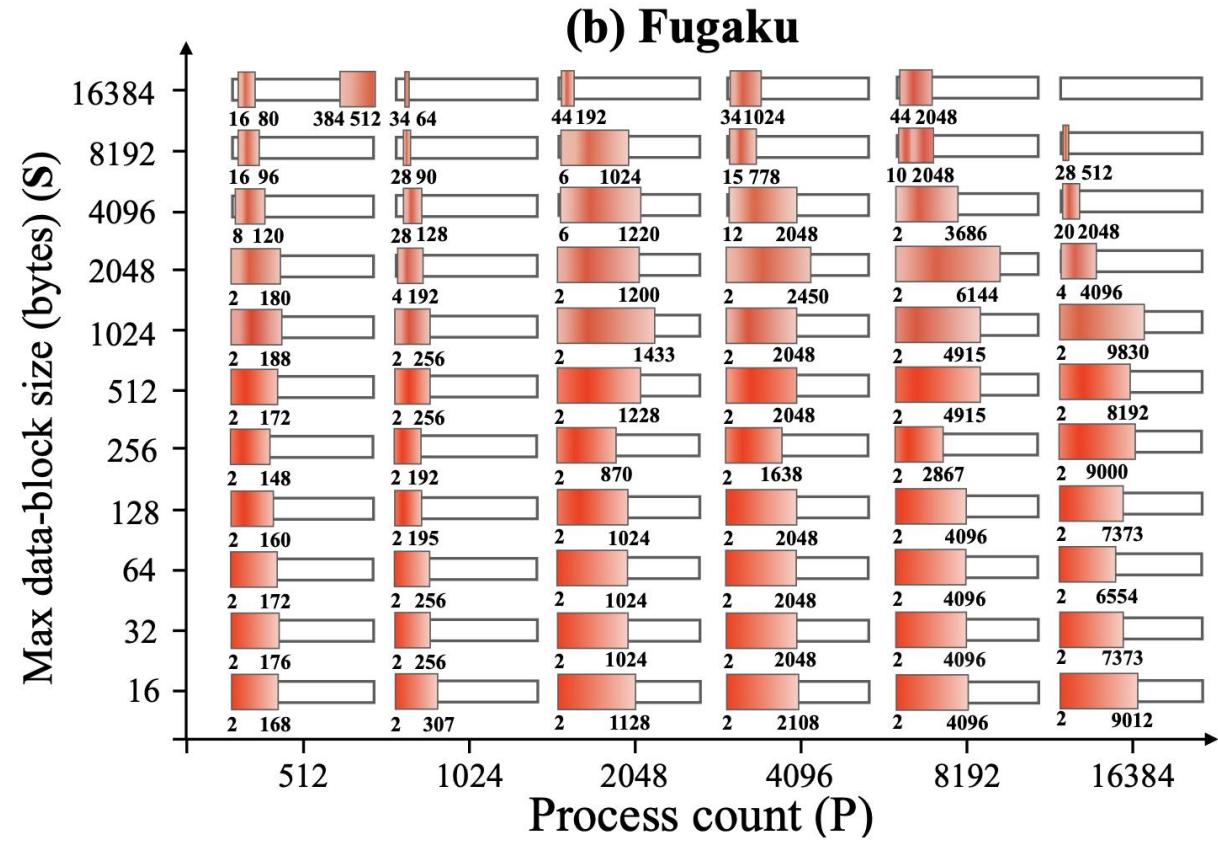
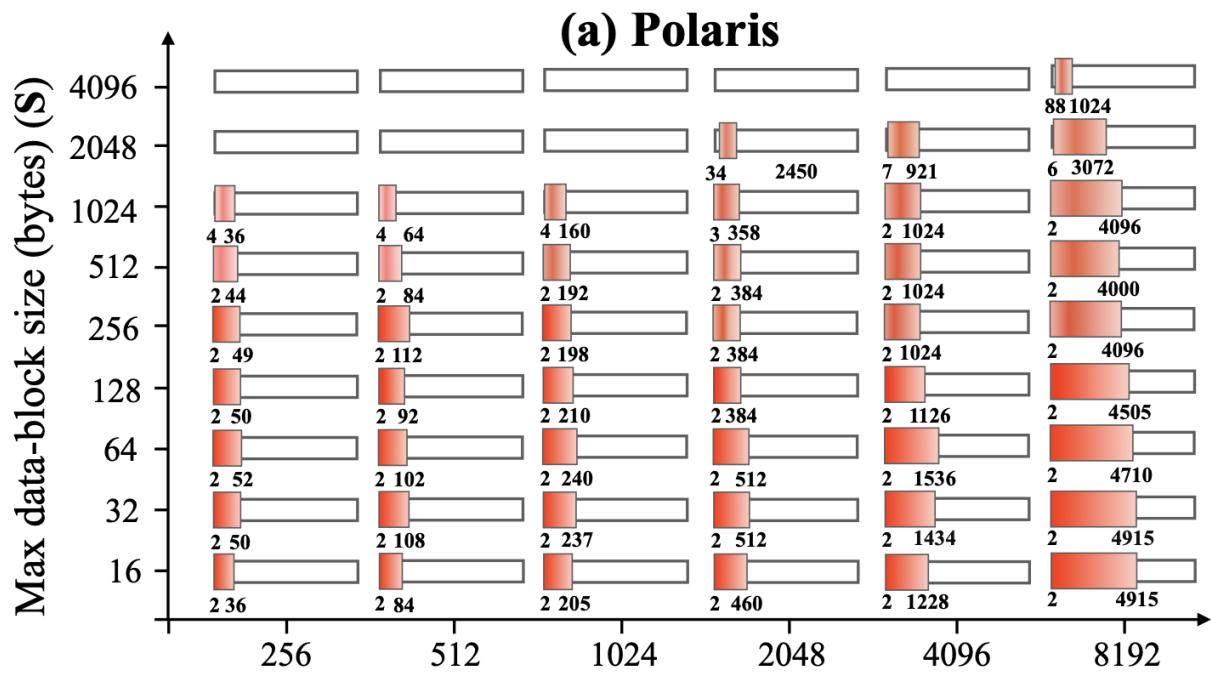
# Comparing ParLogNa with MPI\_Alltoallv



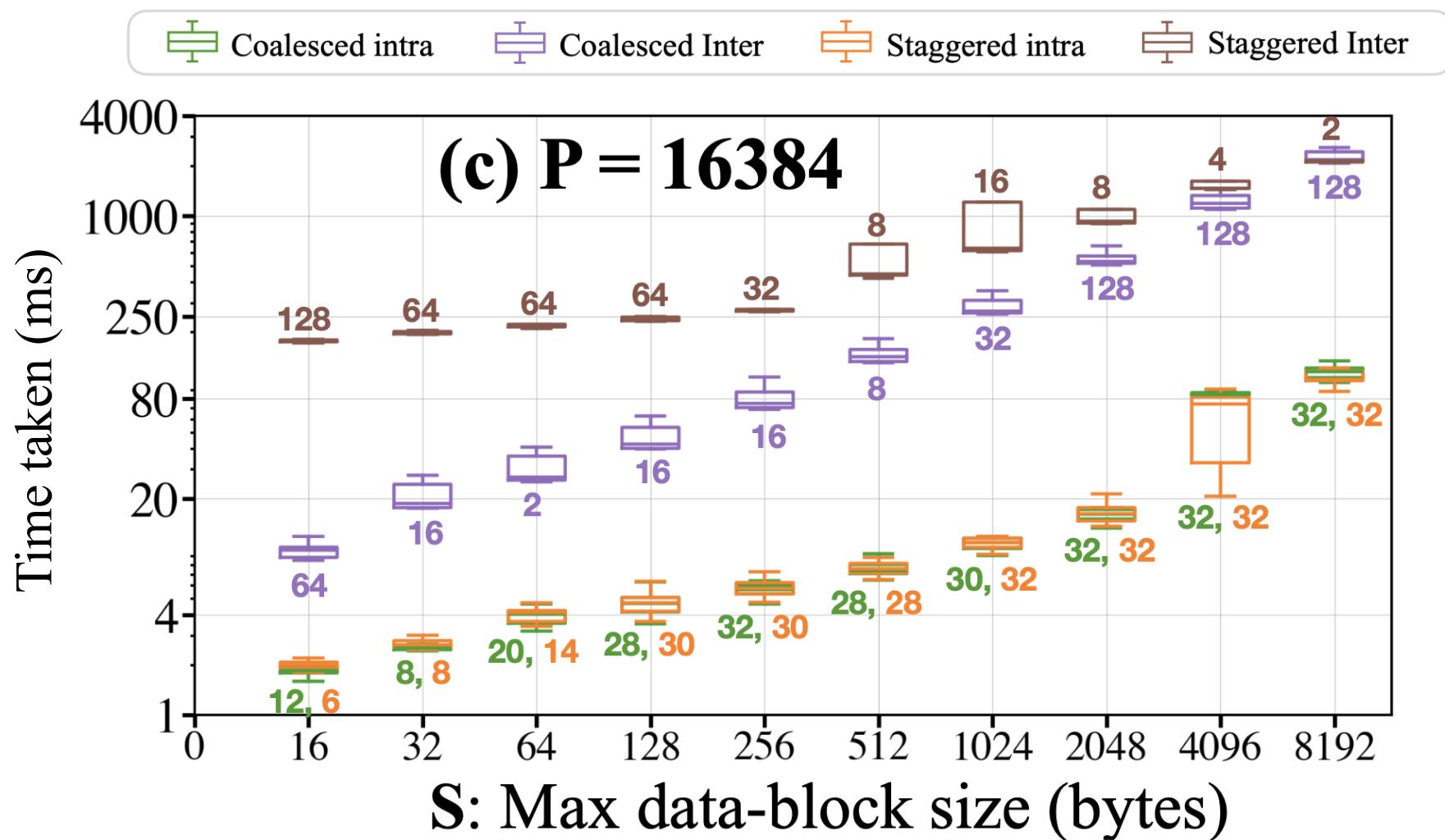
# Comparing ParLogNa with MPI\_Alltoallv



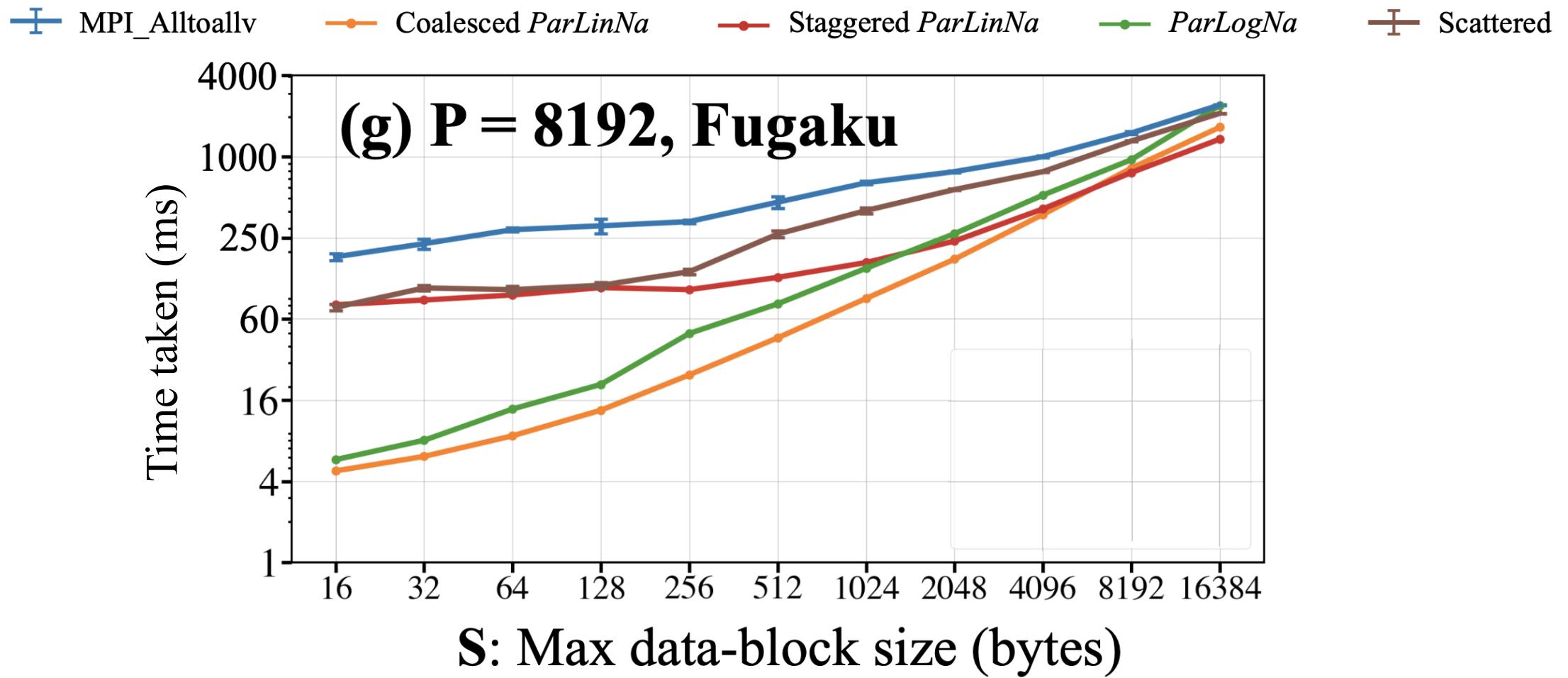
# Ranges of radix where ParLogNa outperforms MPI\_Alltoallv



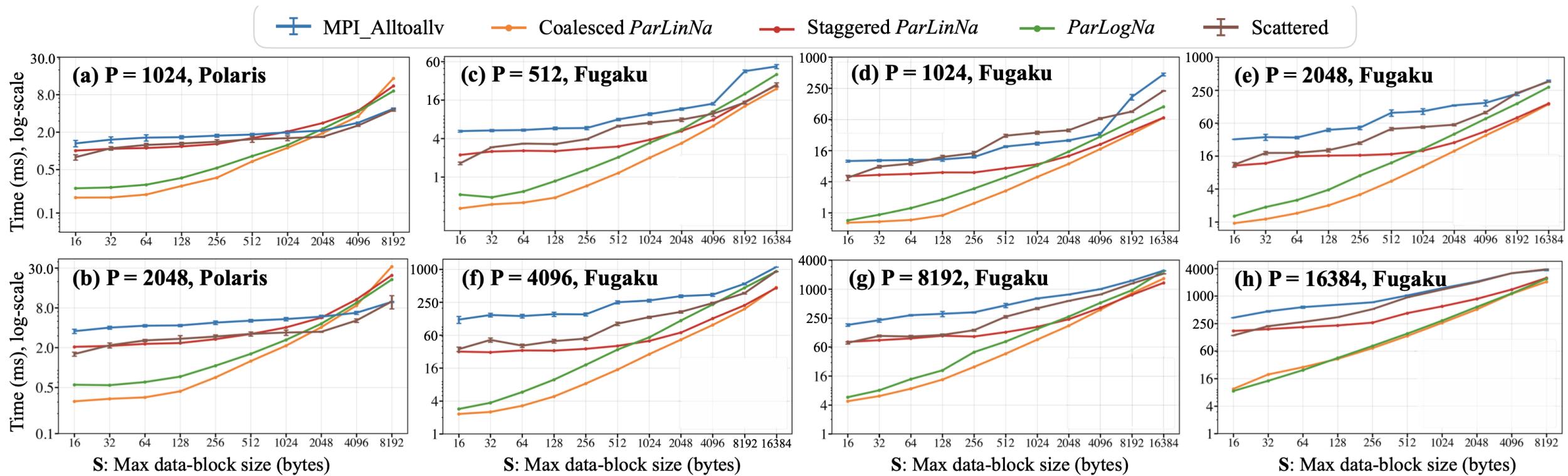
# Comparing Coalesced and Staggered ParLinNa



# Comparing Proposed Algorithms with the Top-performing Benchmark



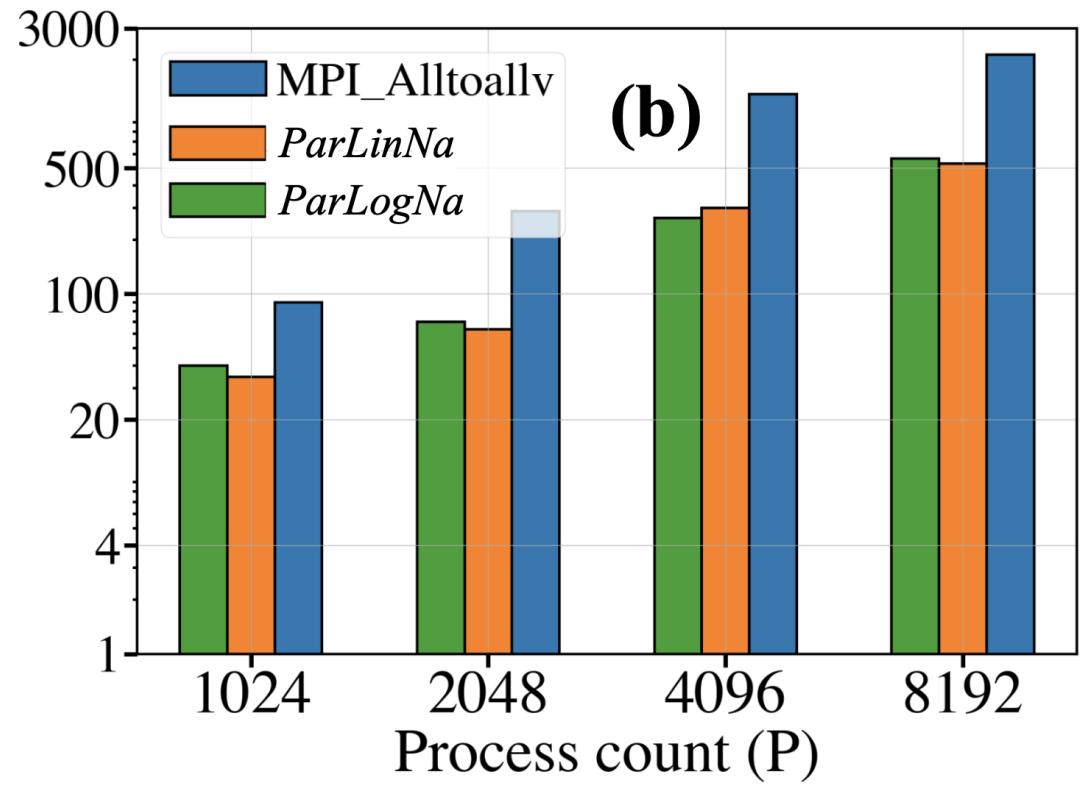
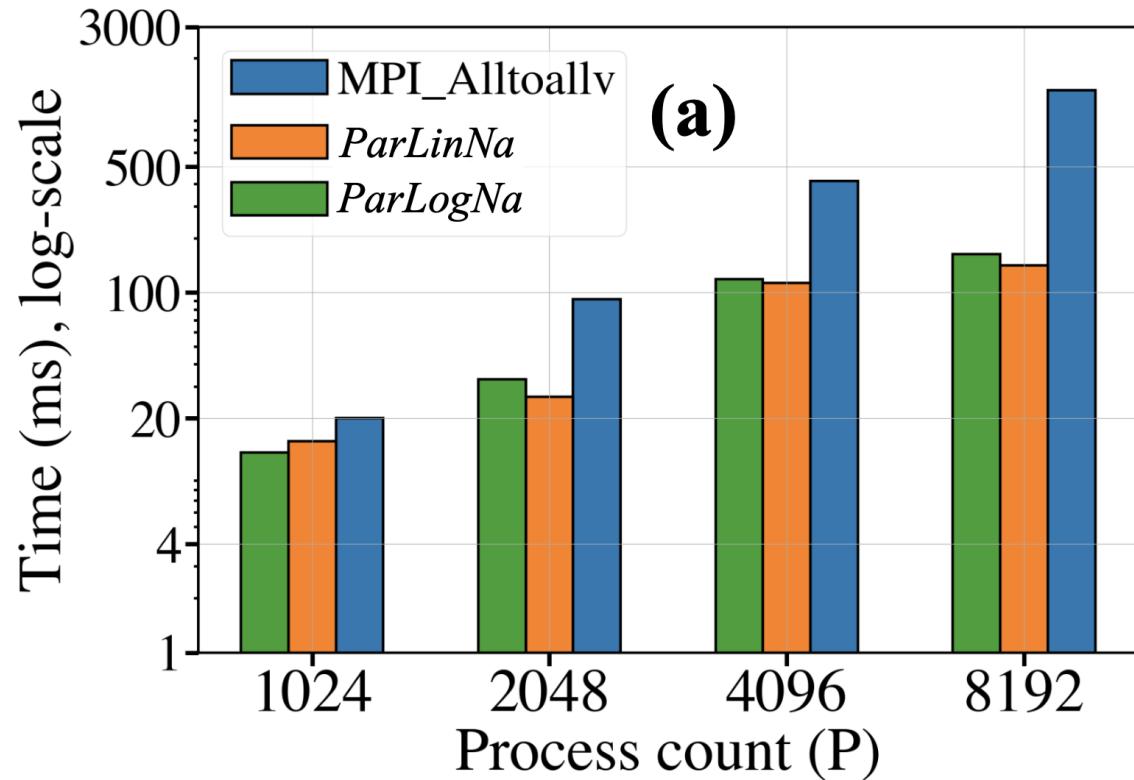
# Comparing Proposed Algorithms with the Top-performing Benchmark



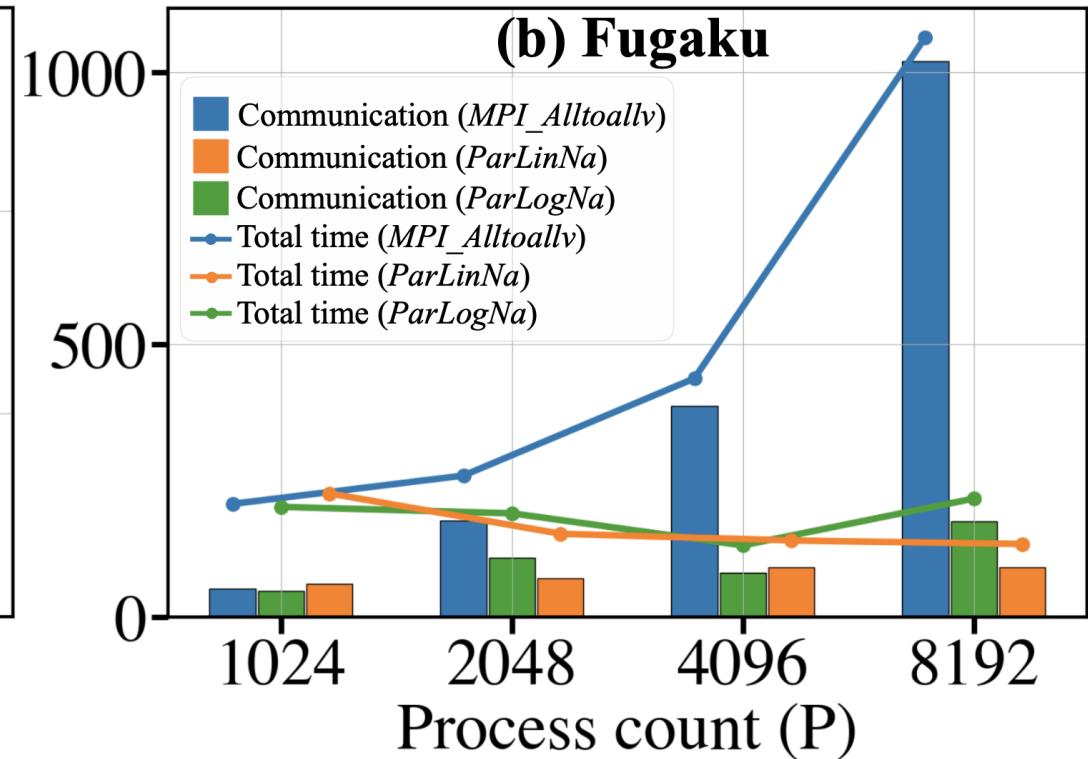
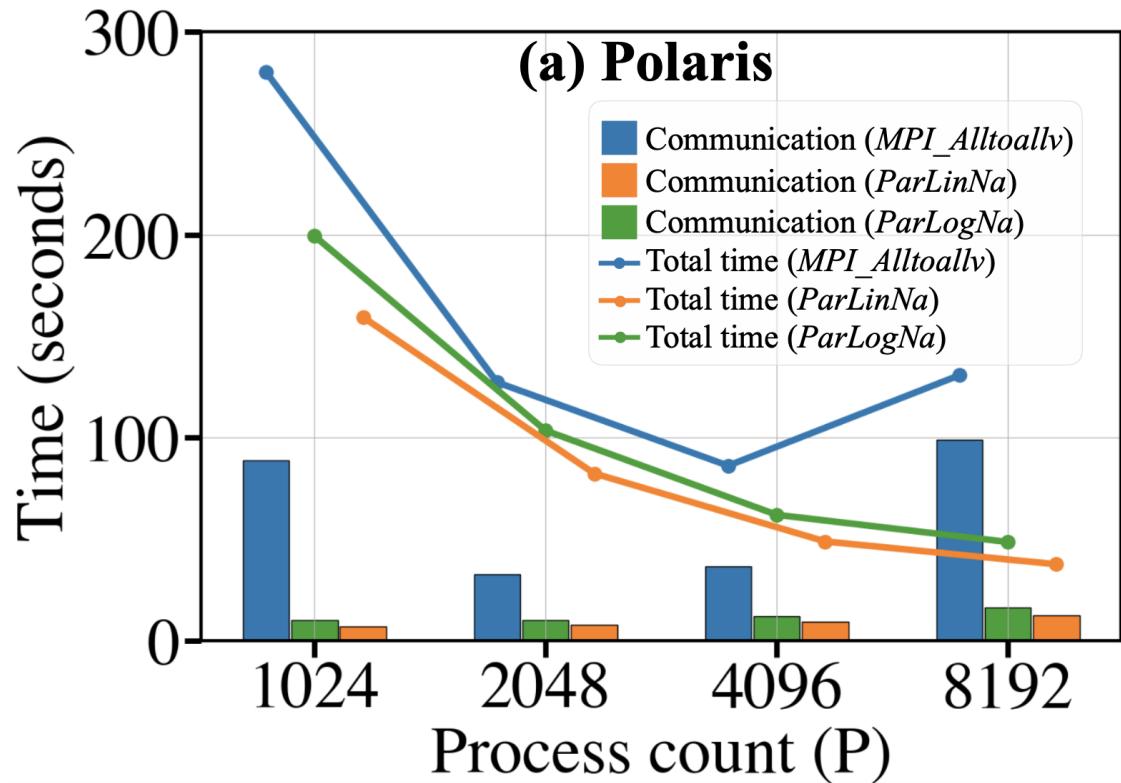
# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- Parameterized Logarithmic Algorithm
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

# Performance of Applying Our Algorithms to FFT



# Performance of Applying Our Algorithms to Path Finding



# Agenda

- Introduction and Background
  - Inter-process Communication
  - All-to-all Data Exchange
  - Standard MPI All-to-all Implementations
- Challenges and Solutions
  - Tunability of Performance
  - Logarithmic Non-uniform All-to-all
- Parameterized Logarithmic Algorithm
- Parameterized Linear Algorithm
- Evaluation
- Application
- Conclusion

# Conclusion

- We have presented algorithms (ParLogNa and ParLinNa) that can improve an important class of collective functions.
- The work can help vendors further optimize their collective routines.
- The work can have a direct impact on a range of applications that uses all to all communication.

# Thank you!