

CS168, Fall 2024

**Beyond Client-Server (part 2)**  
**Collective Communication**

Sylvia Ratnasamy

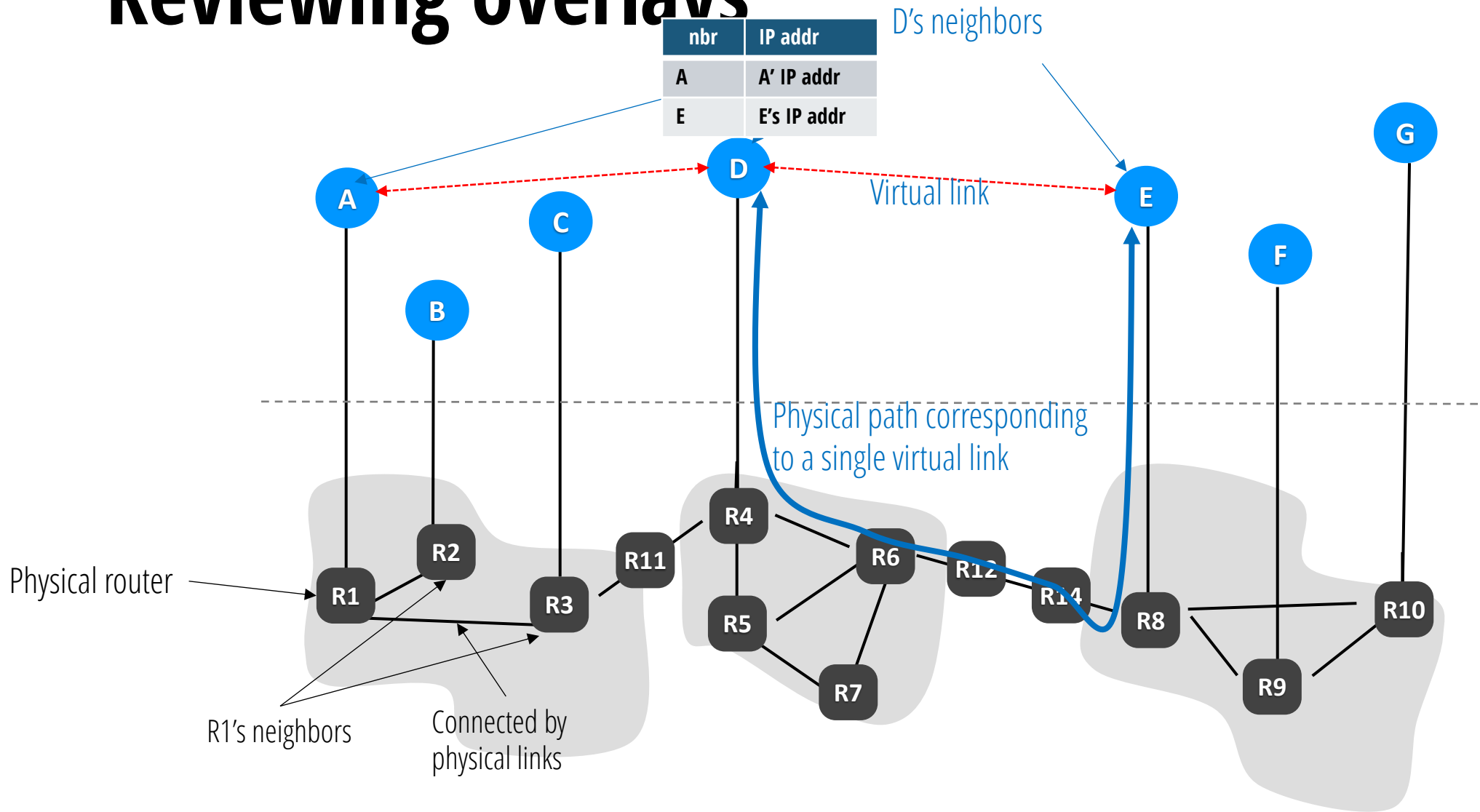
# Recall

- Last time: apps based on group vs. client-server interactions
  - Multicast as a communication primitive
  - Implemented in the network (L3) or as an overlay (L7)
  - Our context: the Internet at large
- Today: distributed training as our app
  - Also based on groups of hosts that communicate
  - But very different context: a single datacenter, different goals, *etc.*
  - Leading to new communication primitives: **collectives**

# Outline

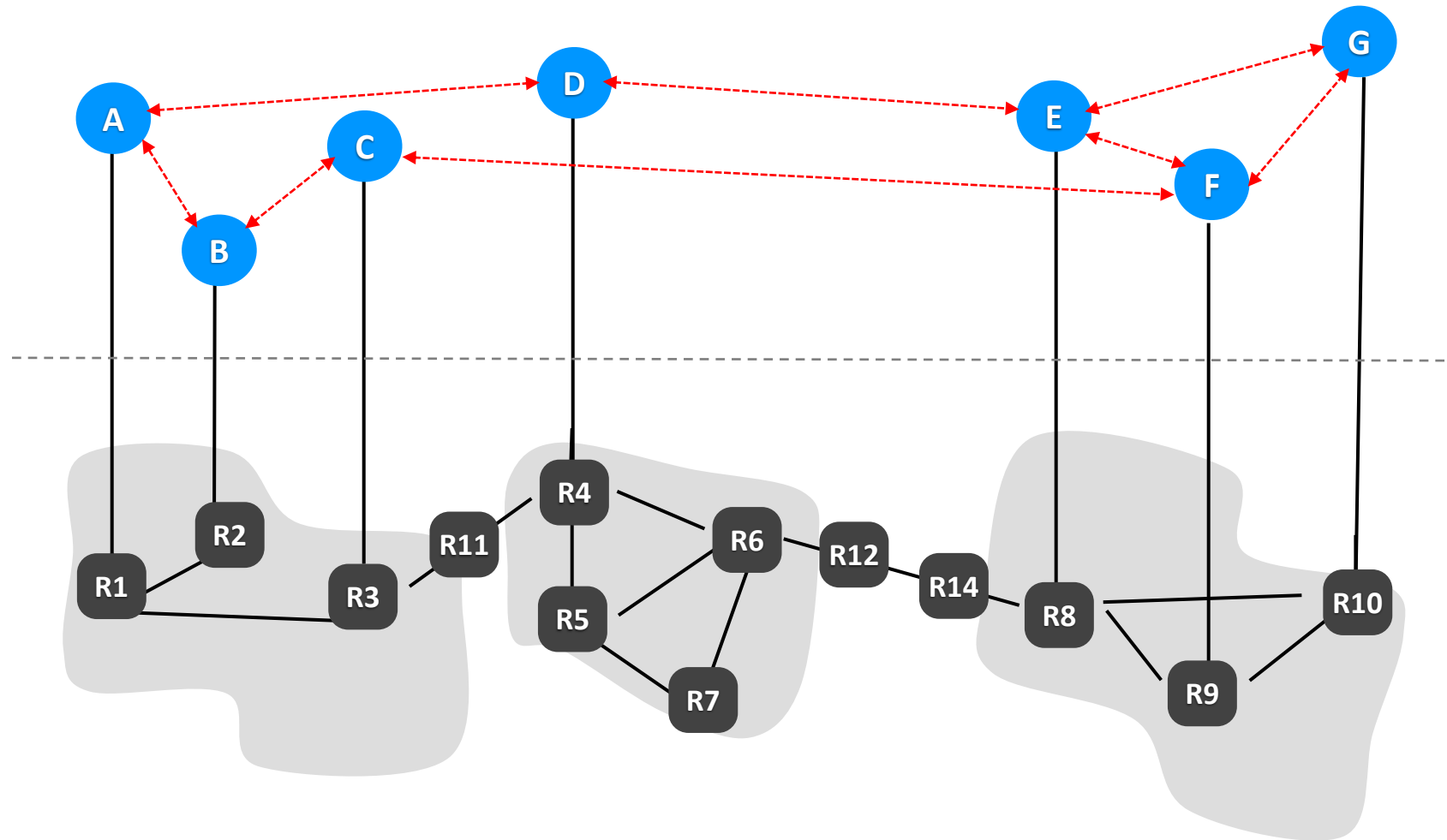
- Quick review: Overlays
- Context on distributed training: app and infrastructure
- Communication collectives: definitions and the function they provide
- How collectives are implemented
  - Focus on one collective: AllReduce
  - Implementation at overlay and underlay level

# Reviewing overlays



**A virtual network (the "overlay") on top of an underlying physical network (the "underlay")**

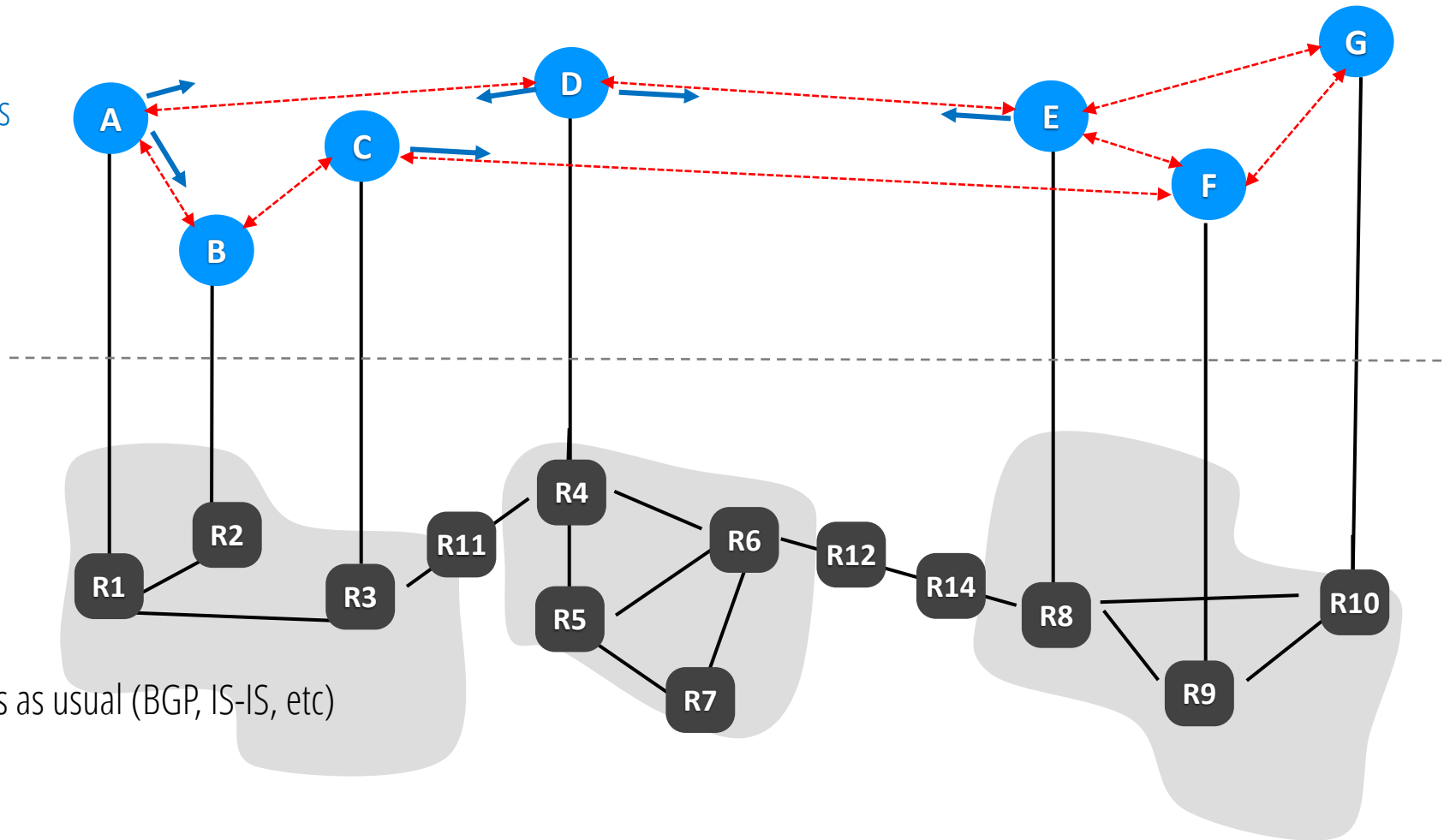
# Reviewing overlays



**A virtual network (the "overlay") on top of an underlying physical network (the "underlay")**

# Reviewing overlays

Run a routing protocol between virtual routers



Run routing protocols as usual (BGP, IS-IS, etc)

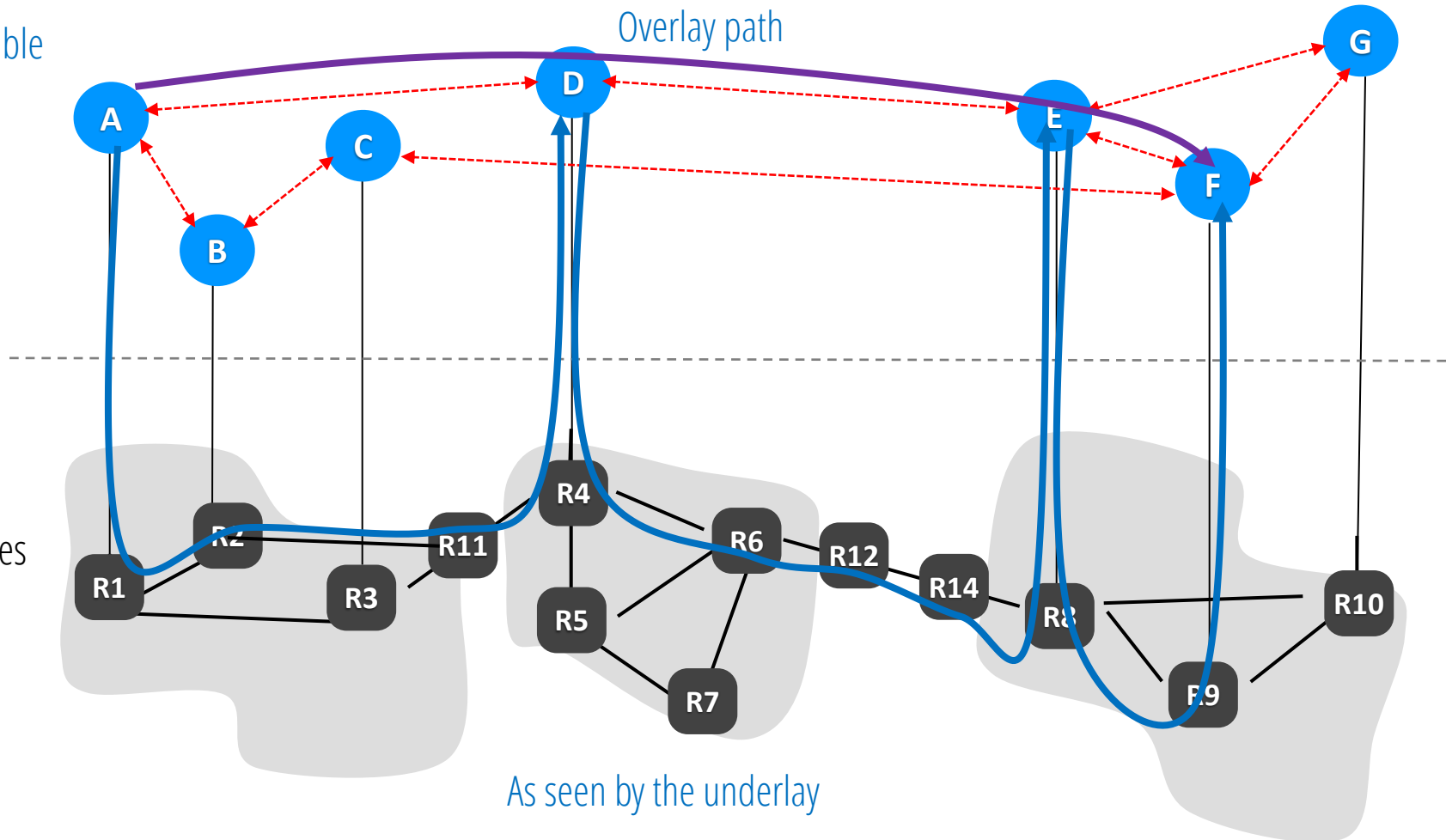
# Reviewing overlays

Overlay routing table

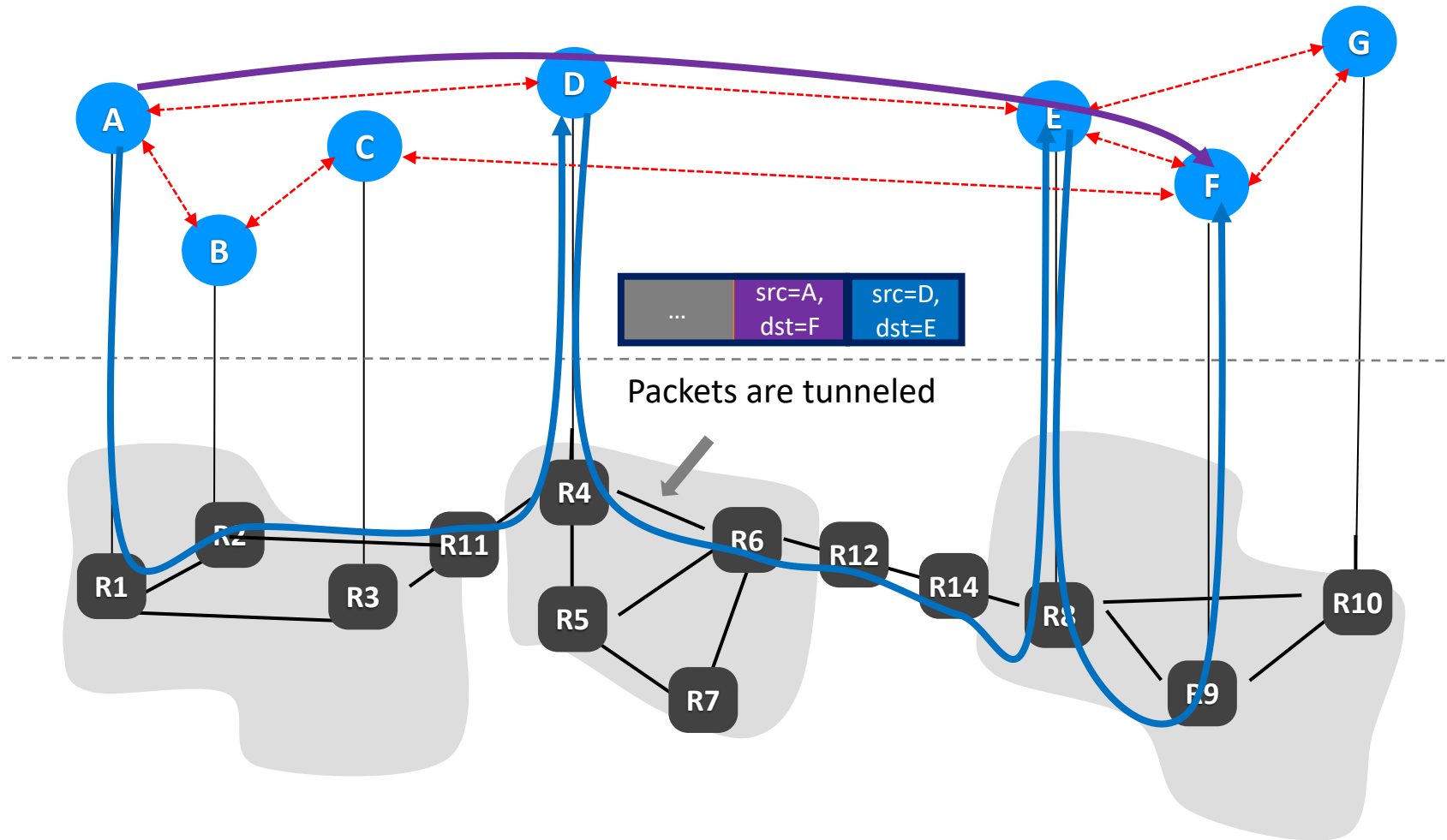
dst	Next-hop
G	D
C	B
F	D
...	...

Underlay routing tables

dst	Next-hop
prefix1	R2
prefix2	R3
prefix3	R2
...	...



# Reviewing overlays

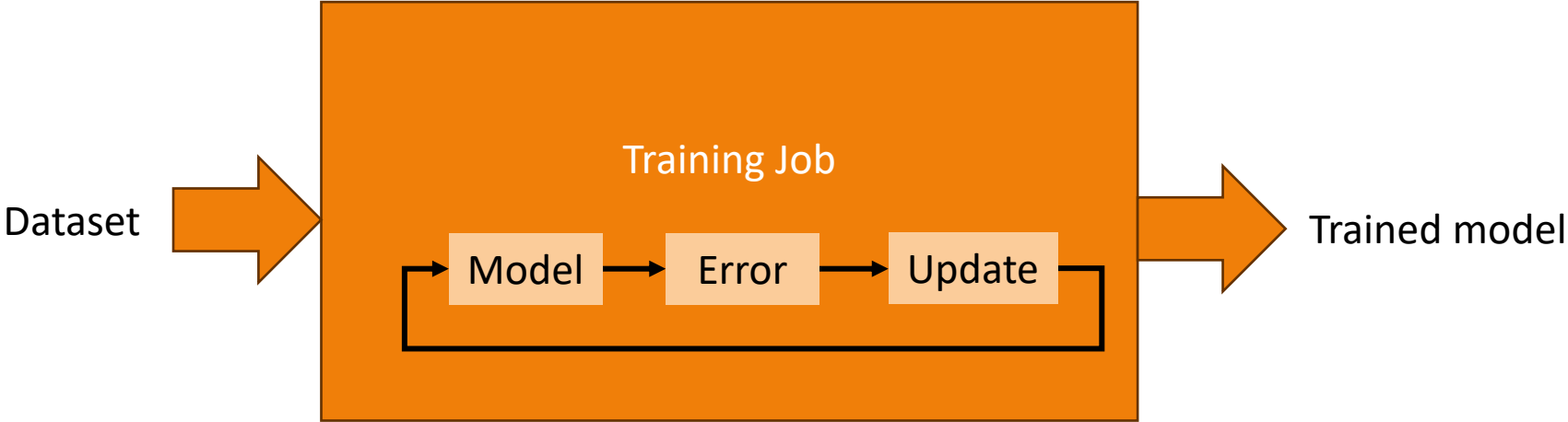




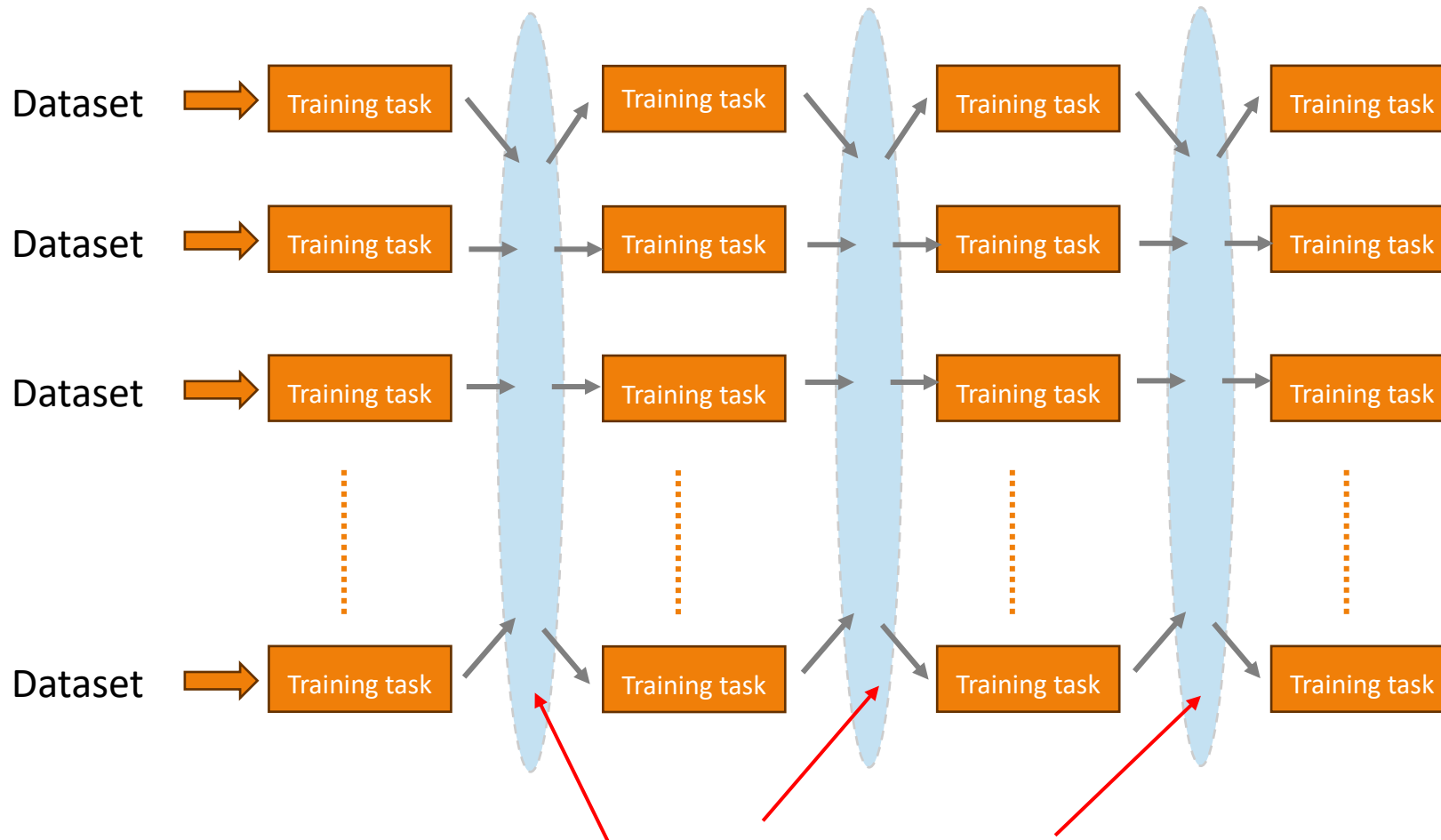
# Outline

- Quick review: Overlays
- Context on distributed training: app and infrastructure
- Communication collectives: definitions and the function they provide
- How collectives are implemented
  - Focus on one collective: AllReduce
  - Implementation at overlay and underlay level

# Training: 10,000 ft. view (for CS168)



# Distributed Training: 10,000 ft. view (for CS168)



Many options for how we paralleliz

Nodes exchange state before proceeding

ita, model, pipeline, hybrid, etc.)

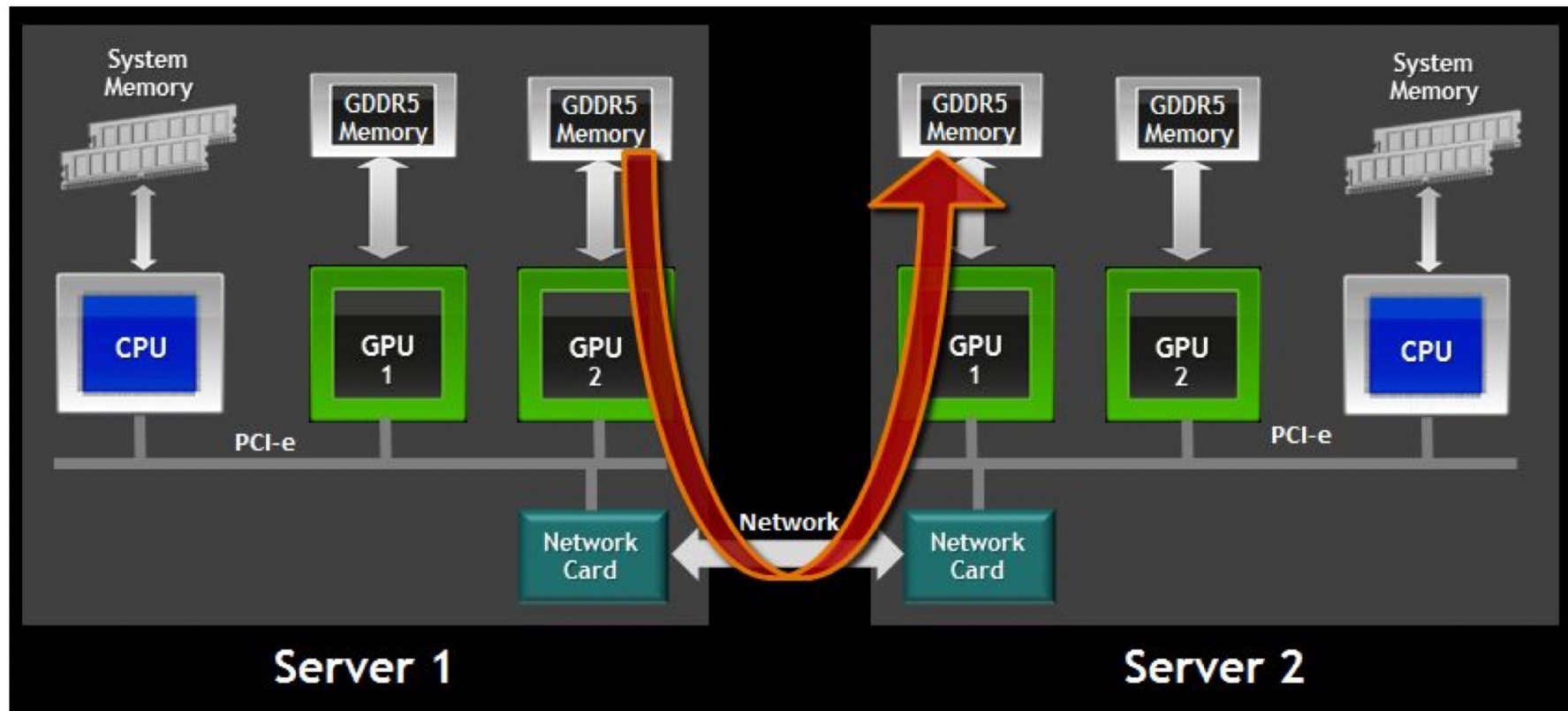
# Collective Communication

- A group of nodes that exchange data in a coordinated manner as part of a group computation
- Nature of this exchange depends on the details of the training and its parallelization
- Distilled into a foundational set of communication patterns referred to as “collectives”
- This lecture: **what** are these collectives and **how are they implemented** in the network
  - We won't get into **why** training/parallelization leads to these particular patterns

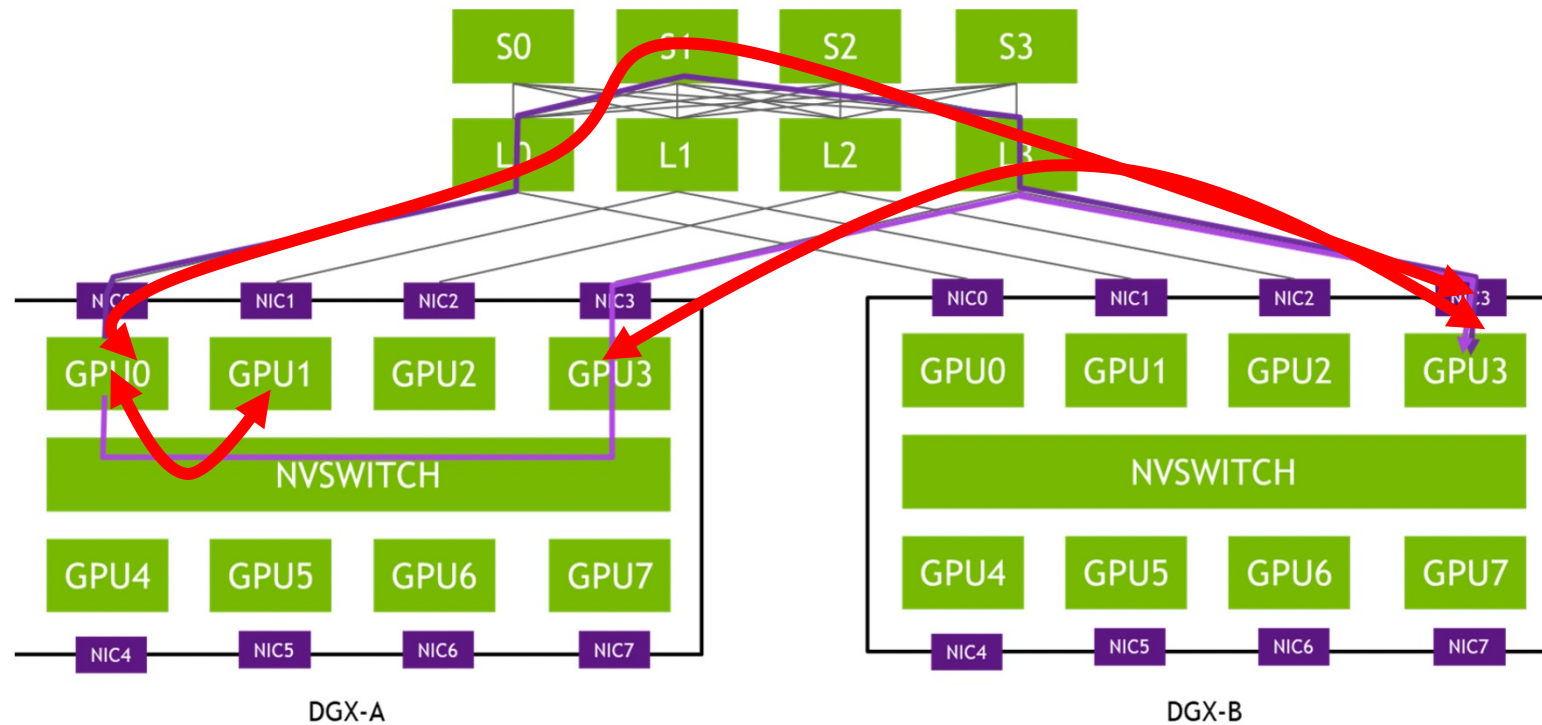
# Distributed Training Infrastructure

- Nodes are specialized Graphics Processing Units (GPUs) or Tensor Processing Units (TPUs)
  - Typical scale: few 100s up to tens of 1000s
- Typically interconnected by a **dedicated datacenter-like** network
  - Dedicated → not shared with other jobs/apps
  - Datacenter-like →
    - Physically close
    - Regular topology (Clos, Torus)
    - Homogeneous
    - High bandwidth

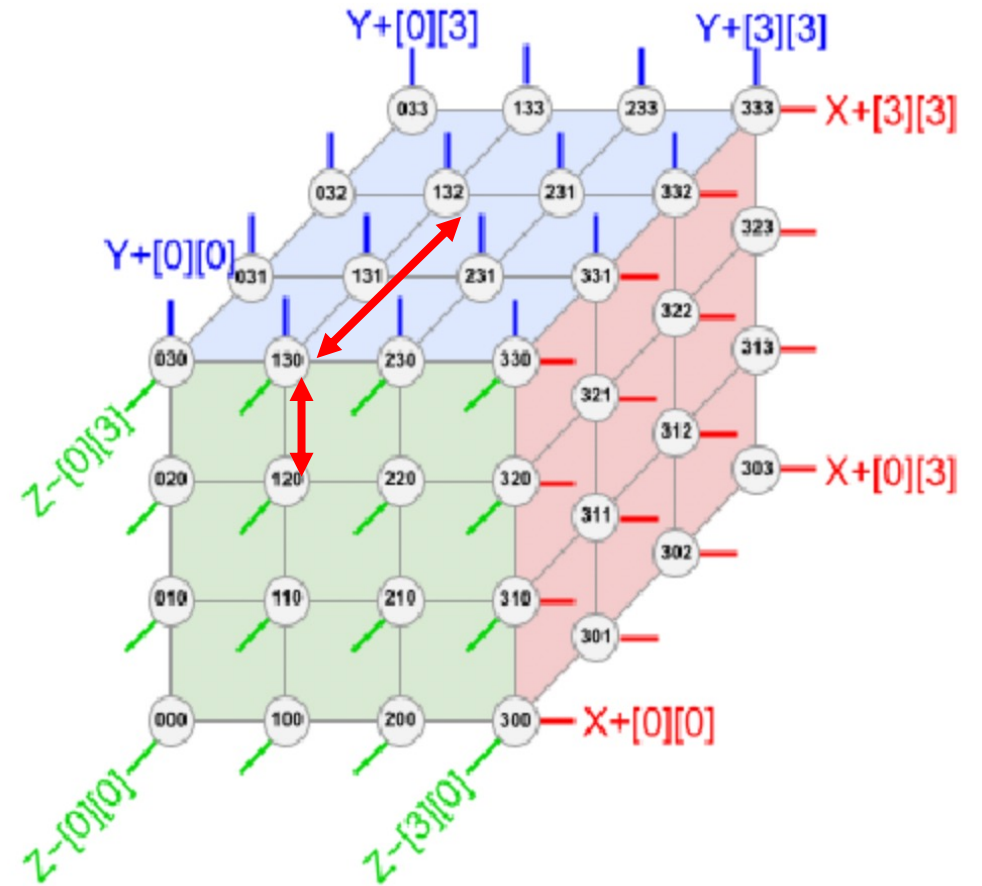
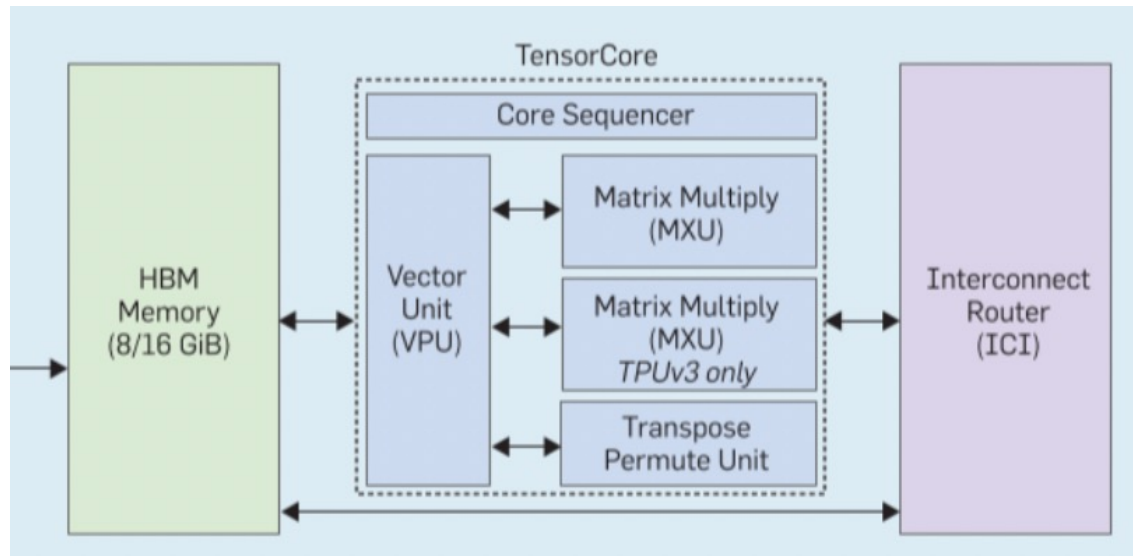
# Distributed Training Infrastructure



# Distributed Training Infrastructure: GPUs



# Distributed Training Infrastructure: TPUs





# Outline

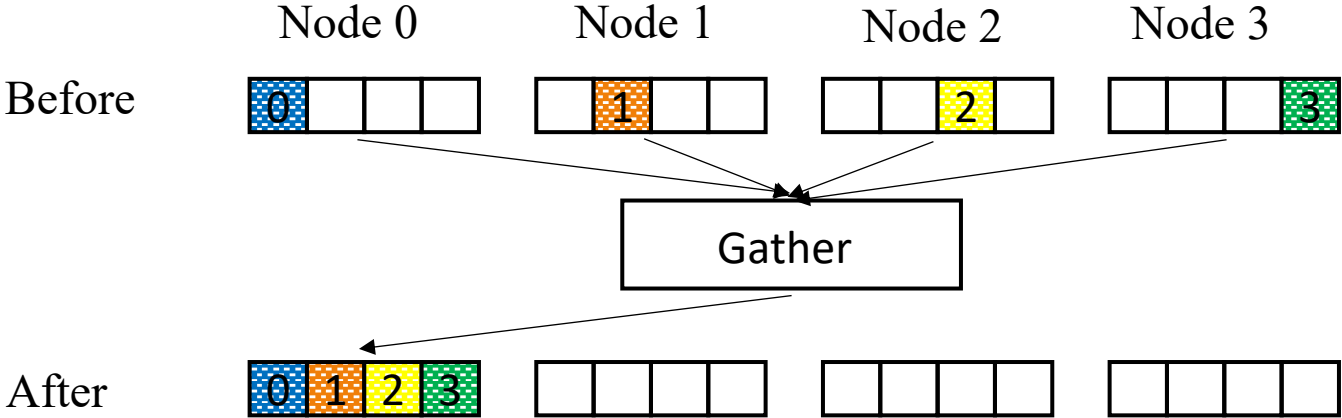
- Quick review: Overlays
- Context on distributed training: app and infrastructure
- Communication collectives: definitions and the function they provide
- How collectives are implemented
  - Focus on one collective: AllReduce
  - Implementation at overlay and underlay level

# Collective Communication

- Coordinated exchange of data between a group of nodes (as part of a group computation)
- Originally designed for supercomputers running high-performance computing (HPC) apps
- Now widely used for AI workloads
  - [Nvidia Collectives Communication Library \(NCCL\)](#)
  - MS Azure MSCCL
  - TCCL, TE-CCL, ...

# A quick example for intuition

(Will return for a deeper look shortly)



# Collectives: common characteristics

- **Synchronized:** nodes (GPUs/TPUs) invoke and execute the collective operation in parallel

## Machine 1

```
comm = communicator.create()  
a = [1, 2, 3]  
b = comm.allreduce(a, op=sum)
```

---

```
assert b == [2, 2, 4]
```

## Machine 2

```
comm = communicator.create()  
a = [1, 0, 1]  
b = comm.allreduce(a, op=sum)
```

---

```
assert b == [2, 2, 4]
```

# Collectives: common characteristics

- **Synchronized**: nodes (GPUs/TPUs) invoke and execute the collective operation in parallel
- **Orchestrated**: centralized job scheduler distributes membership info, member IDs, roles, *etc.*
- **Homogeneous**: all nodes have same resources (compute, BW) and running same code
- **Blocking** (implicit or explicit): collective must complete at all nodes before proceeding
- Data may be **transformed** at intermediate hops in the communication path

# Collectives: taxonomy and key operations

- Data redistribution (about moving data)
  - Broadcast
  - Scatter
  - Gather
  - AllGather
- Data consolidation (about moving *and* aggregating or “reducing” the data)\*
  - Reduce
  - Reduce-Scatter
  - AllReduce

\*We'll assume that aggregation = sum

# Notation

- Number of nodes:  $p$  (=4 in most of our examples)
- Vector of data being exchanged:  $x$
- For some operations,  $x$  is subdivided into subvectors  $x_i$ ,  $i=0, 1, \dots, p-1$
- Superscript denotes a vector that must be summed with vectors from other nodes:

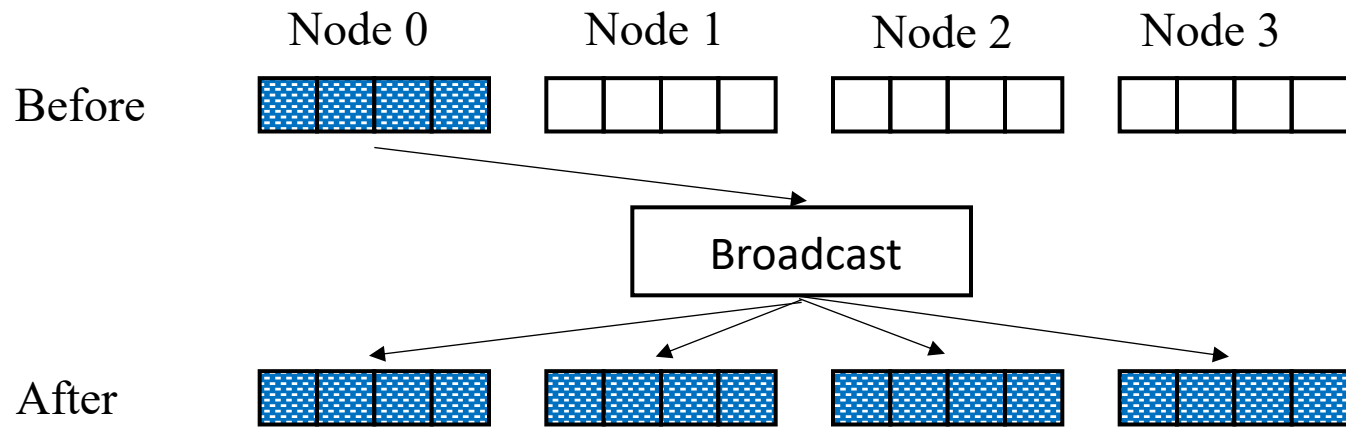
$\sum_j x_i^{(j)}$  indicates the sum of the subvector  $x_i$  from all nodes  $j$

Let's now look at our set of collectives ...



# Broadcast (p=4)

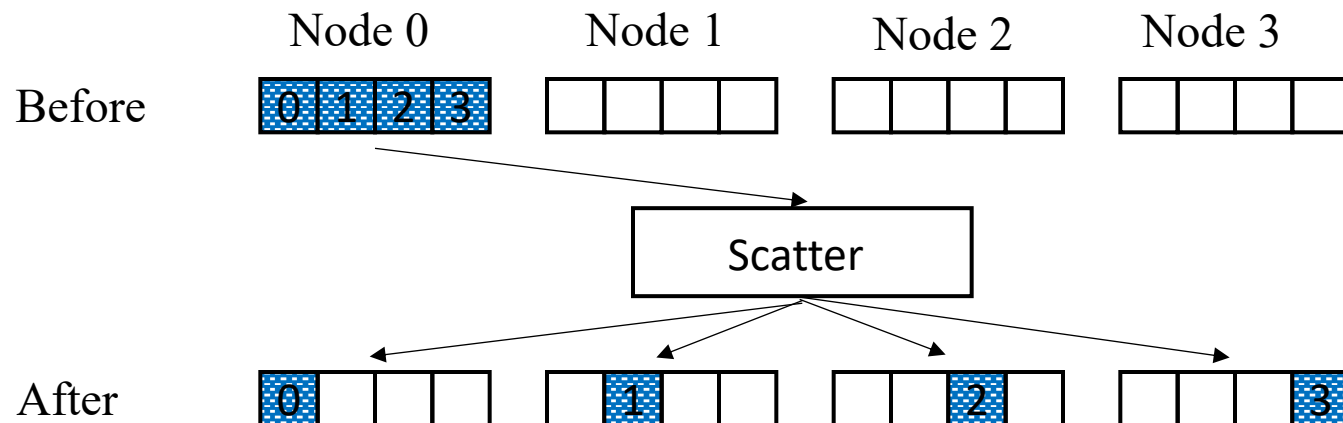
Operation	Before	After																
Broadcast	<table border="1"><tr><td>Node 0</td><td>Node 1</td><td>Node 2</td><td>Node 3</td></tr><tr><td><math>x</math></td><td></td><td></td><td></td></tr></table>	Node 0	Node 1	Node 2	Node 3	$x$				<table border="1"><tr><td>Node 0</td><td>Node 1</td><td>Node 2</td><td>Node 3</td></tr><tr><td><math>x</math></td><td><math>x</math></td><td><math>x</math></td><td><math>x</math></td></tr></table>	Node 0	Node 1	Node 2	Node 3	$x$	$x$	$x$	$x$
Node 0	Node 1	Node 2	Node 3															
$x$																		
Node 0	Node 1	Node 2	Node 3															
$x$	$x$	$x$	$x$															



*“send to all”*

# Scatter (p=4)

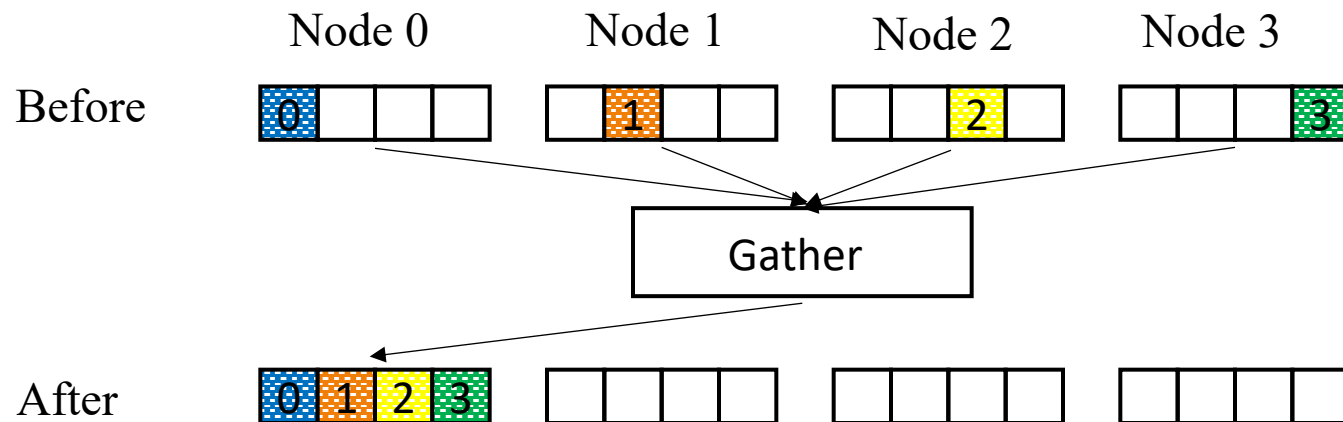
Operation	Before				After			
Scatter	Node 0	Node 1	Node 2	Node 3	Node 0	Node 1	Node 2	Node 3
	$x_0$				$x_0$			
	$x_1$					$x_1$		
	$x_2$						$x_2$	
	$x_3$							$x_3$



*send my  $i^{\text{th}}$  subvector  
to node  $i$*

# Gather (p=4)

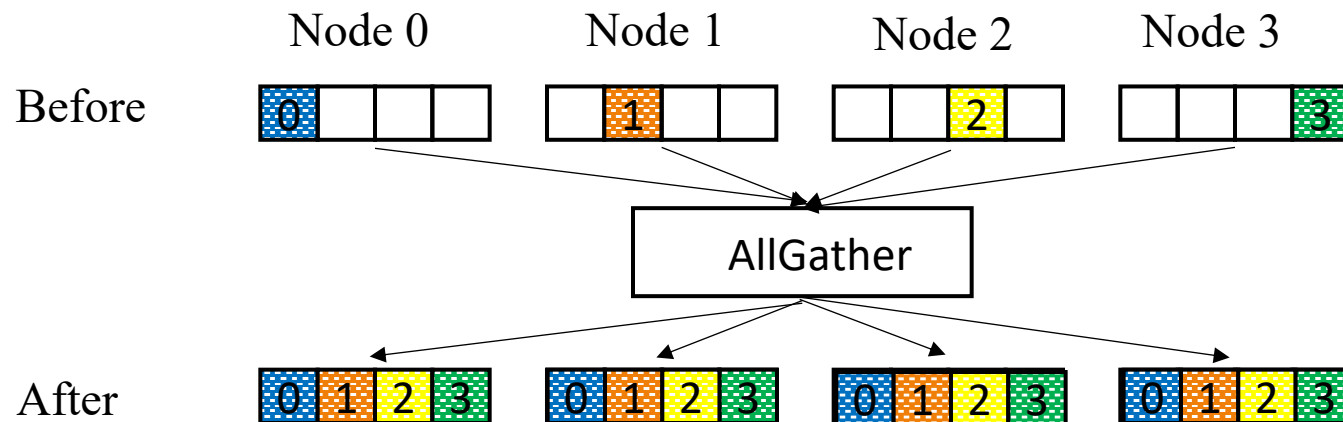
Operation	Before				After			
Gather	Node 0	Node 1	Node 2	Node 3	Node 0	Node 1	Node 2	Node 3
	$x_0$	$x_1$	$x_2$	$x_3$	$x_0$			



*node  $i$  sends its  $i^{\text{th}}$  subvector to a target node*

# AllGather (p=4)

Operation	Before				After				
Allgather	Node 0	Node 1	Node 2	Node 3	Node 0	Node 1	Node 2	Node 3	
	$x_0$				$x_0$	$x_0$	$x_0$	$x_0$	$x_0$
		$x_1$			$x_1$	$x_1$	$x_1$	$x_1$	$x_1$
			$x_2$		$x_2$	$x_2$	$x_2$	$x_2$	$x_2$
				$x_3$	$x_3$	$x_3$	$x_3$	$x_3$	$x_3$



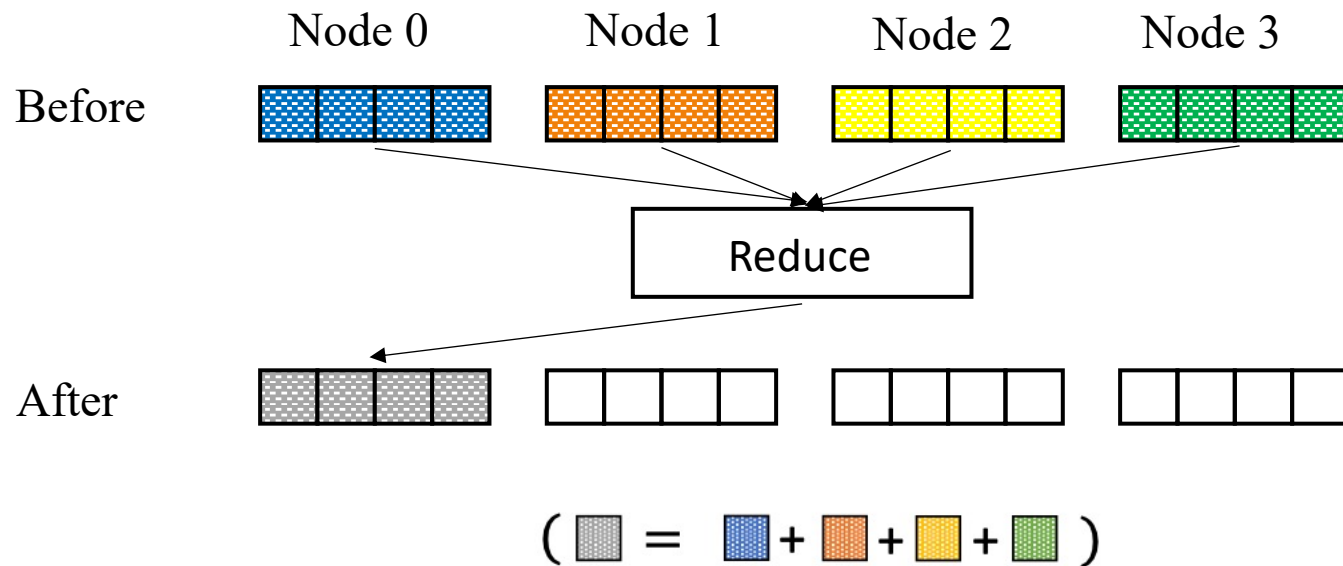
*Node i sends its  
i<sup>th</sup> subvector to all*

# Collectives: taxonomy and key operations

- Data redistribution (about moving data)
  - Broadcast
  - Scatter
  - Gather
  - AllGather
- Data consolidation (about moving *and* aggregating the data)\*
  - Reduce
  - Reduce-Scatter
  - AllReduce

# Reduce (p=4)

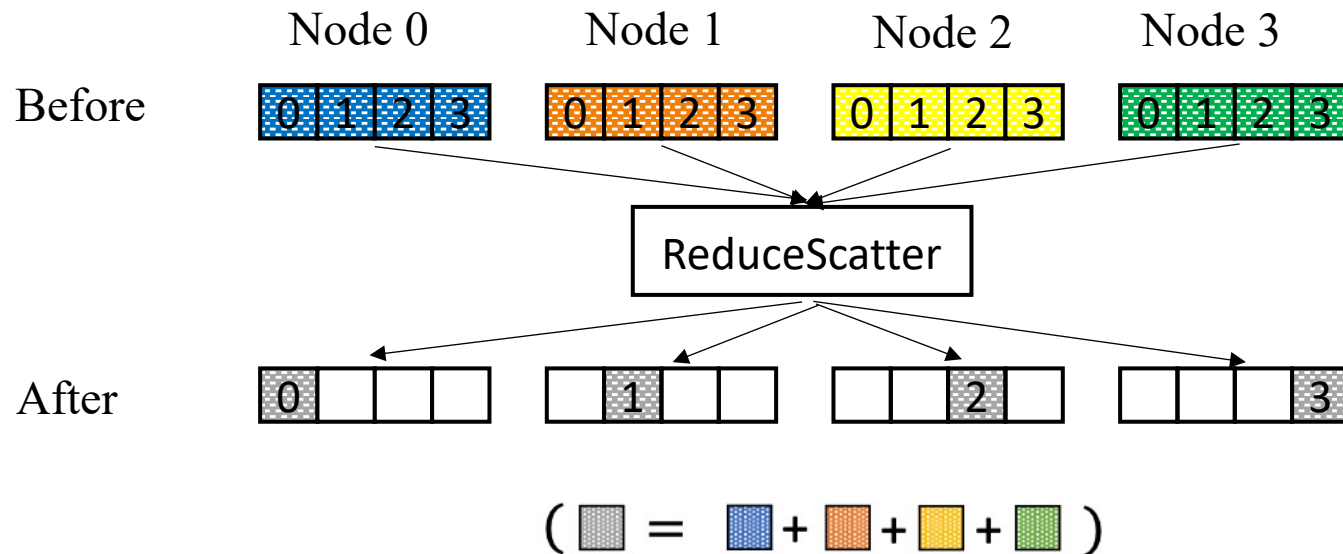
Operation	Before				After			
Reduce(-to-one)	Node 0 $x^{(0)}$	Node 1 $x^{(1)}$	Node 2 $x^{(2)}$	Node 3 $x^{(3)}$	Node 0 $\sum_j x^{(j)}$	Node 1	Node 2	Node 3



*All nodes' vectors are summed at a target node*

# ReduceScatter (p=4)

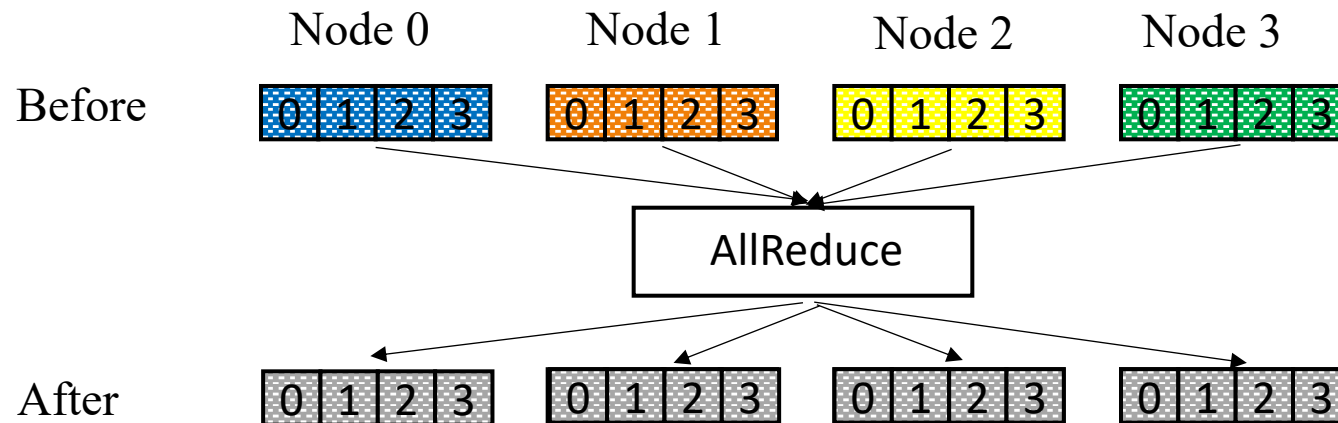
Operation	Before				After			
Reduce-scatter	Node 0	Node 1	Node 2	Node 3	Node 0	Node 1	Node 2	Node 3
	$x_0^{(0)}$	$x_0^{(1)}$	$x_0^{(2)}$	$x_0^{(3)}$	$\sum_j x_0^{(j)}$			
	$x_1^{(0)}$	$x_1^{(1)}$	$x_1^{(2)}$	$x_1^{(3)}$		$\sum_j x_1^{(j)}$		
	$x_2^{(0)}$	$x_2^{(1)}$	$x_2^{(2)}$	$x_2^{(3)}$			$\sum_j x_2^{(j)}$	
	$x_3^{(0)}$	$x_3^{(1)}$	$x_3^{(2)}$	$x_3^{(3)}$				$\sum_j x_3^{(j)}$



*The  $i^{\text{th}}$  subvector from all nodes is summed at node  $i$*

# AllReduce (p=4)

Operation	Before				After			
Allreduce	Node 0 $x^{(0)}$	Node 1 $x^{(1)}$	Node 2 $x^{(2)}$	Node 3 $x^{(3)}$	Node 0 $\sum_j x^{(j)}$	Node 1 $\sum_j x^{(j)}$	Node 2 $\sum_j x^{(j)}$	Node 3 $\sum_j x^{(j)}$



*All nodes have the sum of all nodes' vectors*

$$(\text{grey box}) = (\text{blue box}) + (\text{orange box}) + (\text{yellow box}) + (\text{green box})$$



# Some collectives are duals of each other

- Collectives C1 and C2 are duals if *reversing* the communication in C1 yields C2
  - Reversing means  $A \rightarrow B$  becomes  $B \rightarrow A$
  - We ignore any computation when determining duality
- For a given topology/routing scheme, a collective and its dual see equivalent performance

# Broadcast and Reduce are duals

Operation	Before				After			
Broadcast	Node 0 $x$	Node 1	Node 2	Node 3	Node 0 $x$	Node 1 $x$	Node 2 $x$	Node 3 $x$
Reduce(- to-one)	Node 0 $x^{(0)}$	Node 1 $x^{(1)}$	Node 2 $x^{(2)}$	Node 3 $x^{(3)}$	Node 0 $\sum_j x^{(j)}$	Node 1	Node 2	Node 3

# Scatter and Gather are duals

Operation	Before				After			
Scatter	Node 0	Node 1	Node 2	Node 3	Node 0	Node 1	Node 2	Node 3
	$x_0$				$x_0$	$x_1$	$x_2$	$x_3$
Gather	Node 0	Node 1	Node 2	Node 3	Node 0	Node 1	Node 2	Node 3
	$x_0$	$x_1$	$x_2$	$x_3$	$x_0$			

# AllGather and ReduceScatter are duals

Operation	Before				After			
Allgather	Node 0	Node 1	Node 2	Node 3	Node 0	Node 1	Node 2	Node 3
	$x_0$	$x_1$	$x_2$	$x_3$	$x_0$	$x_0$	$x_0$	$x_0$
Reduce-scatter	Node 0	Node 1	Node 2	Node 3	Node 0	Node 1	Node 2	Node 3
	$x_0^{(0)}$	$x_0^{(1)}$	$x_0^{(2)}$	$x_0^{(3)}$	$\sum_j x_0^{(j)}$	$\sum_j x_1^{(j)}$	$\sum_j x_2^{(j)}$	$\sum_j x_3^{(j)}$
Reduce-scatter	$x_1^{(0)}$	$x_1^{(1)}$	$x_1^{(2)}$	$x_1^{(3)}$	$\sum_j x_1^{(j)}$	$\sum_j x_2^{(j)}$	$\sum_j x_3^{(j)}$	$\sum_j x_3^{(j)}$
	$x_2^{(0)}$	$x_2^{(1)}$	$x_2^{(2)}$	$x_2^{(3)}$				
Reduce-scatter	$x_3^{(0)}$	$x_3^{(1)}$	$x_3^{(2)}$	$x_3^{(3)}$	$\sum_j x_3^{(j)}$	$\sum_j x_3^{(j)}$	$\sum_j x_3^{(j)}$	$\sum_j x_3^{(j)}$
	$x_3^{(0)}$	$x_3^{(1)}$	$x_3^{(2)}$	$x_3^{(3)}$				

# AllReduce does not have a dual

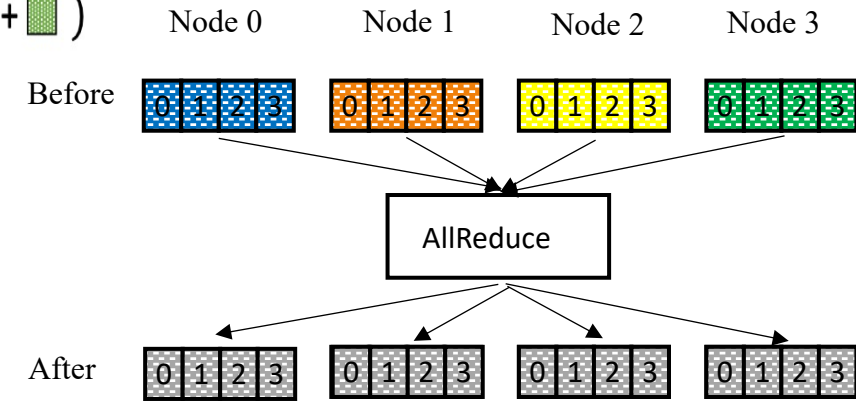
**Note:** AllReduce = ReduceScatter + AllGather (Convince yourself of this later!)

# Outline

- Context on distributed training: app and infrastructure
- Communication collectives: definitions and the function they provide
- How collectives are implemented
  - We'll focus on one collective: AllReduce

# Implementing AllReduce

$$(\text{Grey} = \text{Blue} + \text{Orange} + \text{Yellow} + \text{Green})$$



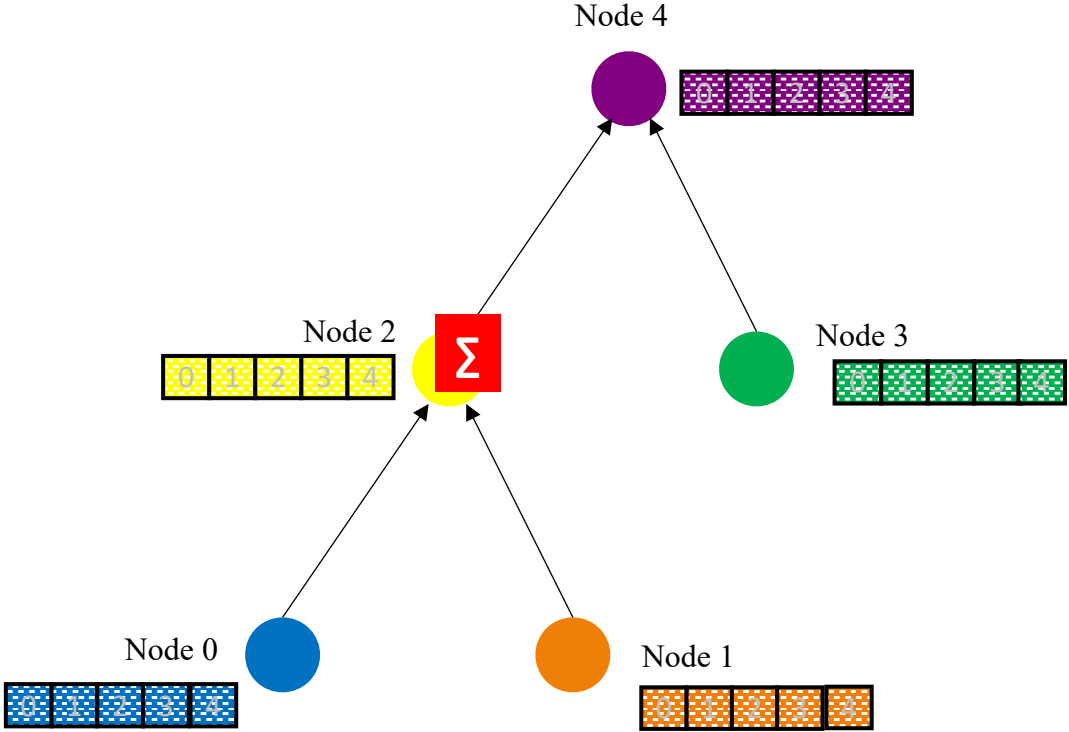
- Let's consider the simplest implementation: a full mesh of node-to-node exchanges
  - Every node sends its vector  $x$  directly to every other node
  - Each node sums all the vectors it receives
- Assuming  $p$  nodes and that vector  $x$  is of size  $D$  bytes:
  - Each node implements  $p-1$  transmit, receive, and summation operations
  - Each on a vector of size  $D$  bytes
  - Total traffic in the network:  $O(p^2 \times D)$
  - Consider  $D \sim 100$ s of GB and  $p \sim 1000$ s of nodes
- Need a more scalable solution!

# Implementing AllReduce ... using overlays

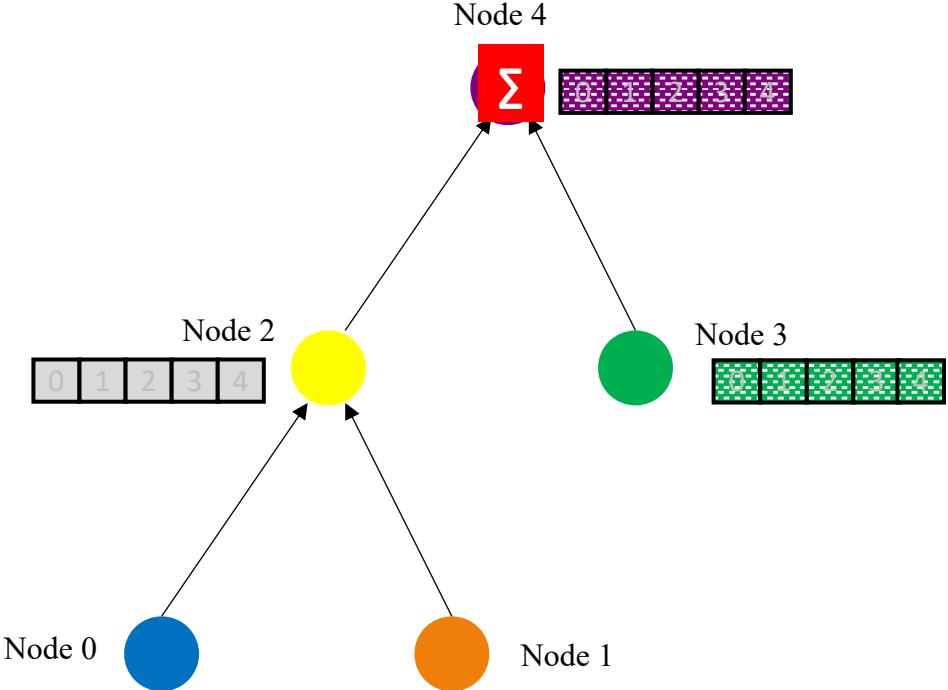
- Idea: construct a virtual topology between the  $p$  nodes; “reduce” (sum) data as it traverses the overlay
- Details vary depending on the virtual topology selected
- Two typical choices: tree or ring



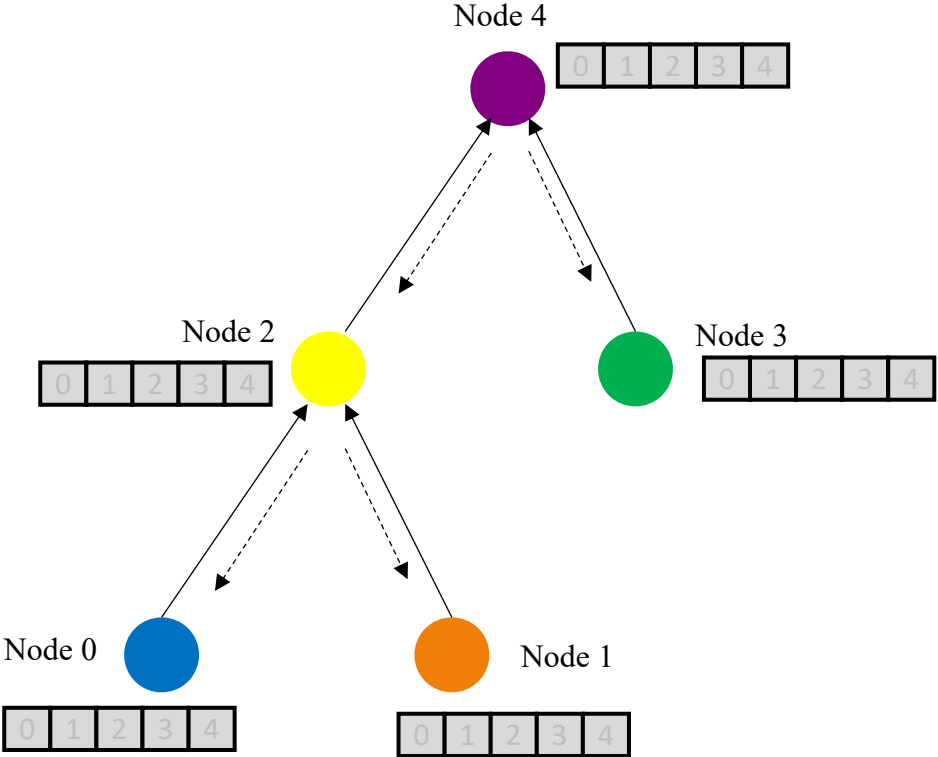
# Tree-based AllReduce



# Tree-based AllReduce



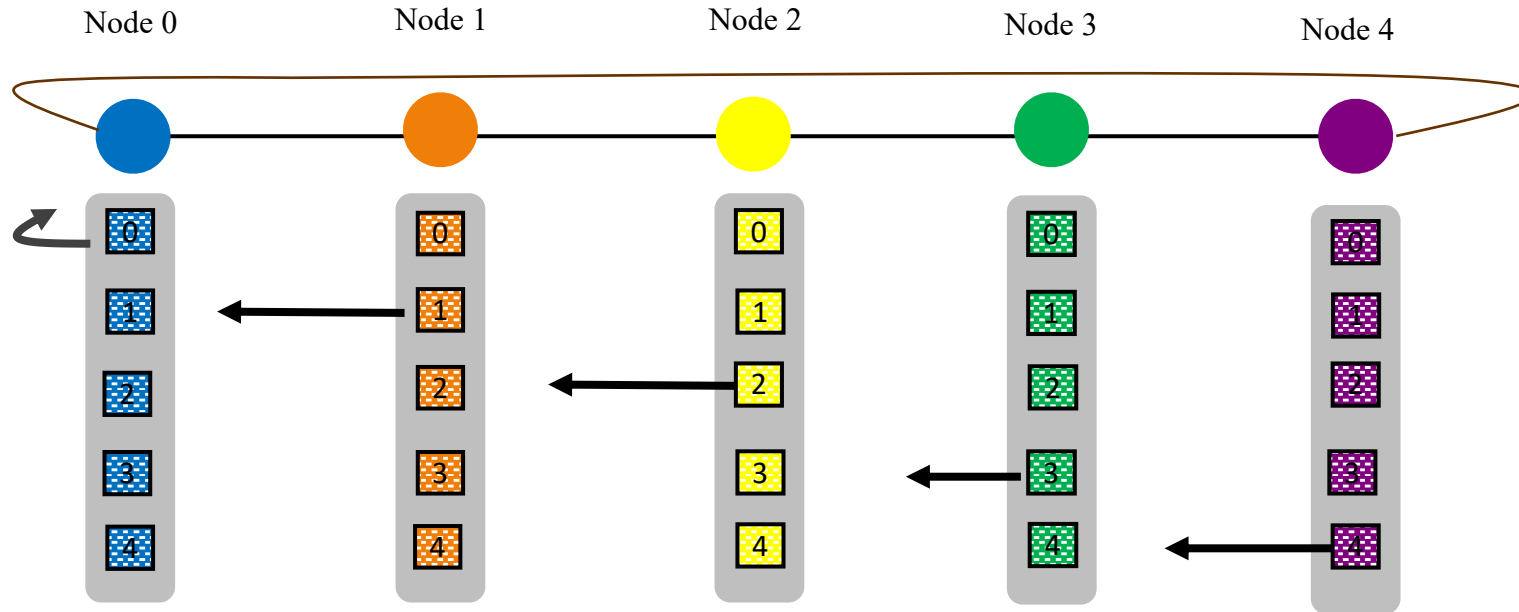
# Tree-based AllReduce



# Tree-based AllReduce

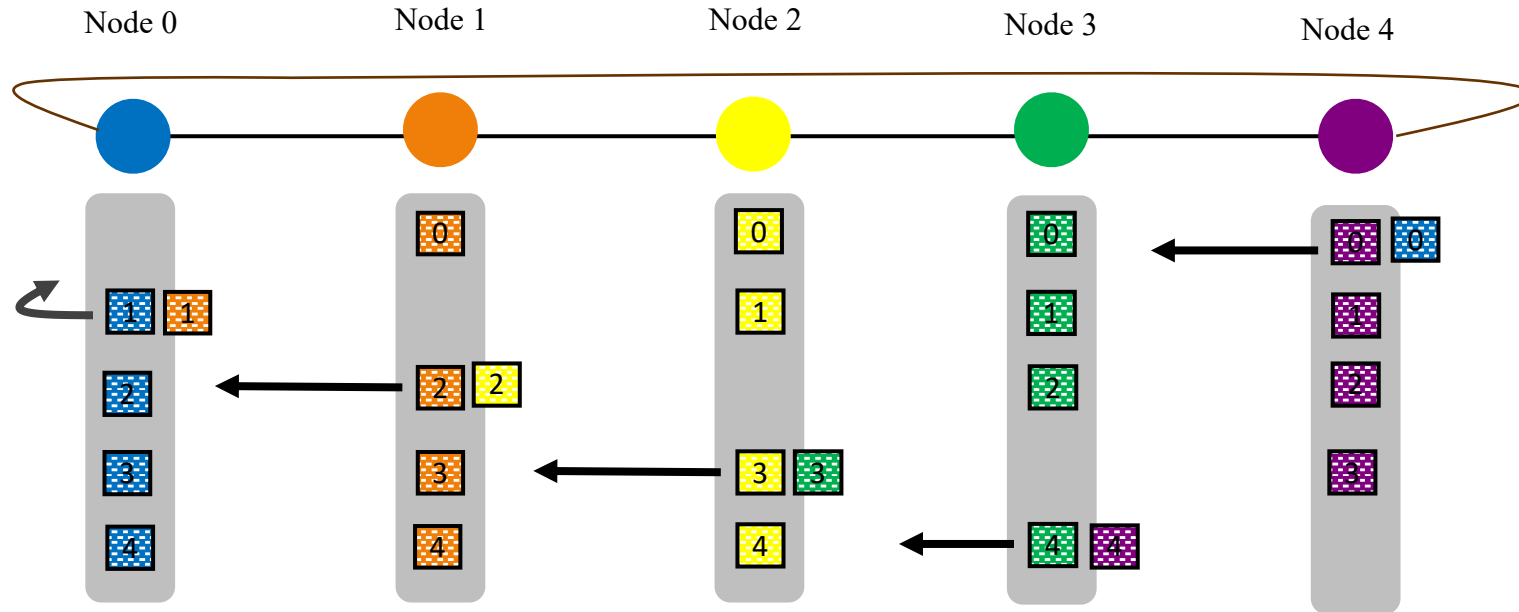
- Nodes form a logical tree; aggregate to the root and then broadcast
  - Leaf node: transmits vector to its parent
  - Every intermediate node
    - aggregates (sums) its own vector with that from each child node
    - send the resulting (aggregated) vector to its parent
  - Root broadcasts the final aggregated vector down the tree
- Assuming a binary tree with  $p$  nodes and vector of  $D$  bytes:
  - AllReduce involves  $O(\log P)$  steps – *to travel up and down the tree*
  - Each node implements 3 transmit, receive, and summation operations, each on a vector of size  $D$  bytes
  - Total traffic in the network:  **$O(p \times D)$  across all steps** – *much better than with the mesh!*

# Ring-based AllReduce



- Initial step: node  $i$  sends  $i^{\text{th}}$  subvector to its predecessor on the ring

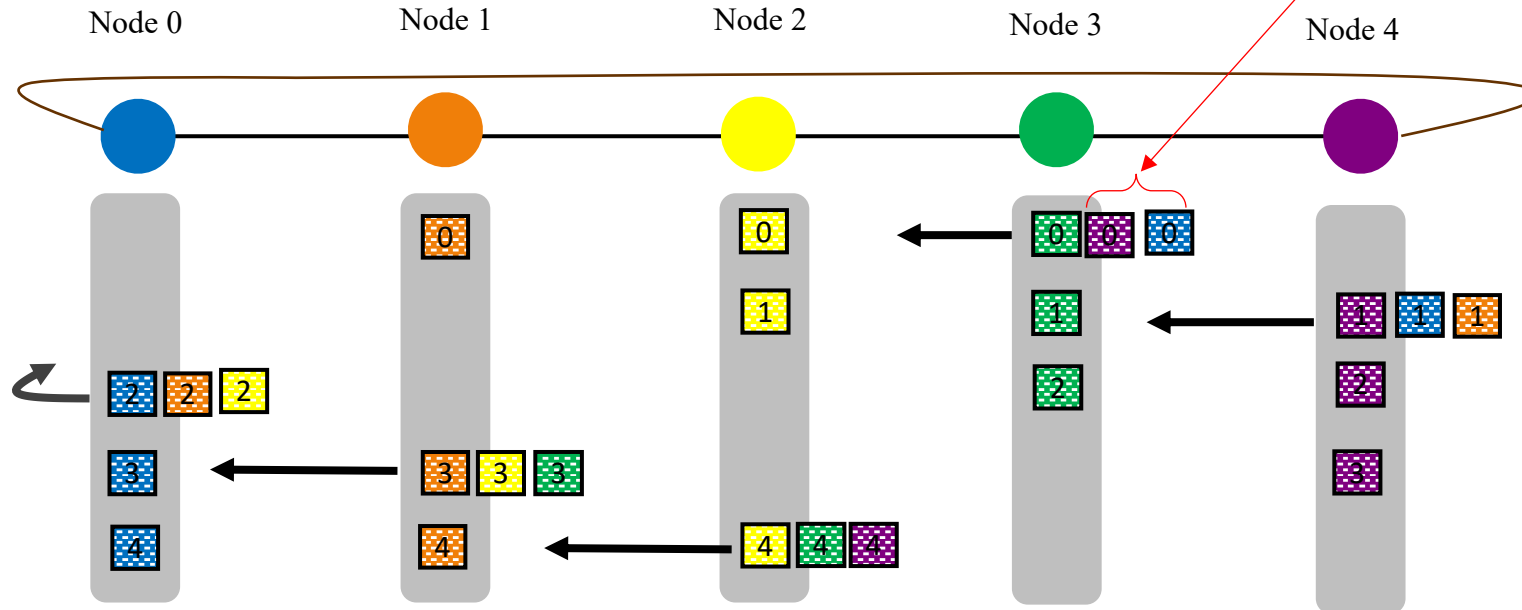
# Ring-based AllReduce



- Initial step: node  $i$  sends  $i^{\text{th}}$  subvector to its predecessor on the ring
- When a node receives the  $k^{\text{th}}$  subvector from its successor, it adds this to its own  $k^{\text{th}}$  subvector and in the next step, sends this (aggregated)  $k^{\text{th}}$  subvector to its predecessor

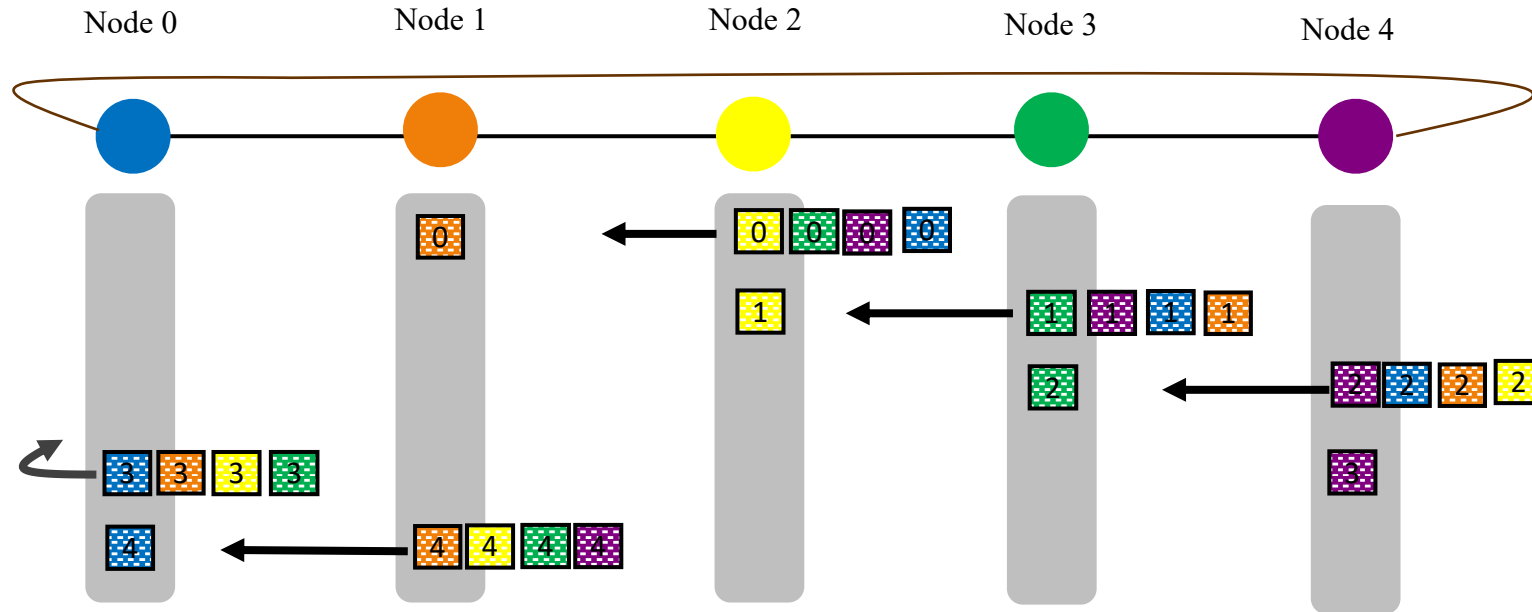
# Ring-based AllReduce

In reality: we send a single subvector equal to the sum of these subvectors (not showing the aggregated subvector for illustration purposes only)



- Initial step: node  $i$  sends  $i^{\text{th}}$  subvector to its predecessor on the ring
- When a node receives the  $k^{\text{th}}$  subvector from its successor, it adds this to its own  $k^{\text{th}}$  subvector and in the next step, sends this (aggregated)  $k^{\text{th}}$  subvector to its predecessor

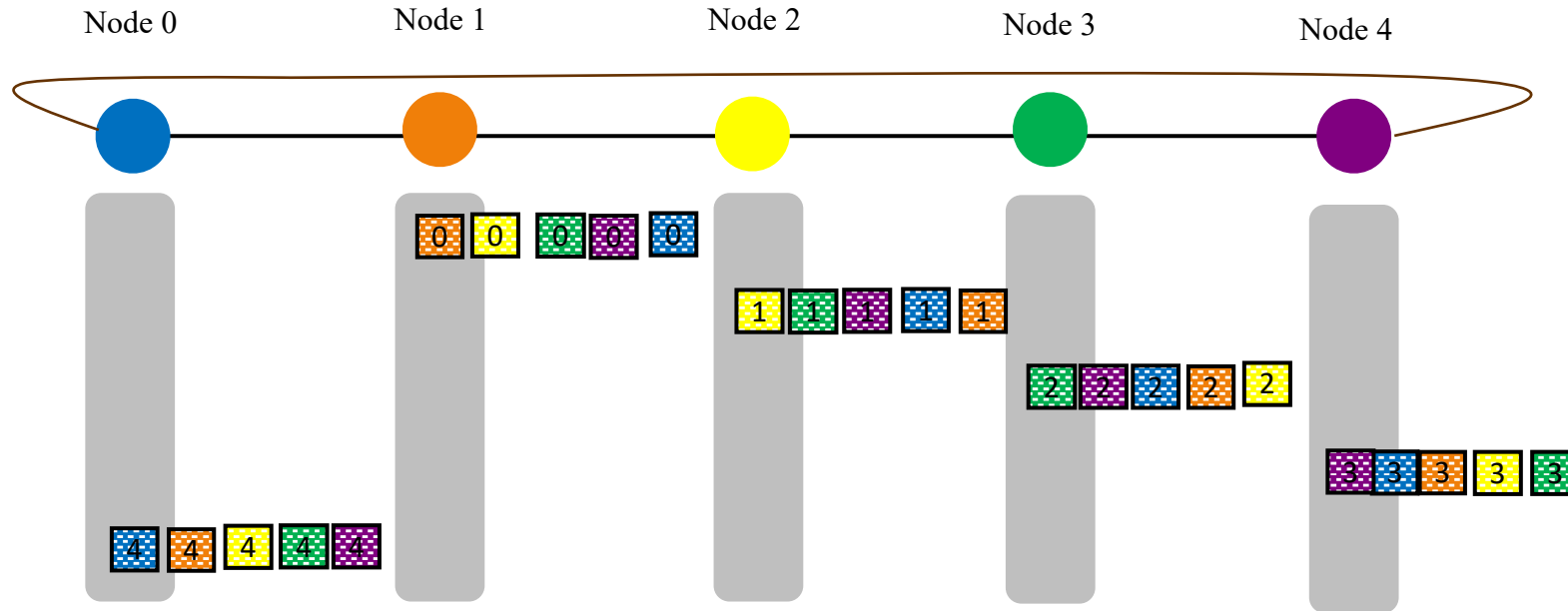
# Ring-based AllReduce



- Initial step: node  $i$  sends  $i^{\text{th}}$  subvector to its predecessor on the ring
- When a node receives the  $k^{\text{th}}$  subvector from its successor, it adds this to its own  $k^{\text{th}}$  subvector and in the next step, sends this (aggregated)  $k^{\text{th}}$  subvector to its predecessor



# Ring-based AllReduce

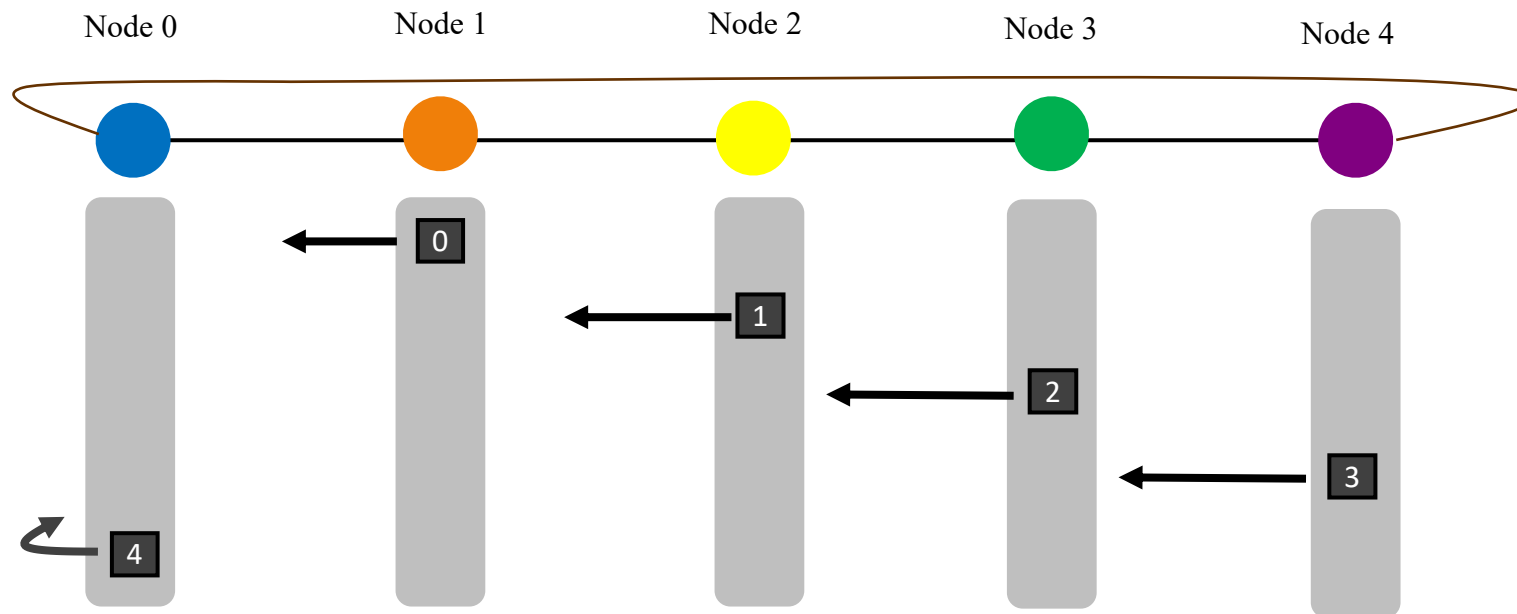


- Initial step: node  $i$  sends  $i^{\text{th}}$  subvector to its predecessor on the ring
- When a node receives the  $k^{\text{th}}$  subvector from its successor, it adds this to its own  $k^{\text{th}}$  subvector and in the next step, sends this (aggregated)  $k^{\text{th}}$  subvector to its predecessor
- After  $p-1$  steps, the successor of node  $i$  will have the fully aggregated  $i^{\text{th}}$  subvector



Denotes sum of  $i^{\text{th}}$  subvector from all nodes

# Ring-based AllReduce



- Initial step: node  $i$  sends  $i^{\text{th}}$  subvector to its predecessor on the ring
- When a node receives the  $k^{\text{th}}$  subvector from its successor, it adds this to its own  $k^{\text{th}}$  subvector and in the next step, sends this (aggregated)  $k^{\text{th}}$  subvector to its predecessor
- After  $p-1$  steps, the successor of node  $i$  will have the fully aggregated  $i^{\text{th}}$  subvector
- Repeat (without aggregation); after  $p-1$  steps all nodes have the entire reduced vector!

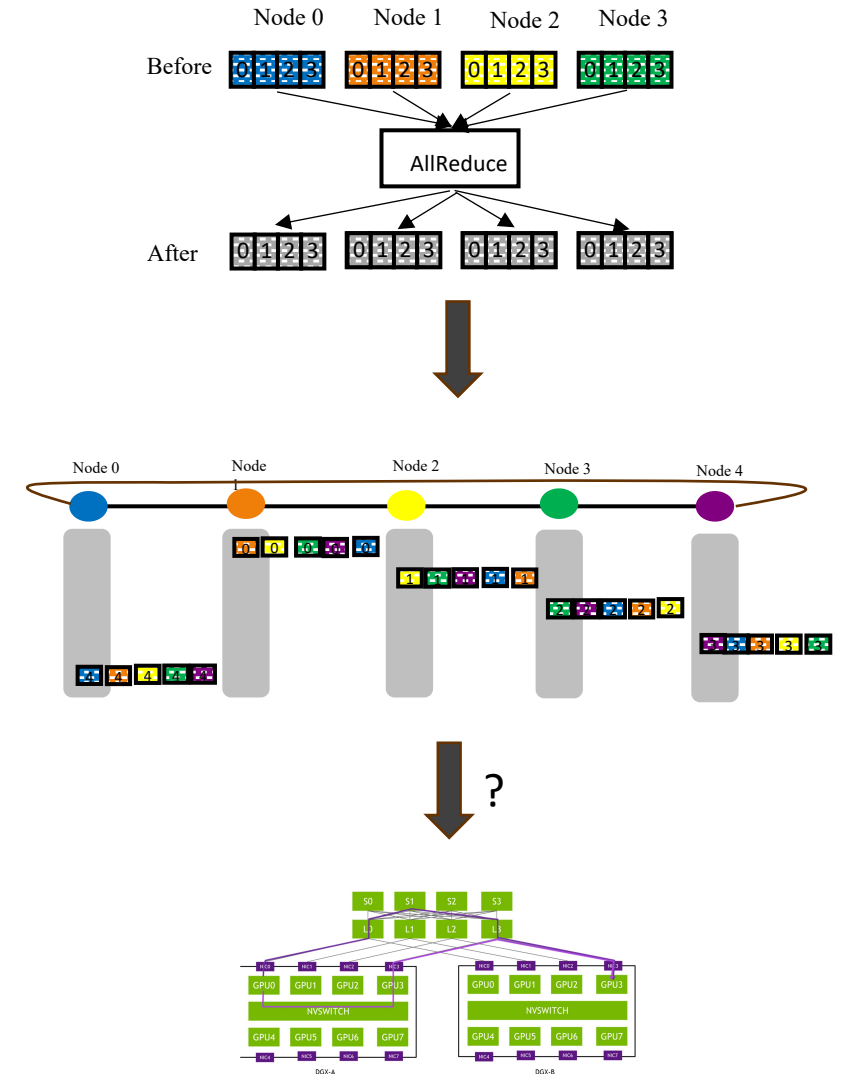
# Ring-based AllReduce

- Process
  - Initial step: node  $i$  sends  $i^{\text{th}}$  chunk to its predecessor on the logical ring
  - When a node receives the  $k^{\text{th}}$  chunk from its successor, it adds this to its own  $k^{\text{th}}$  chunk and in the next step, sends this (aggregated)  $k^{\text{th}}$  chunk to its predecessor
  - After  $p-1$  steps, the successor of node  $i$  will have the fully aggregated  $i^{\text{th}}$  chunk
  - Repeat but now without aggregation; after  $p-1$  steps all nodes have the entire reduced vector!
- Assuming a ring with  $p$  nodes and vector of  $D$  bytes:
  - Takes  $\sim 2p$  steps
  - In each step, a node transmits/receives/sums a subvector of size  $D/p$  bytes
  - Total traffic in the network across all steps:  $O(p \times D)$

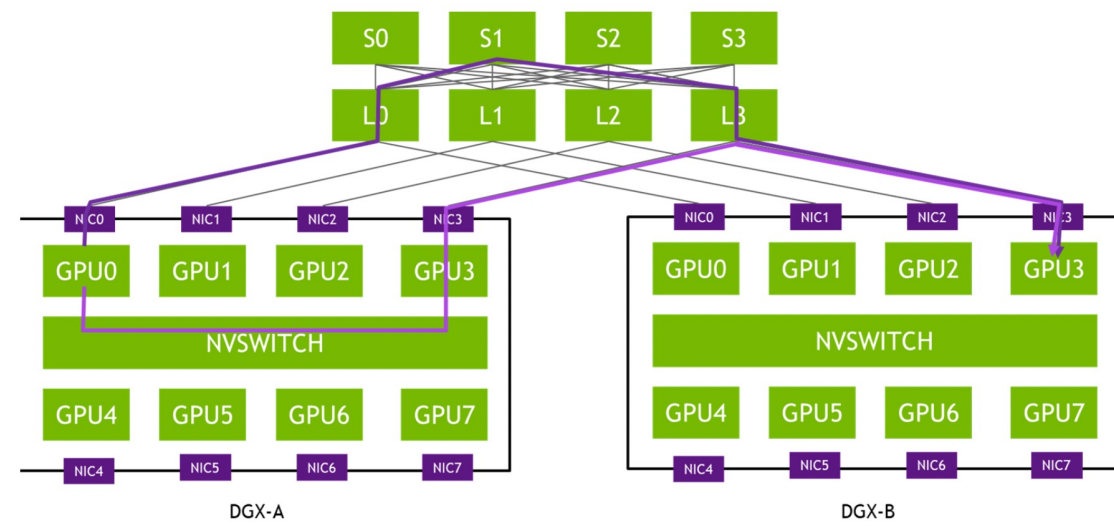
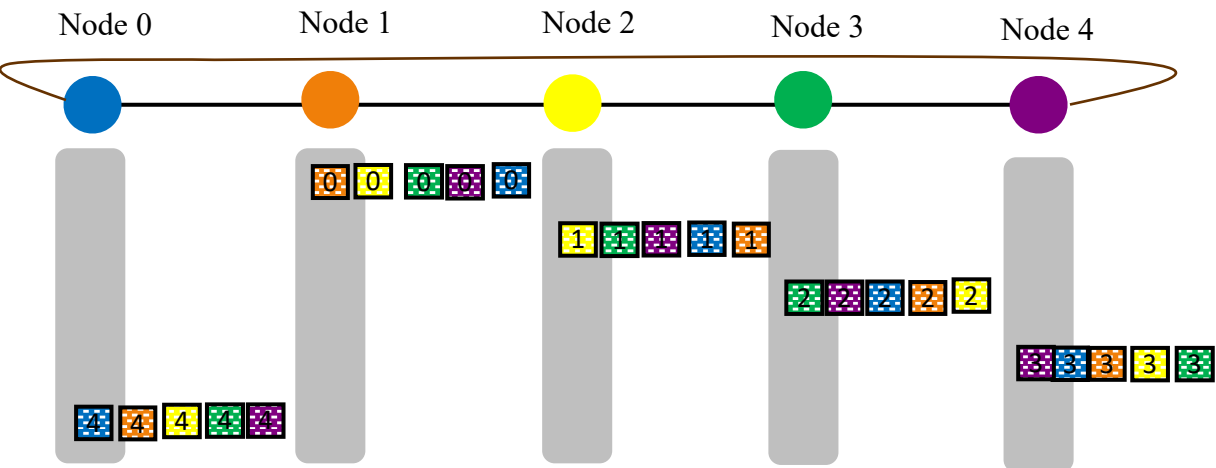
# Taking Stock

Three levels at which we can view communication collectives:

- Definition or “service model”
- Overlay viewpoint
- Underlay (physical network) viewpoint



# Closing point:



How do we assign overlay nodes to physical GPUs so as to achieve low “stretch” (last lecture)?

**Questions?**