# ASTRA-sim and Chakra Tutorial:
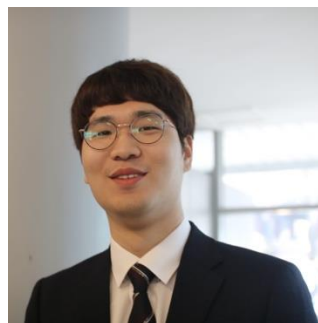## *System Layer*

Will Won
Ph.D. Student
School CS, Georgia Institute of Technology
william.won@gatech.edu

# ASTRA-sim Tutorial - Agenda

| Time (PDT) | Topic | Presenter |
|---|---|---|
| 3:00 – 3:30 pm | **Introduction to Distributed ML** | Tushar Krishna |
| 3:30 – 3:45 pm | **Overview of Chakra and ASTRA-sim** | Tushar Krishna |
| 3:45 – 4:35 pm | **Deeper Dive into Chakra and ASTRA-sim** | Will Won |
| | Workload, System, and Network Layers | |
| 4:35 – 4:45 pm | **Demo** | Will Won |
| 4:45 – 5:00 pm | **Closing Remarks** | Tushar Krishna |

**Tutorial Website**
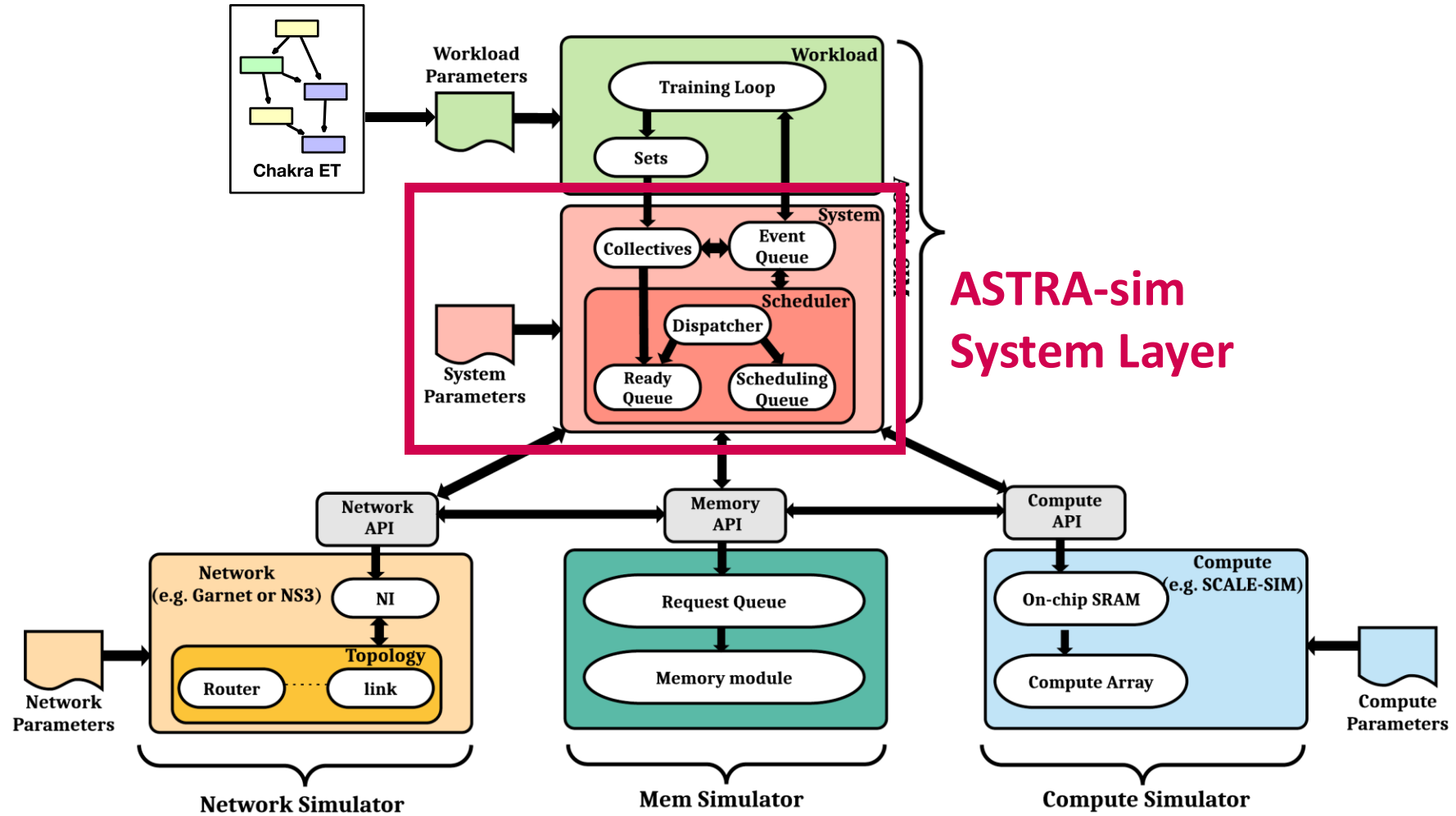 *includes agenda, slides, ASTRA-sim installation instructions (via source + docker image)*
***https://astra-sim.github.io/tutorials/hoti-2024***

**Attention:** Tutorial is being recorded

# Design Space: System

**Workload**

• DNN Model
• Parallelization Strategy
• Framework Policy and Pattern

**System**

• **Compute-Comms Mechanism**
• **Collective Scheduling**
• **Collective Algorithm**

**Compute**

• **Supported Compute Kernels**
• **Mapping and Dataflow**
• **Compute Architecture**

**Network**

• **Network Topology**
• **Connectivity Design**
• **Network Implementation**

**Remote Memory**

• **Memory System Architecture**
• **Memory Topology**
• **Memory Implementation**

# ASTRA-sim: System Layer



**ASTRA-sim System Layer**

# Recall: Workload Layer

```
void Workload::issue_comm(node) {
    hw_resource->occupy(node);

    if (node->comm_type() == ChakraCollectiveCommType::ALL_REDUCE) {
        DataSet* fp = sys->generate_all_reduce(node->comm_size(), ...)

        fp->set_notifier(EventType::CollectiveCommunicationFinished);

    }

    (...)
}
```
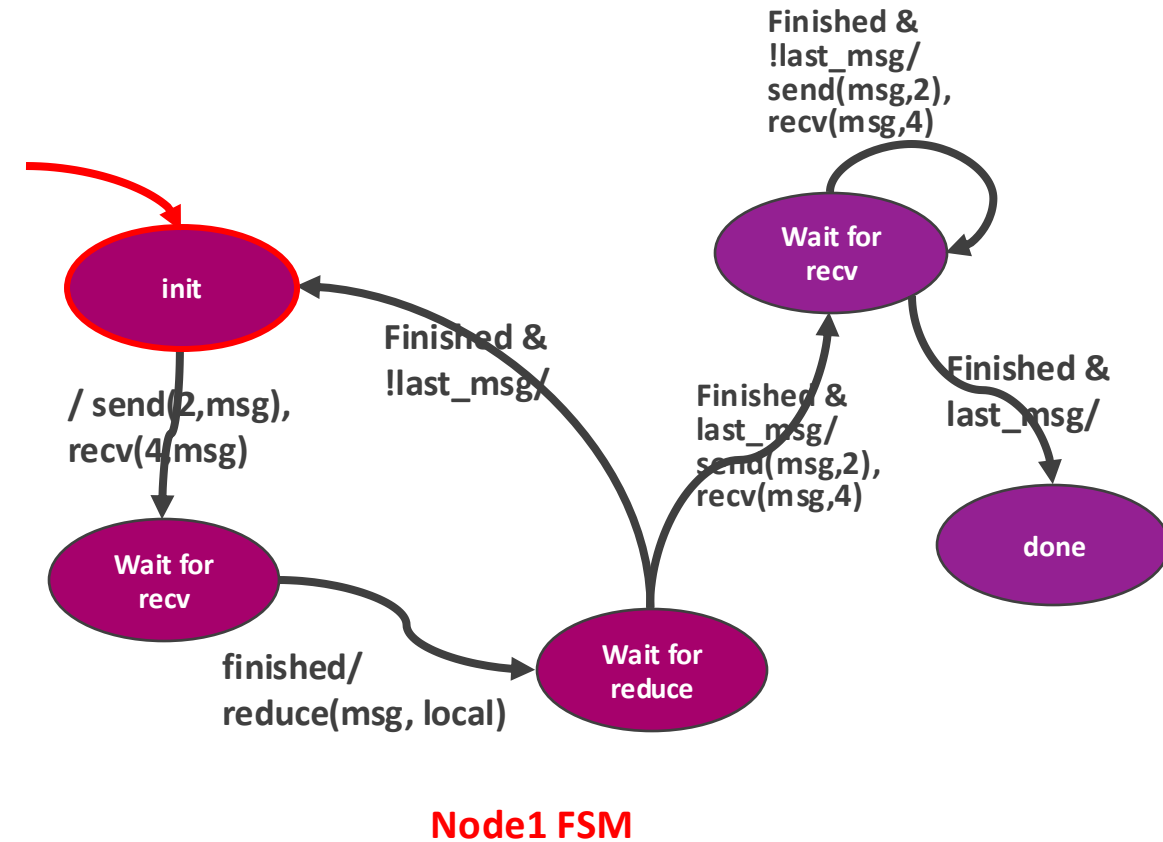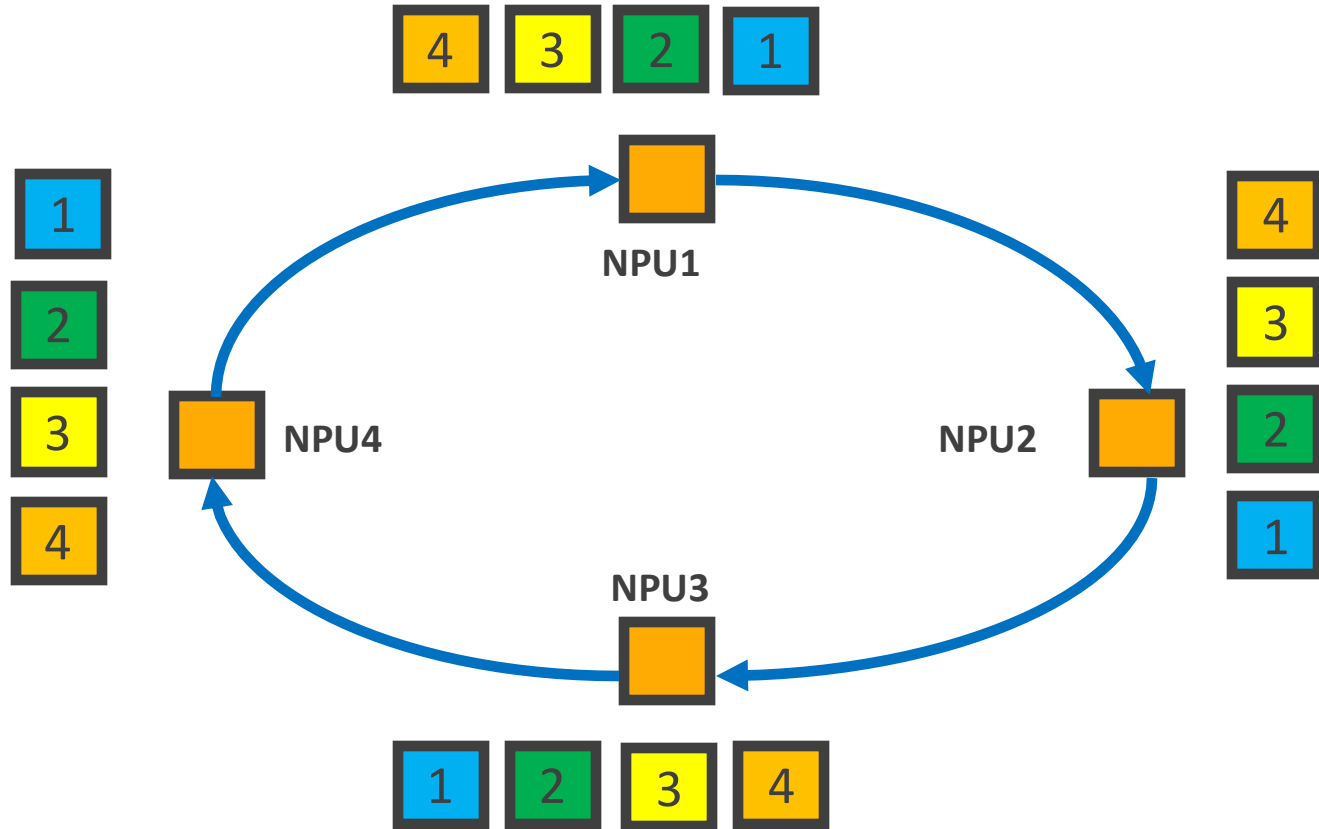
**Managed by the System Layer**

# System Layer

- **Collective algorithms** are implemented
  - Collectives are broken down into individual **chunk send-recv**

- Collective algorithms are managed via:
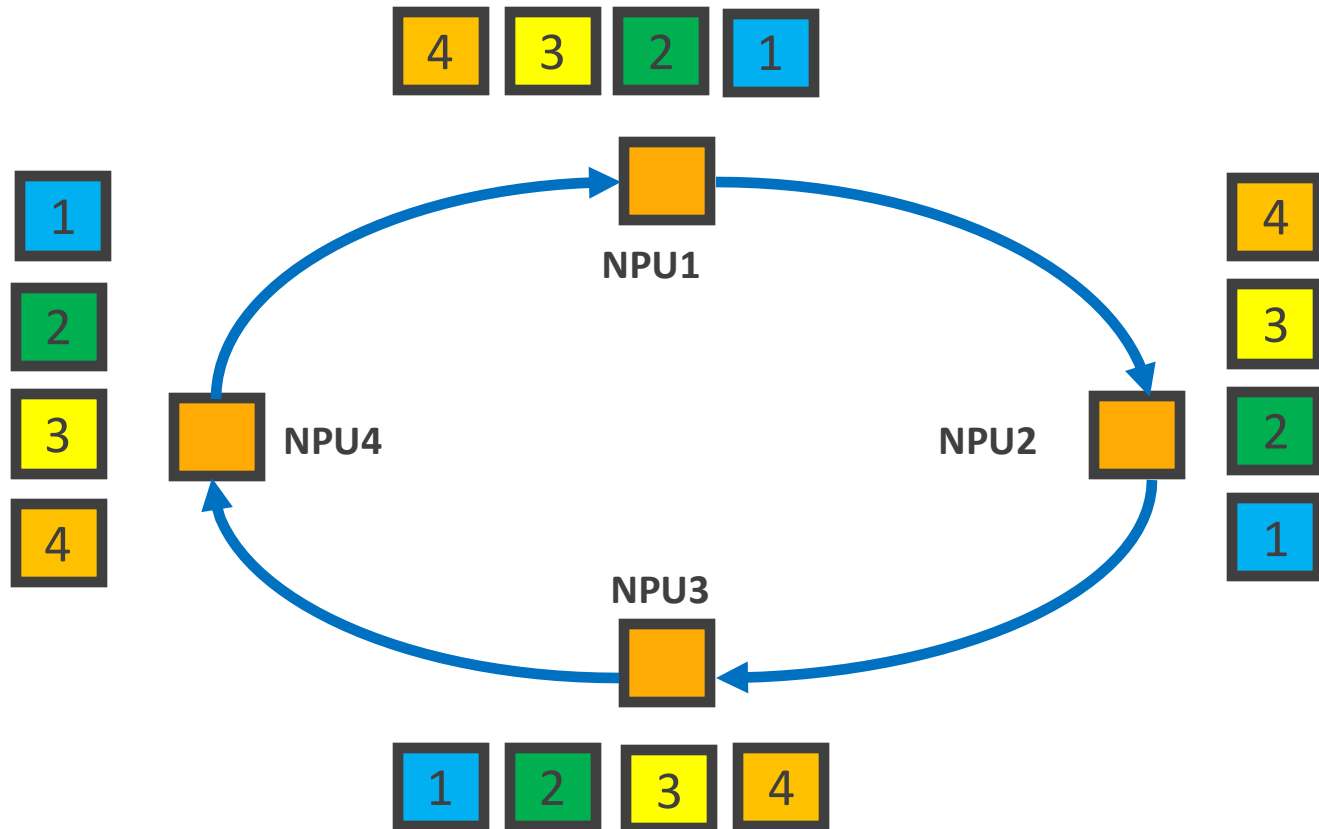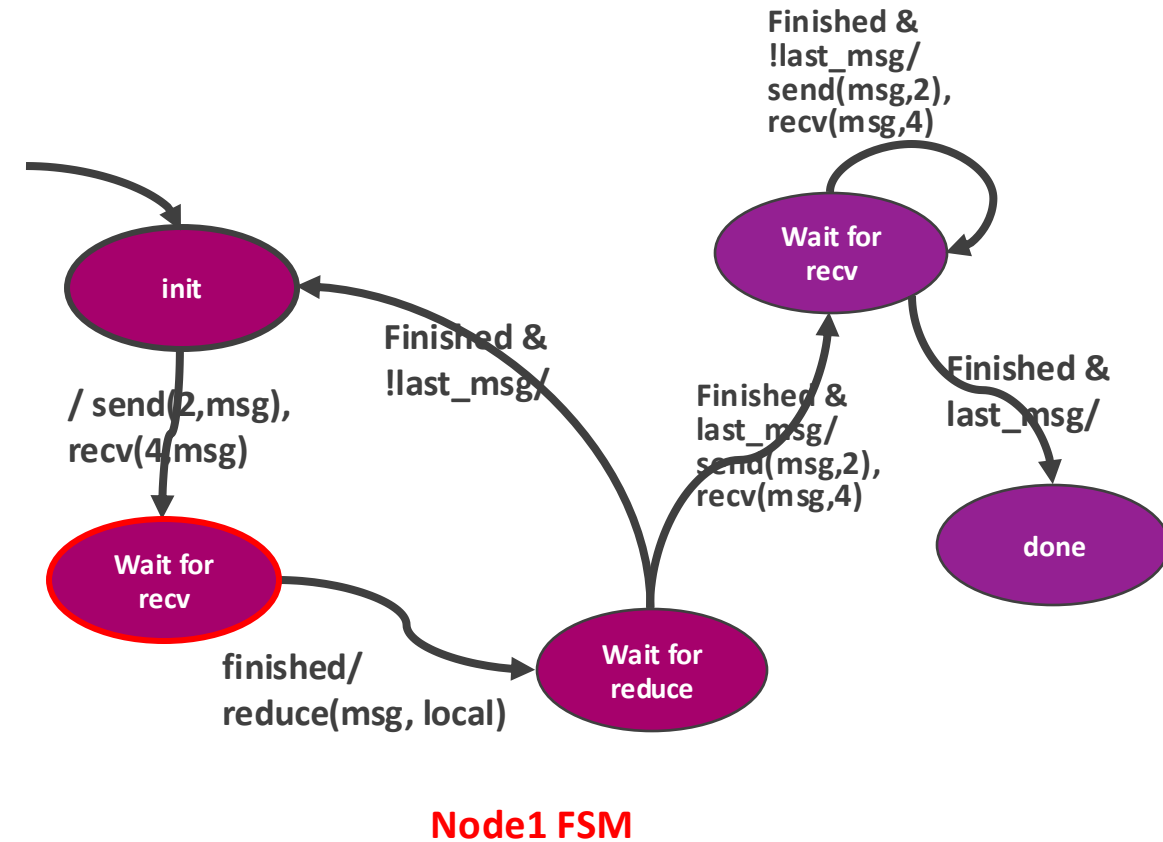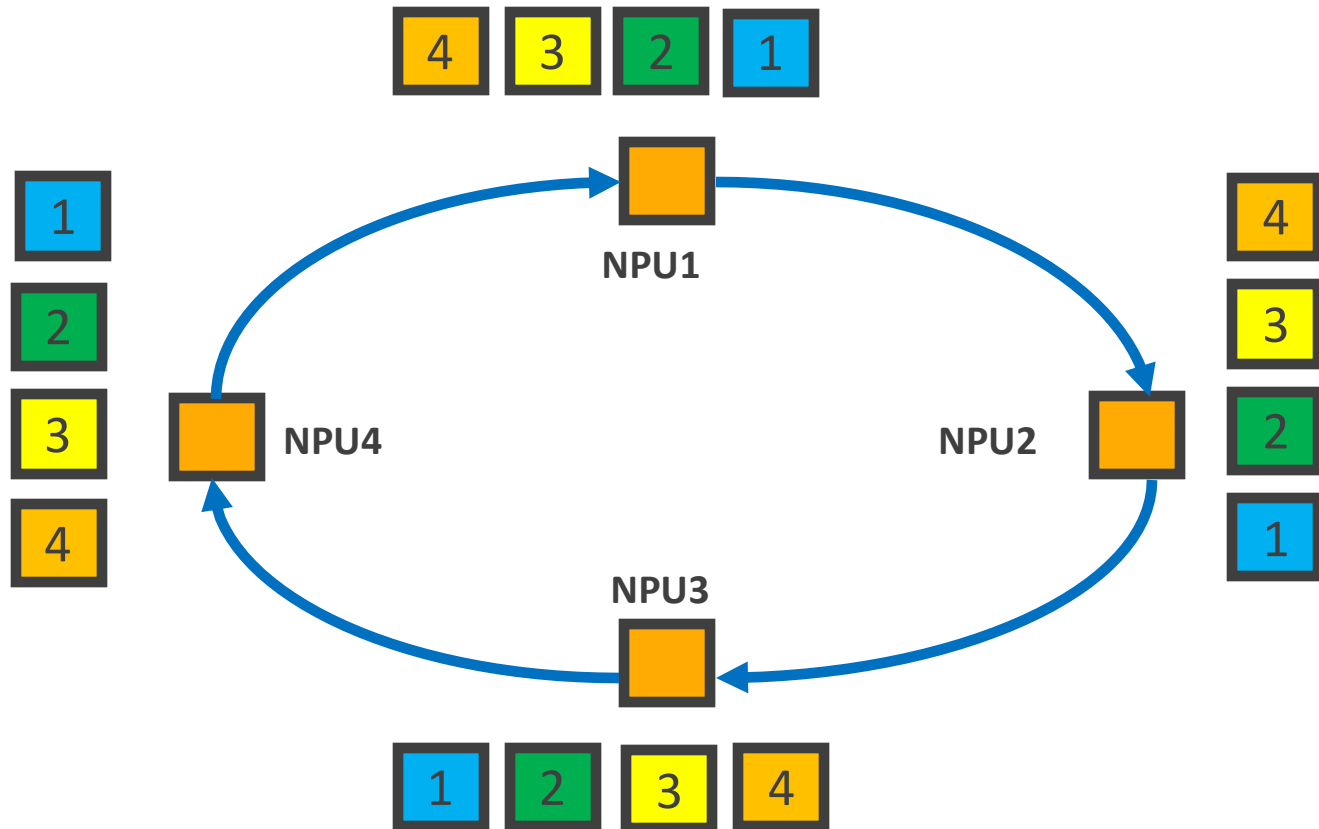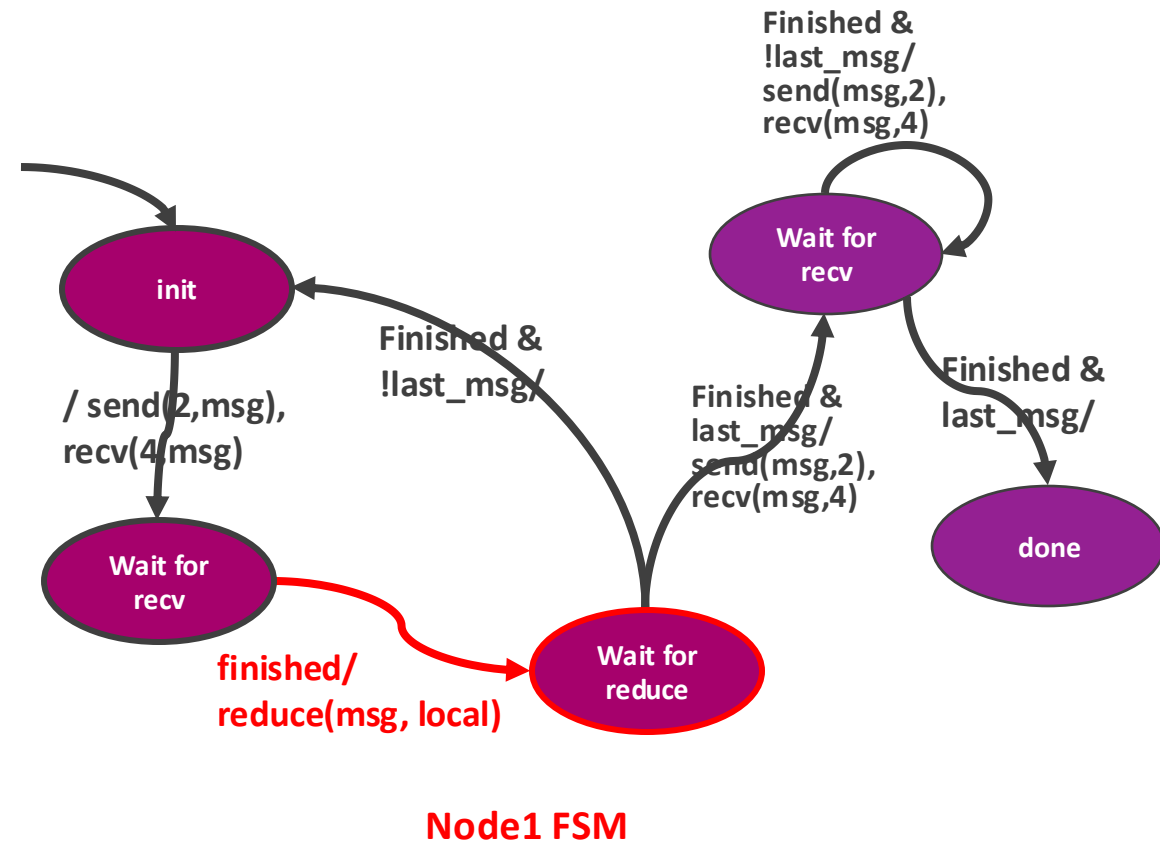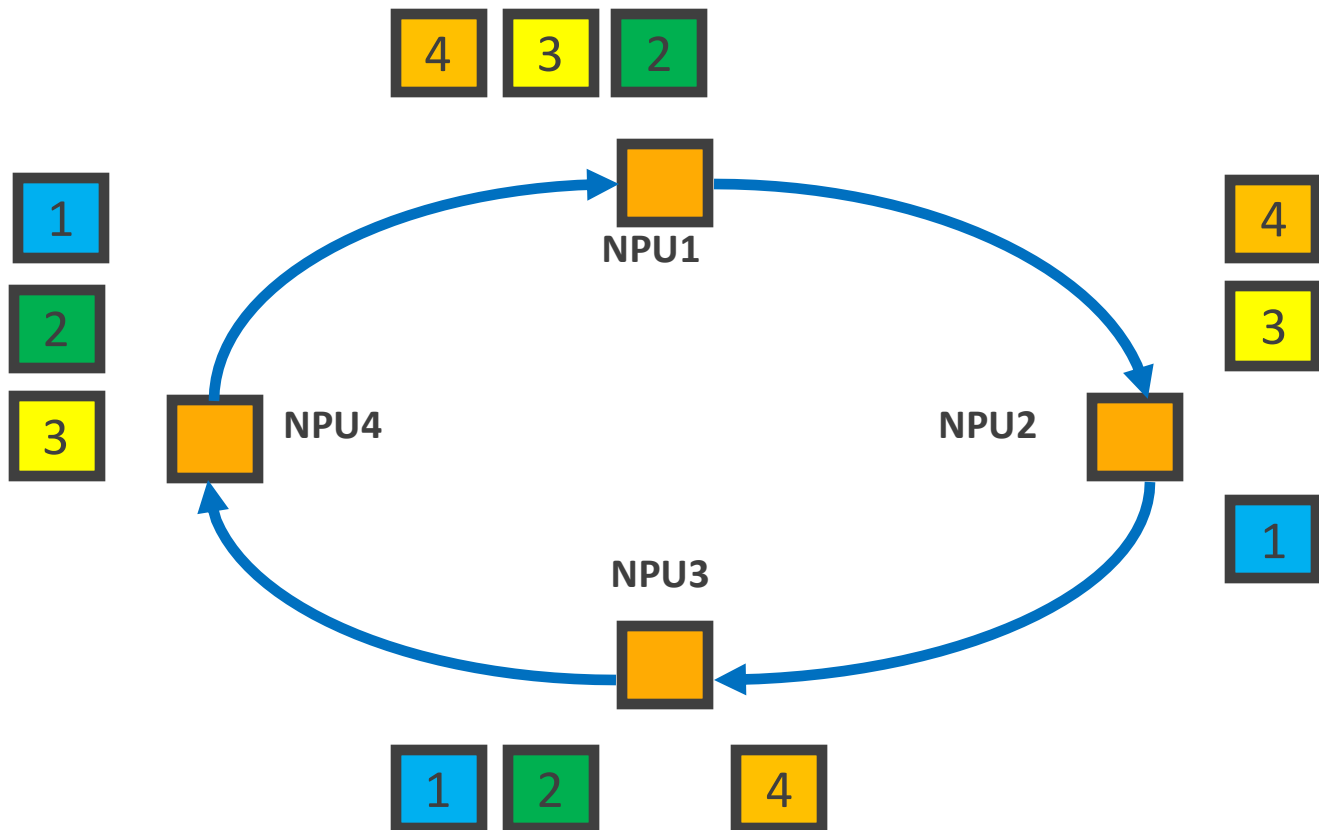  - **Finite State Machine**

# System Layer Collective Implementation

- Collective algorithms can be implemented using **state machines**.



**Node1 FSM**

# System Layer Collective Implementation

- Collective algorithms can be implemented using **state machines**.
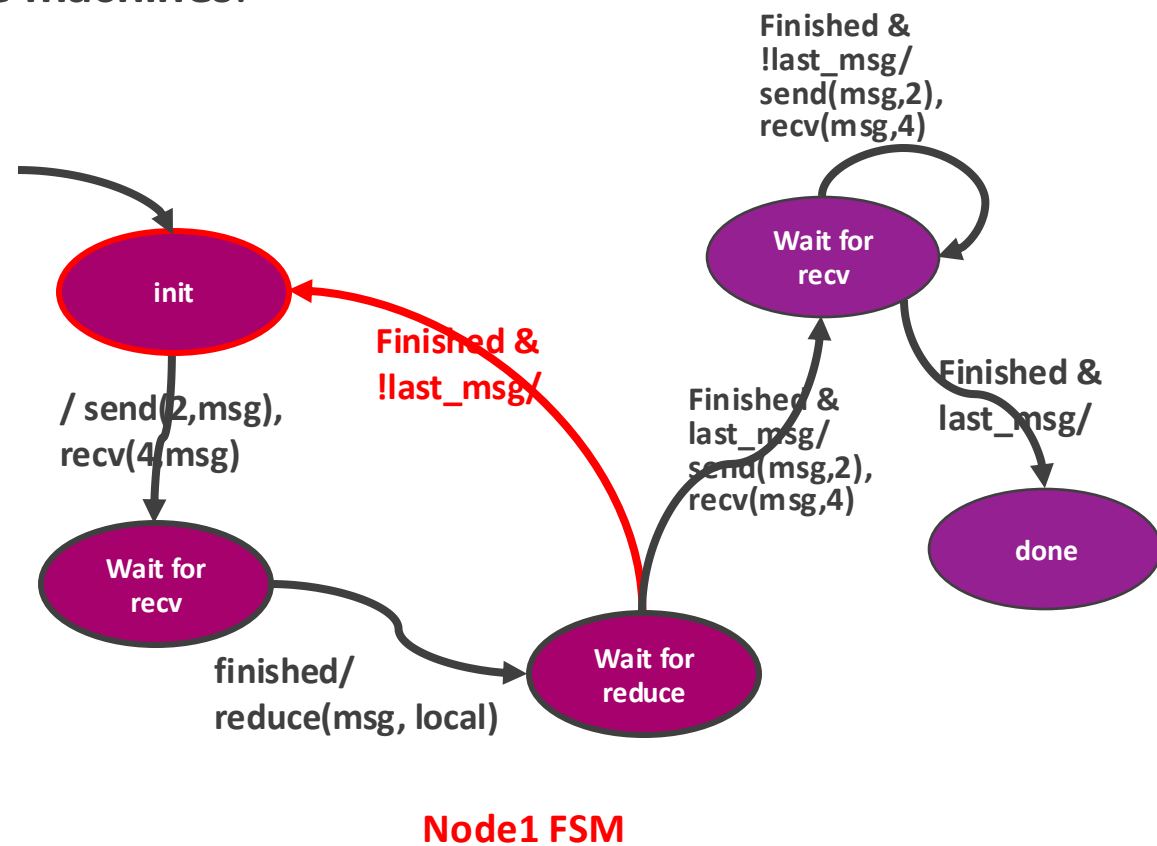


**Node1 FSM**

# System Layer Collective Implementation

- Collective algorithms can be implemented using **state machines**.
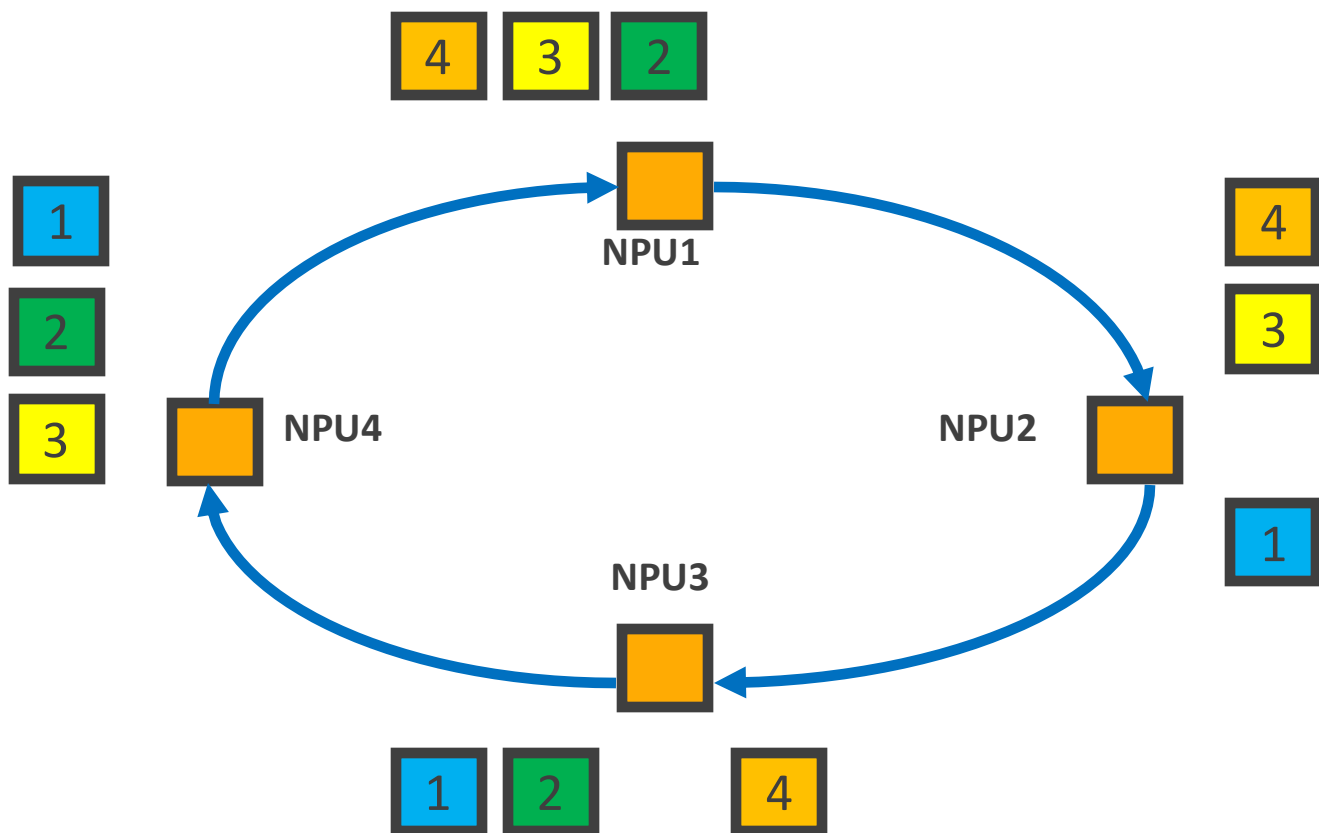


**Node1 FSM**

# System Layer Collective Implementation

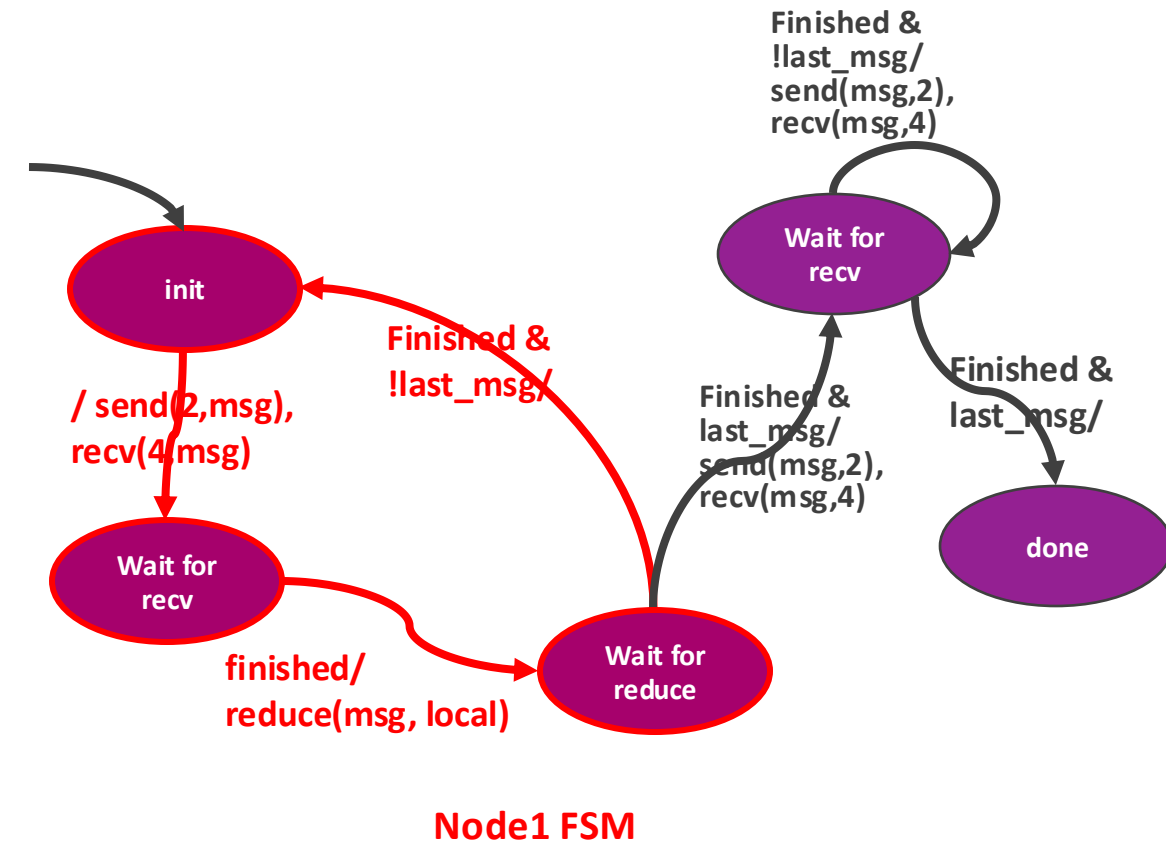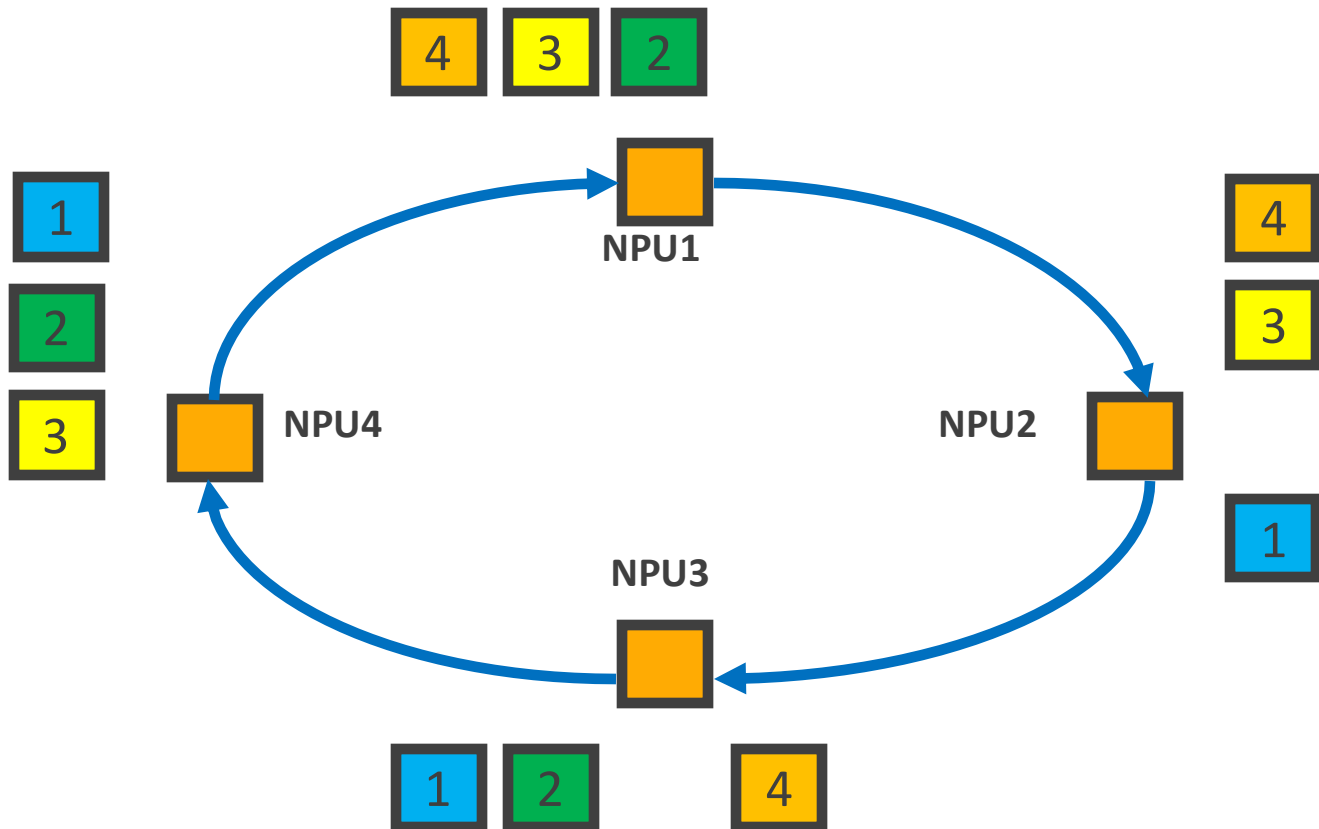- Collective algorithms can be implemented using **state machines**.

# System Layer Collective Implementation

- Collective algorithms can be implemented using **state machines**.

# System Layer Collective Implementation

- Collective algorithms can be implemented using **state machines**.
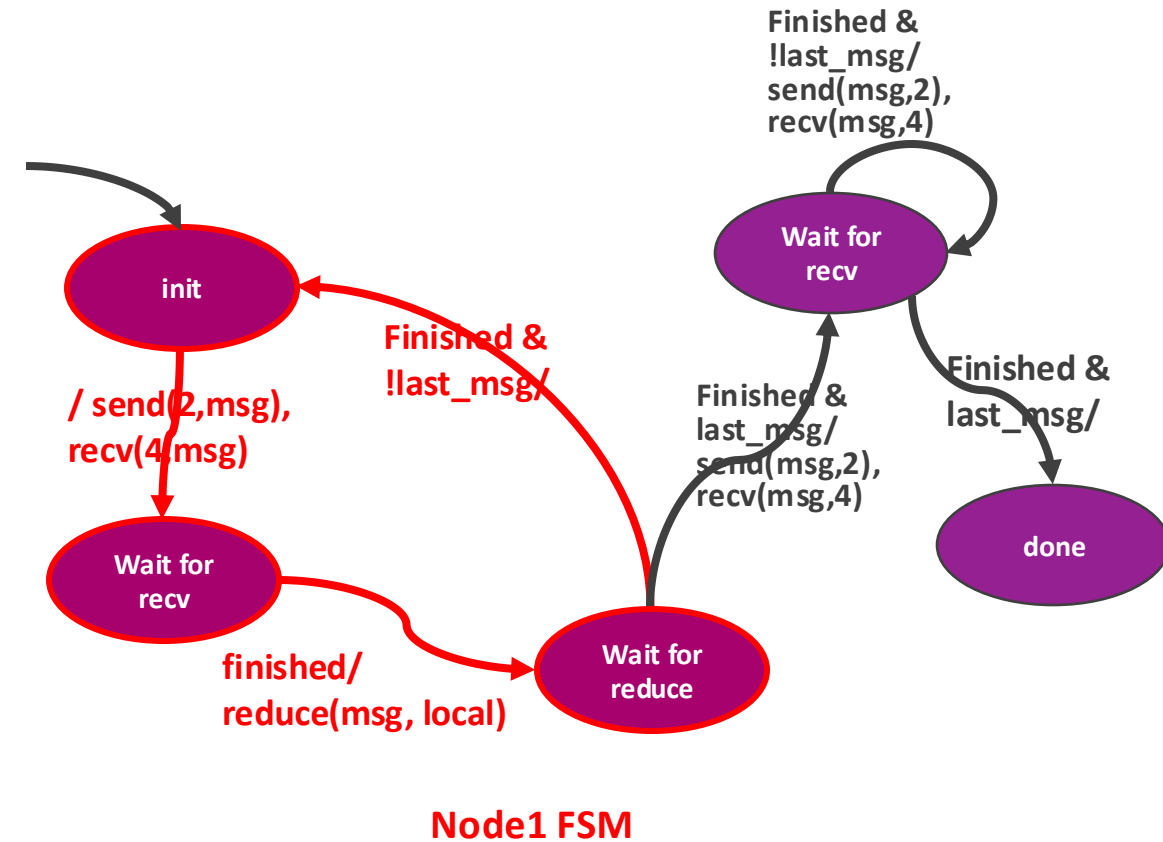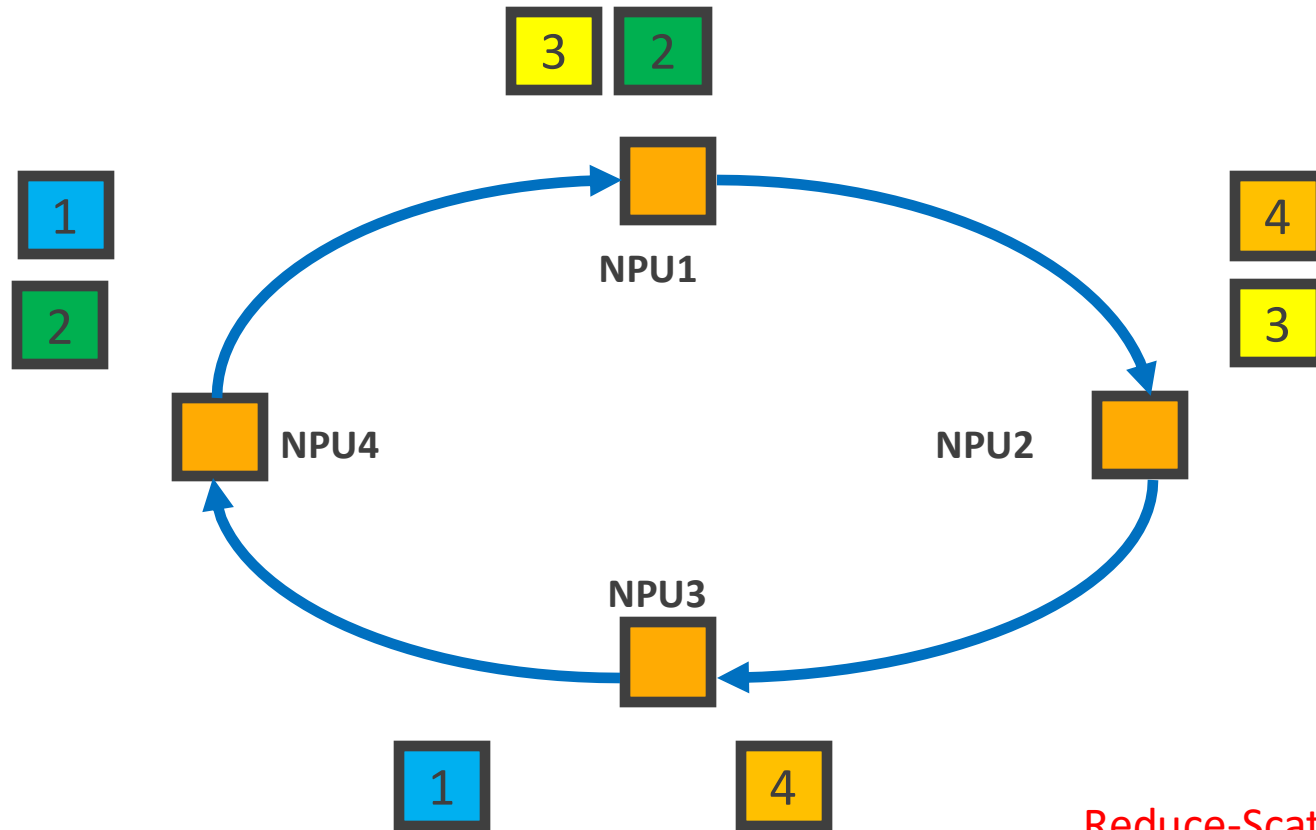


**Node1 FSM**

# System Layer Collective Implementation

- Collective algorithms can be implemented using **state machines**.
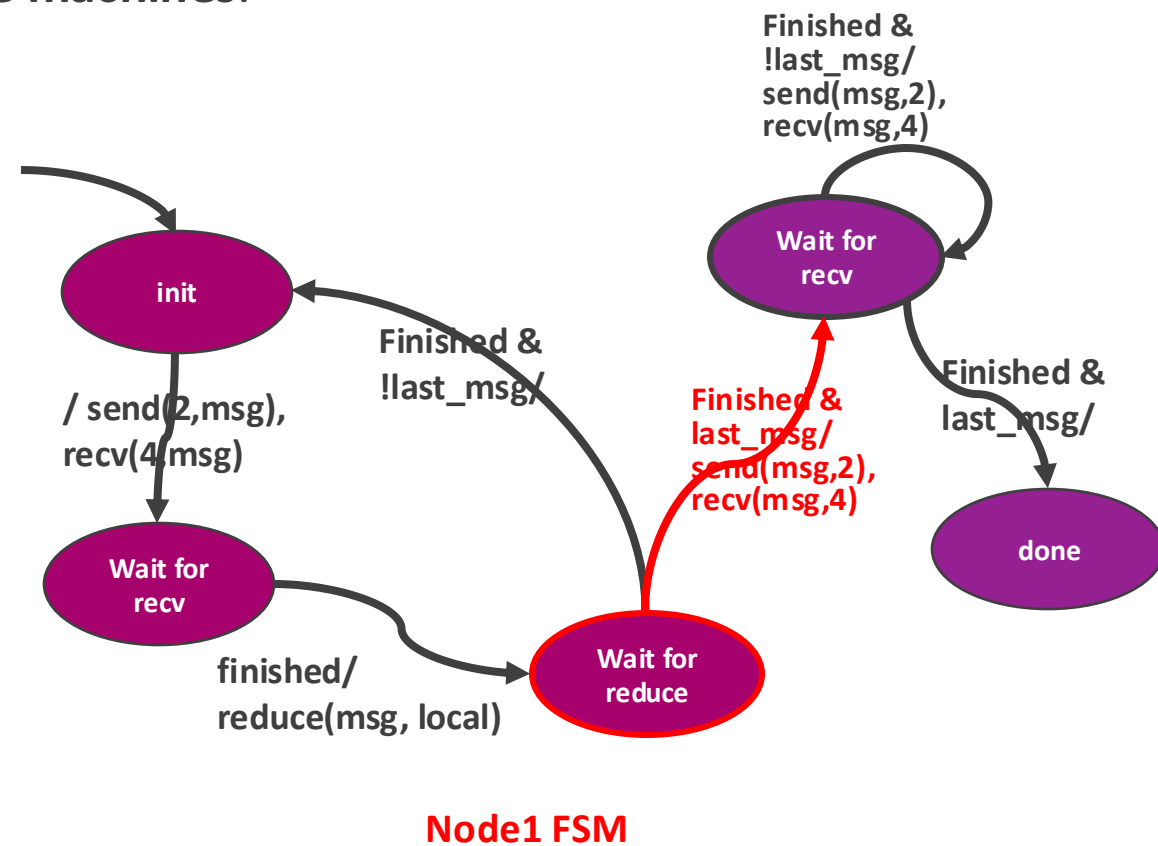


Reduce-Scatter done!

Node1 FSM

# System Layer Collective Implementation

- Collective algorithms can be implemented using **state machines**.



**Node1 FSM**

# System Layer Collective Implementation

- Collective algorithms can be implemented using **state machines**.
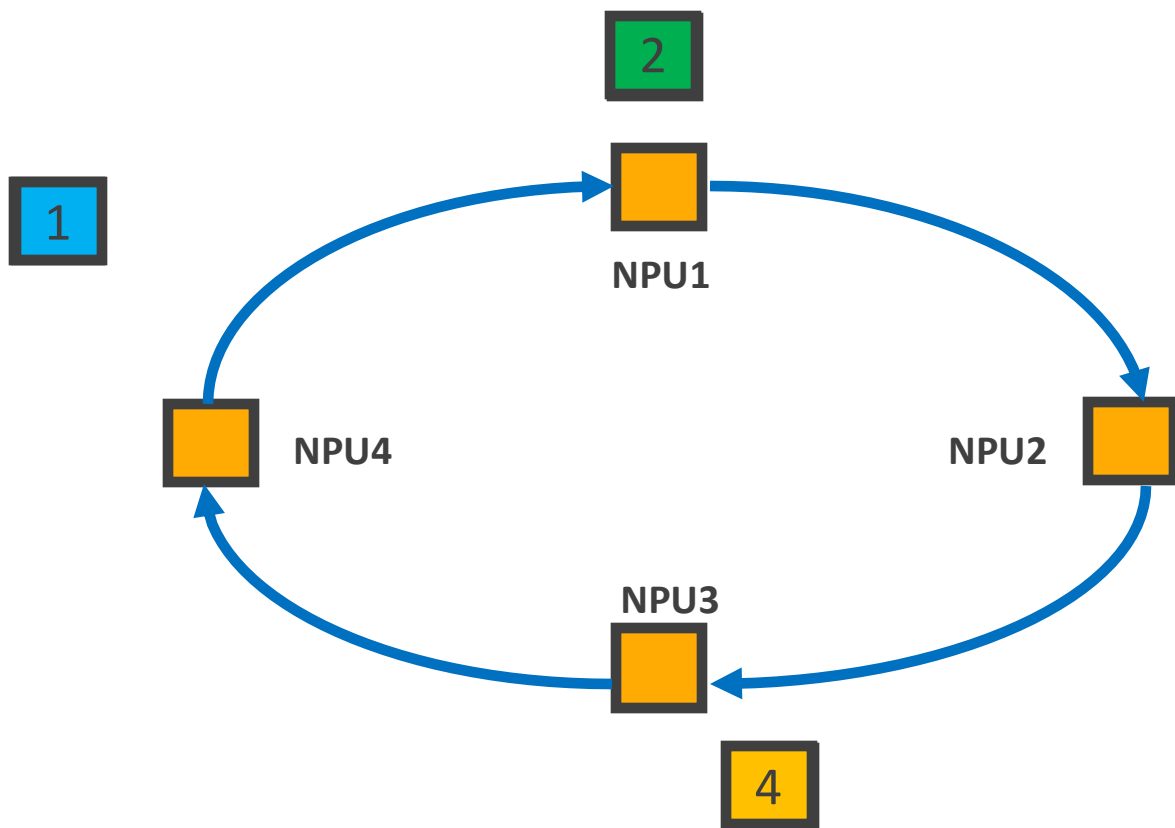


**Node1 FSM**

# System Layer Collective Implementation

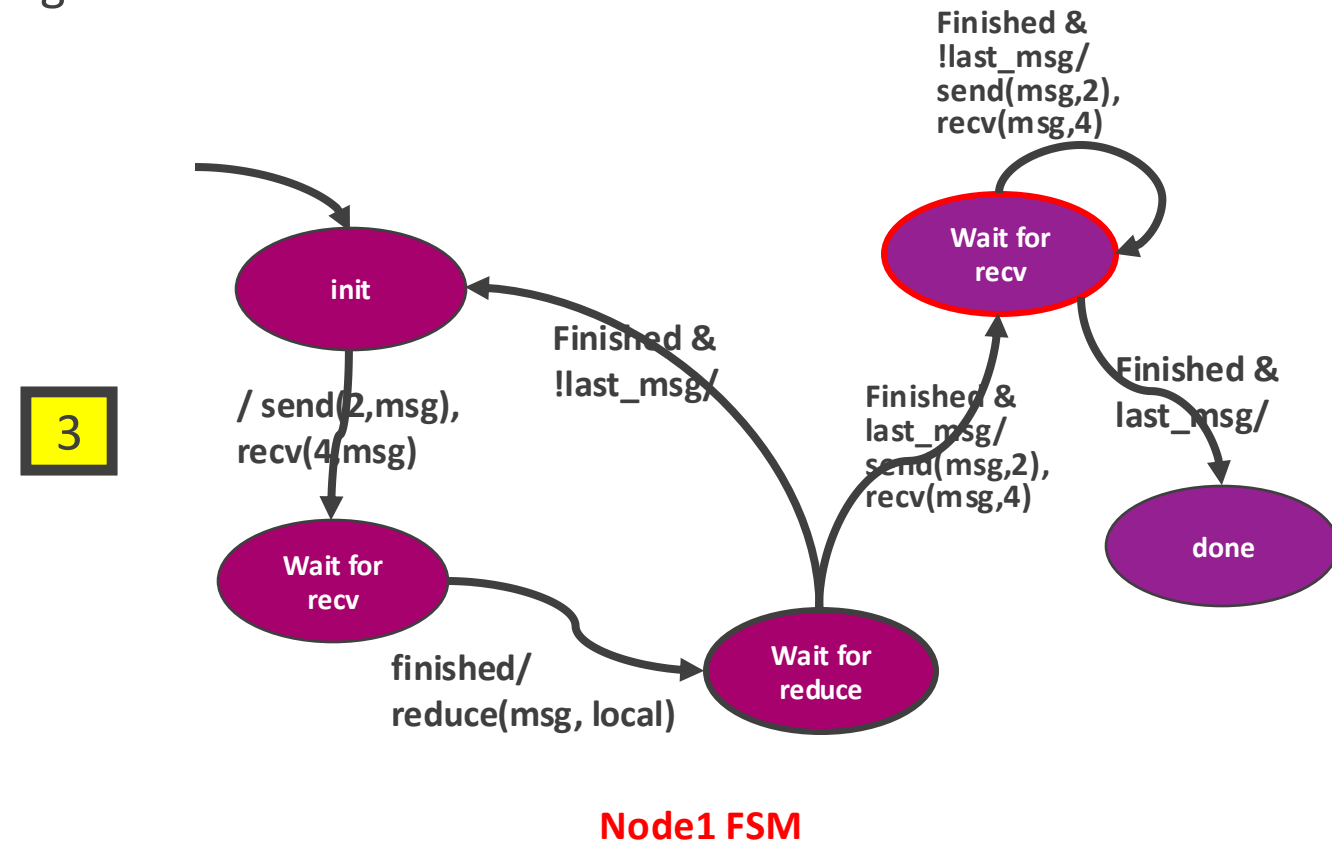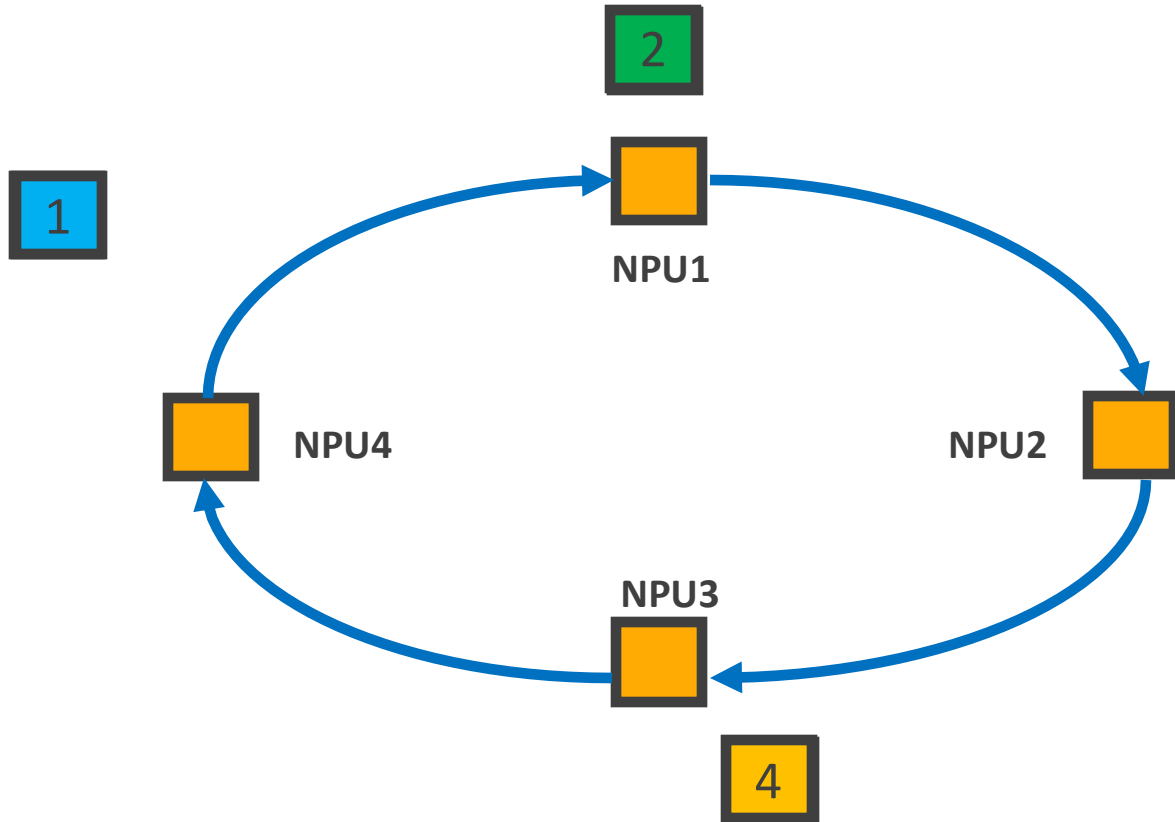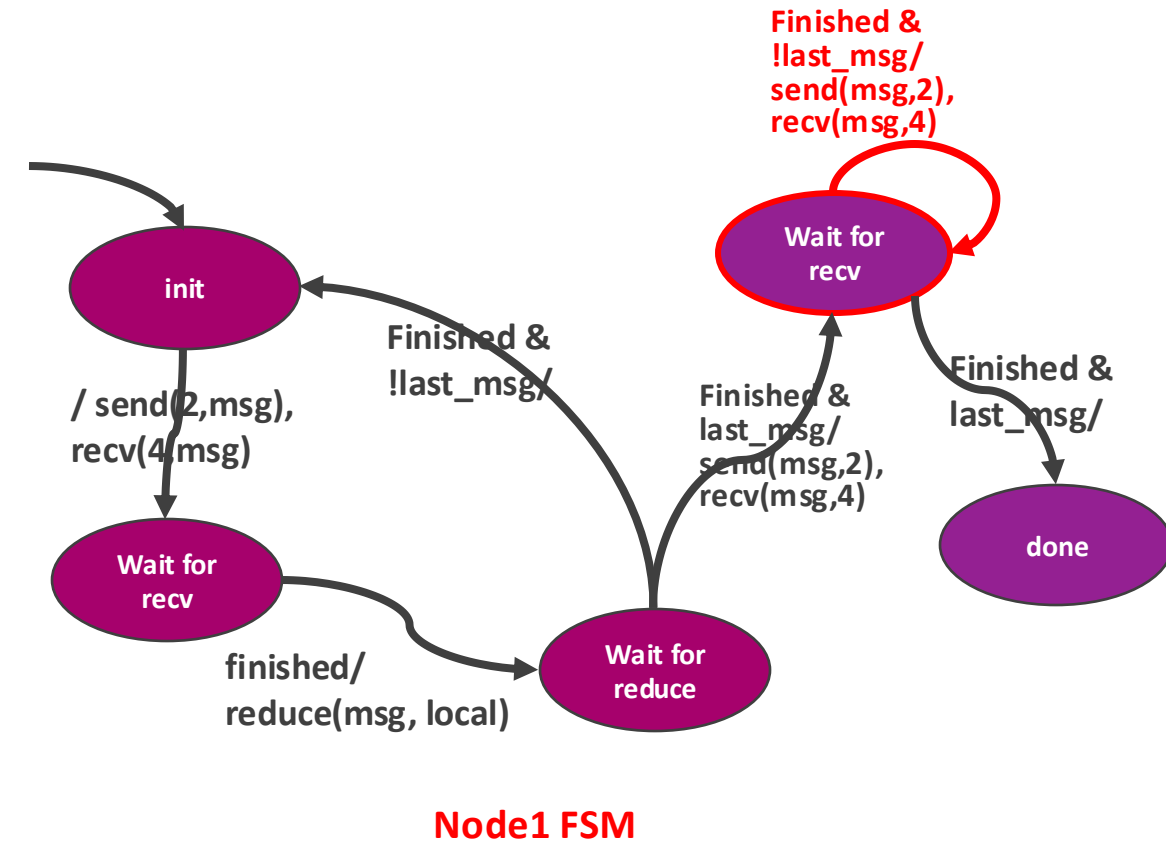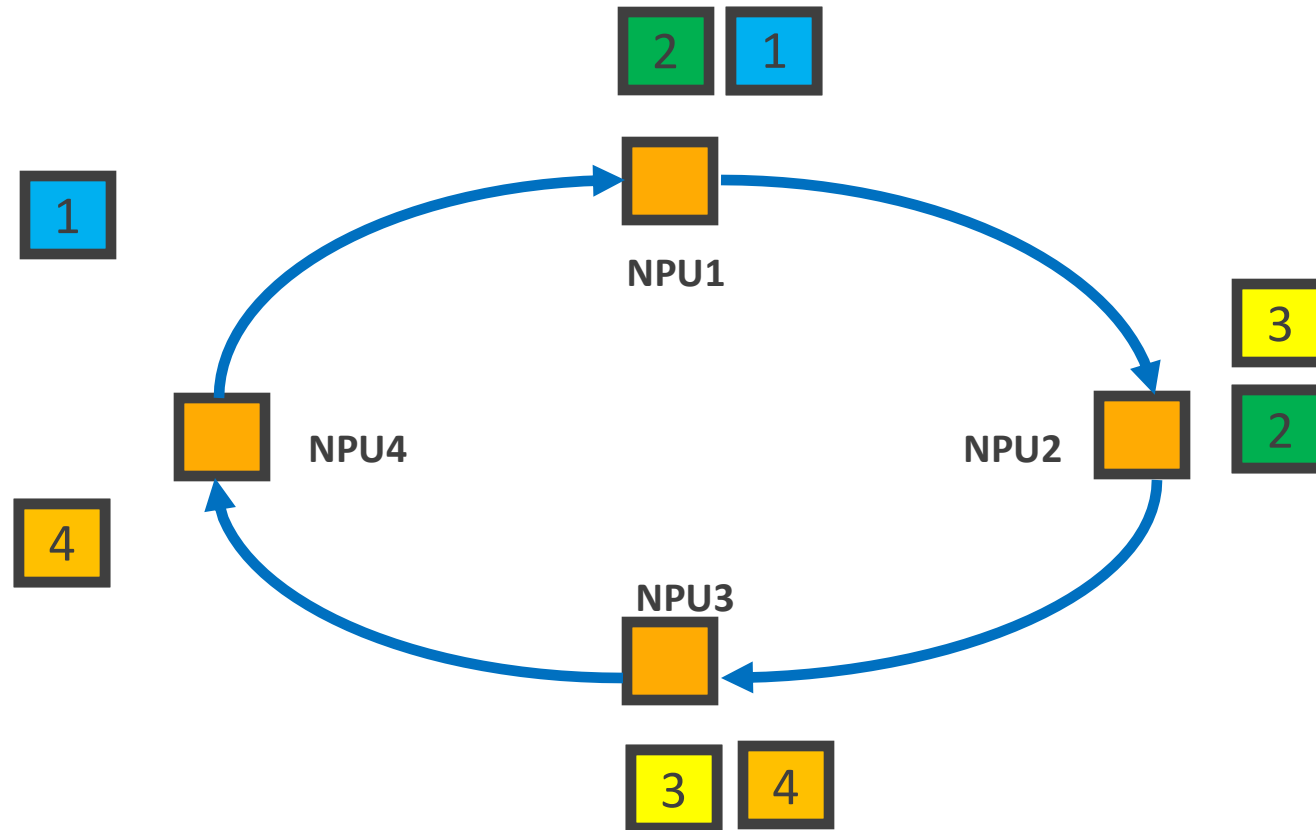- Collective algorithms can be implemented using **state machines**.

# System Layer Collective Implementation

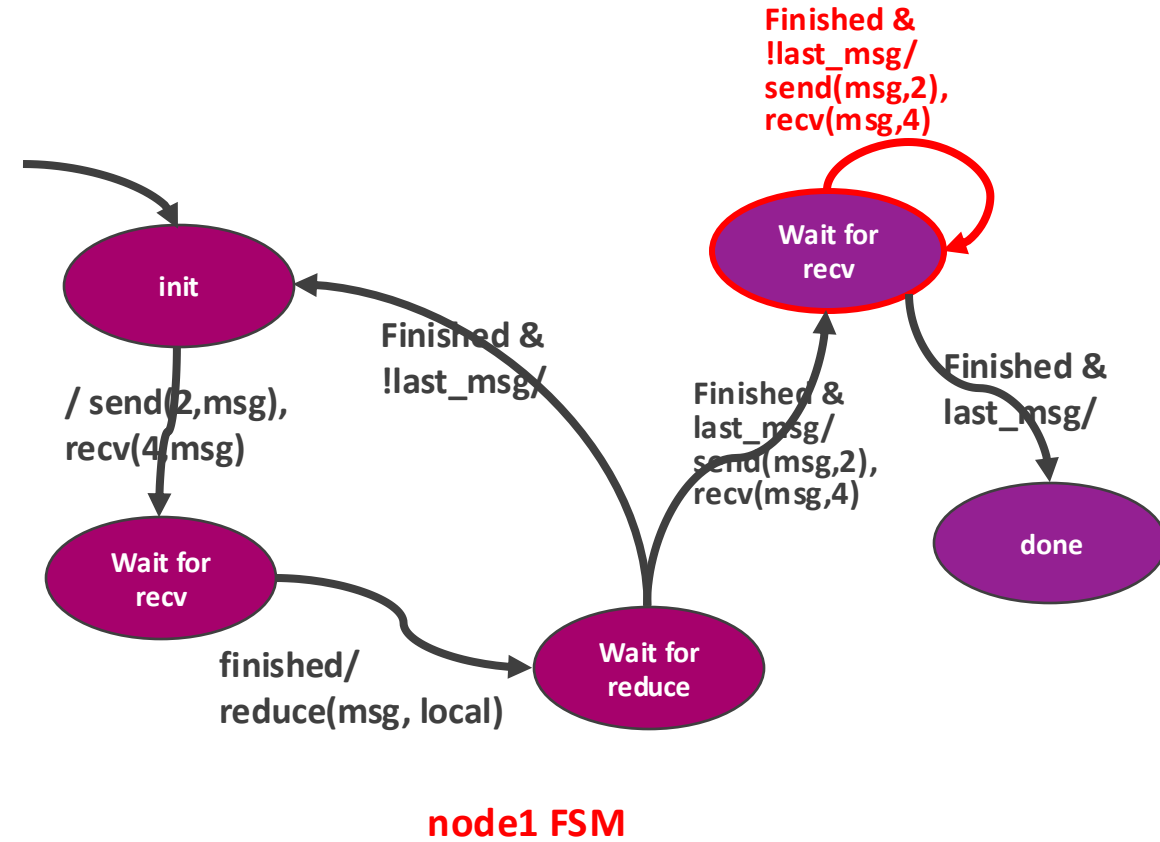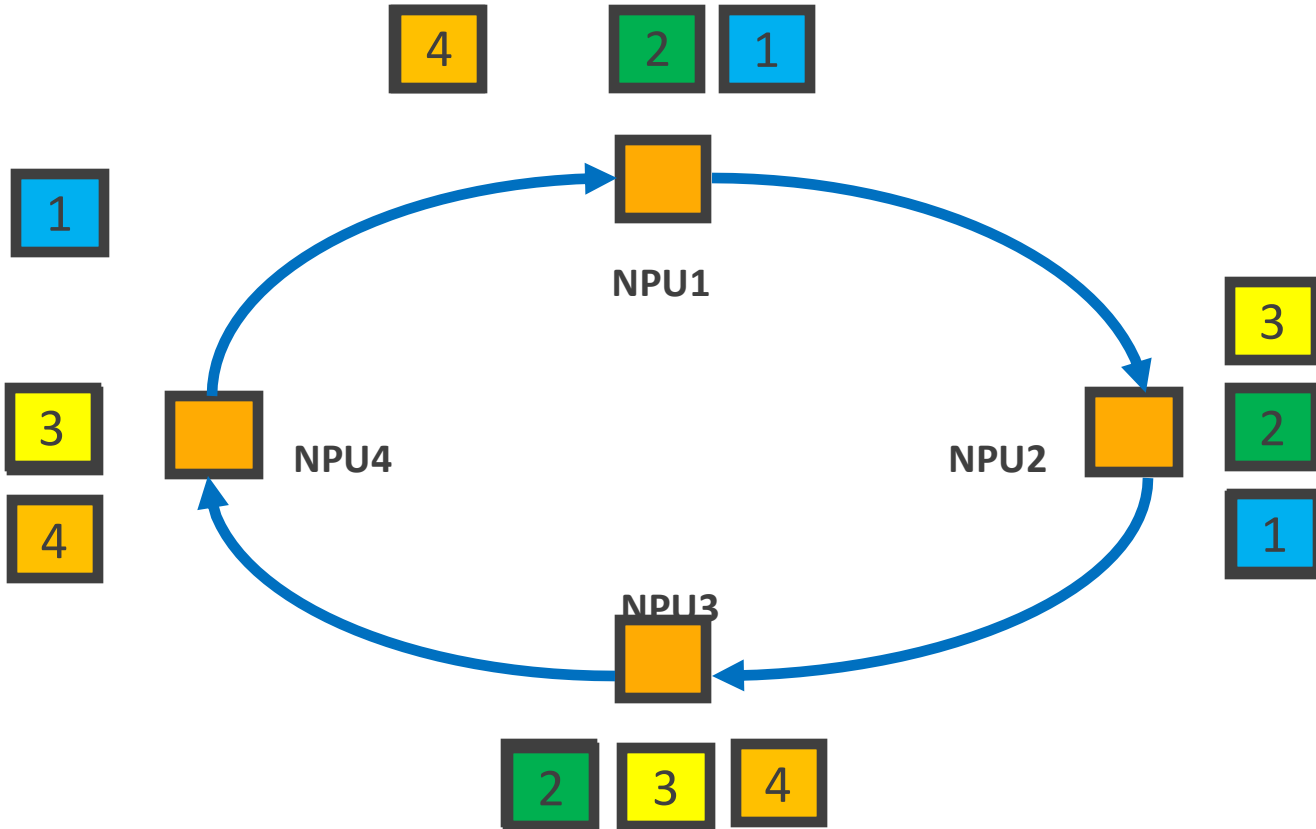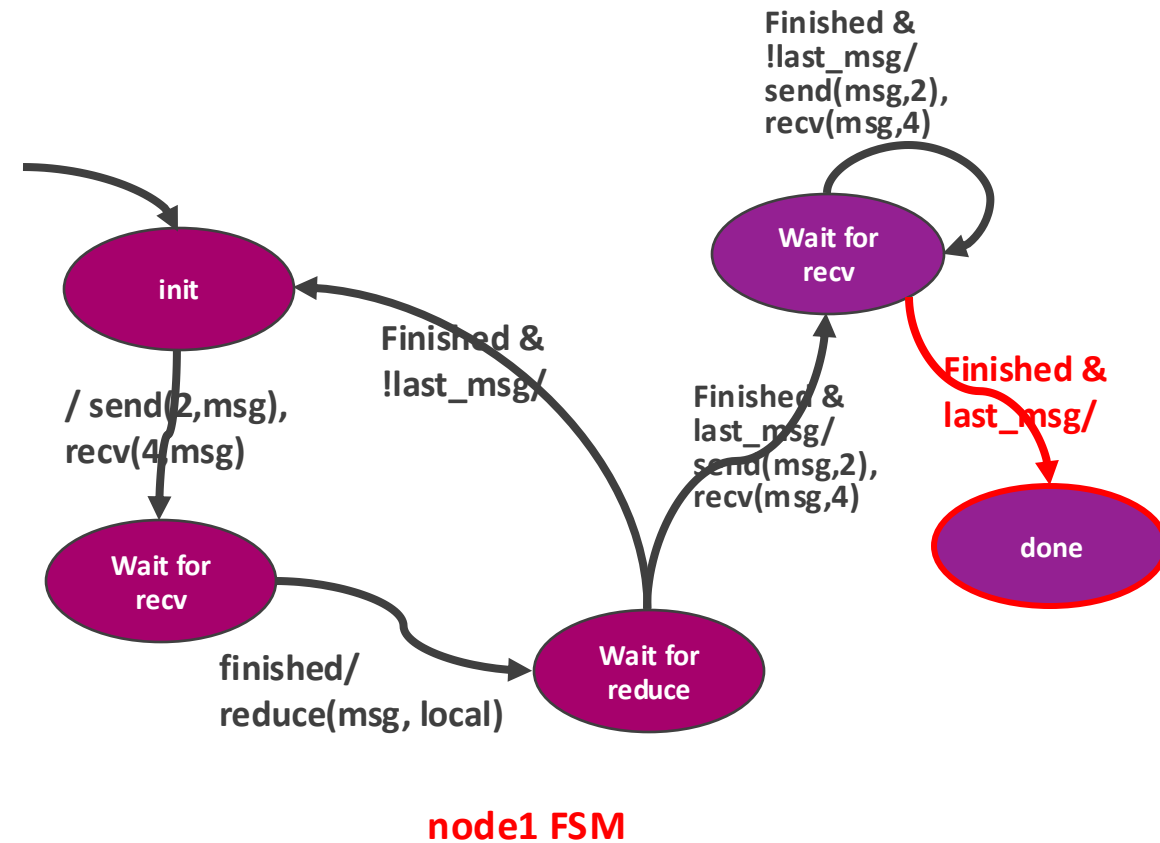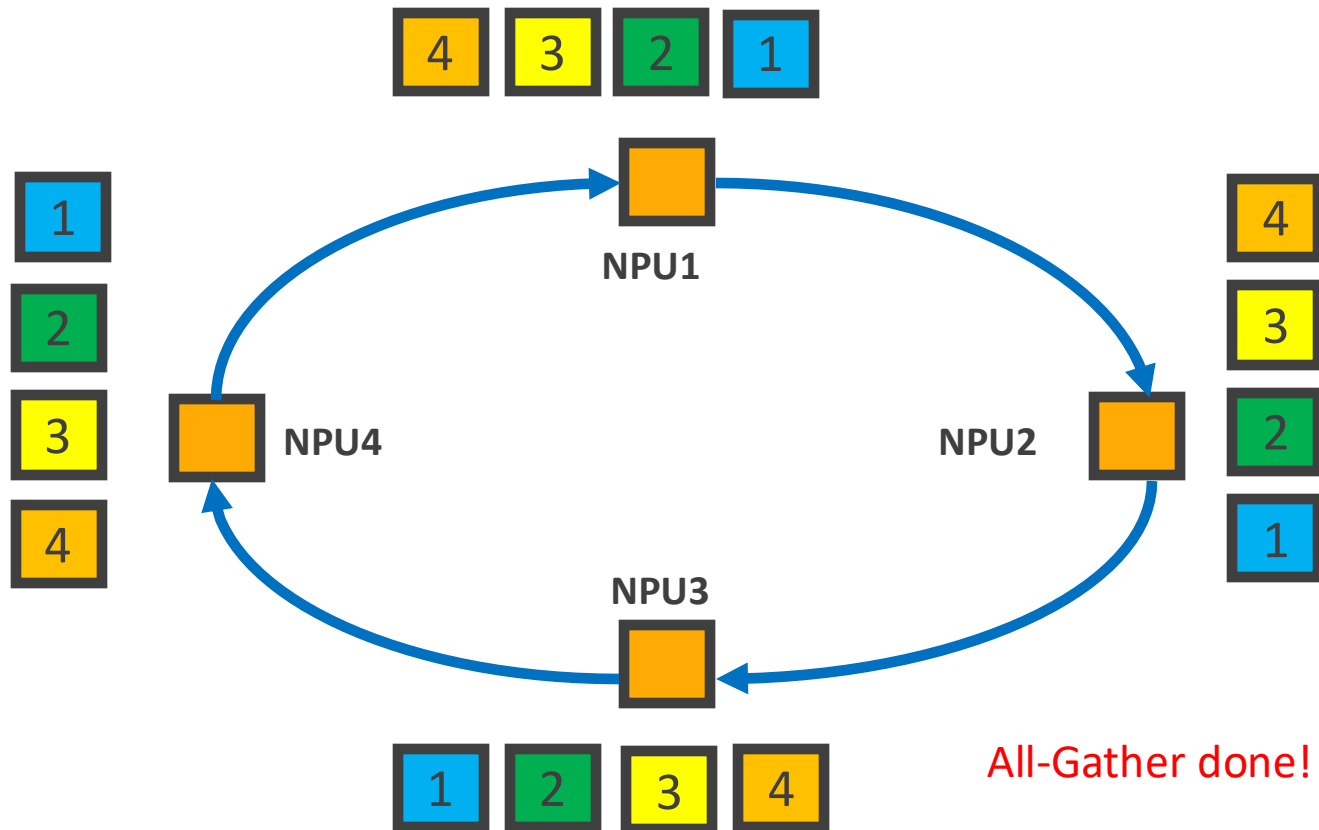- Collective algorithms can be implemented using **state machines**.

# System Layer Collective Implementation

- Collective algorithms can be implemented using **state machines**.



All-Gather done!

**node1 FSM**

# System Input

Collective Policy for

```
sample_torus_sys.txt        ×

1   scheduling-policy: LIFO
2   endpoint-delay: 1
3   active-chunks-per-dimension: 1
4   preferred-dataset-splits: 4
5   boost-mode: 0
6   all-reduce-implementation: ring_ring_ring
7   all-gather-implementation: ring_ring_ring
8   reduce-scatter-implementation: ring_ring_ring
9   all-to-all-implementation: ring_ring_ring
10  collective-optimization: localBWAware
11
```

NPU    Intra-package scale-up
Inter-package scale-up   Package

# System Input

Constant delay before NPU sending a message



```
sample_torus_sys.txt                    ×

1   scheduling-policy: LIFO
2   endpoint-delay: 1
3   active-chunks-per-dimension: 1
4   preferred-dataset-splits: 4
5   boost-mode: 0
6   all-reduce-implementation: ring_ring_ring
7   all-gather-implementation: ring_ring_ring
8   reduce-scatter-implementation: ring_ring_ring
9   all-to-all-implementation: ring_ring_ring
10  collective-optimization: localBWAware
11
```

# System Input

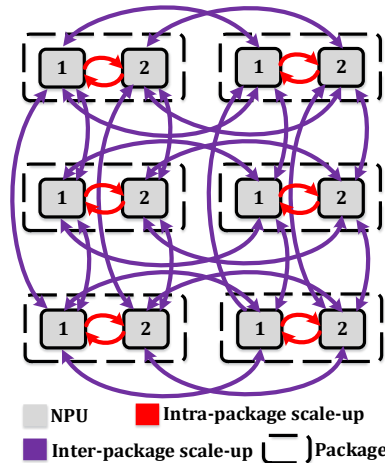Max running chunks per each physical network dimension



```
sample_torus_sys.txt                    ×

 1  scheduling-policy: LIFO
 2  endpoint-delay: 1
 3  active-chunks-per-dimension: 1
 4  preferred-dataset-splits: 4
 5  boost-mode: 0
 6  all-reduce-implementation: ring_ring_ring
 7  all-gather-implementation: ring_ring_ring
 8  reduce-scatter-implementation: ring_ring_ring
 9  all-to-all-implementation: ring_ring_ring
10  collective-optimization: localBWAware
11
```
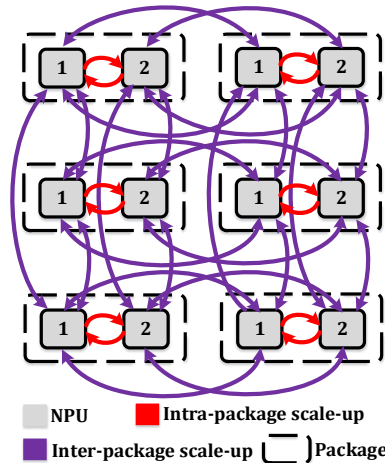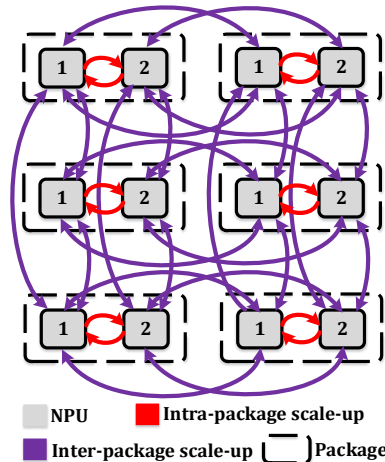
NPU  ▇ Intra-package scale-up
▇ Inter-package scale-up ⌐ ̄ ⌐ Package

# System Input

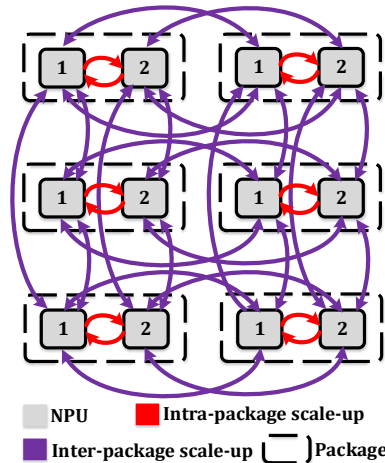# of chunks to split each collective into

```
sample_torus_sys.txt                    ×

 1  scheduling-policy: LIFO
 2  endpoint-delay: 1
 3  active-chunks-per-dimension: 1
 4  preferred-dataset-splits: 4
 5  boost-mode: 0
 6  all-reduce-implementation: ring_ring_ring
 7  all-gather-implementation: ring_ring_ring
 8  reduce-scatter-implementation: ring_ring_ring
 9  all-to-all-implementation: ring_ring_ring
10  collective-optimization: localBWAware
11
```



NPU   Intra-package scale-up
Inter-package scale-up   )Package

# System Input



Speed-up the simulation

```
sample_torus_sys.txt    ×

 1  scheduling-policy: LIFO
 2  endpoint-delay: 1
 3  active-chunks-per-dimension: 1
 4  preferred-dataset-splits: 4
 5  boost-mode: 0
 6  all-reduce-implementation: ring_ring_ring
 7  all-gather-implementation: ring_ring_ring
 8  reduce-scatter-implementation: ring_ring_ring
 9  all-to-all-implementation: ring_ring_ring
10  collective-optimization: localBWAware
11
```
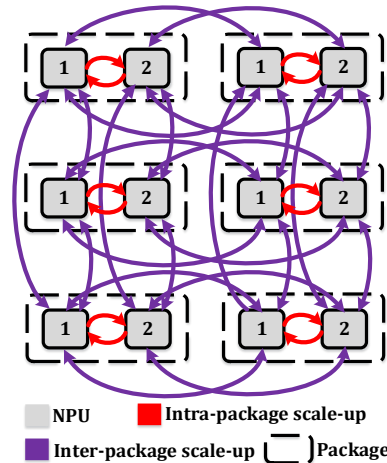
# System Input

Hierarchical collective algorithm implementation



```
sample_torus_sys.txt                    ×

1   scheduling-policy: LIFO
2   endpoint-delay: 1
3   active-chunks-per-dimension: 1
4   preferred-dataset-splits: 4
5   boost-mode: 0
6   all-reduce-implementation: ring_ring_ring
7   all-gather-implementation: ring_ring_ring
8   reduce-scatter-implementation: ring_ring_ring
9   all-to-all-implementation: ring_ring_ring
10  collective-optimization: localBWAware
11
```

NPU    Intra-package scale-up

Inter-package scale-up   Package

# Sneak Peek: Collective API

- ## MSCCLang high-level DSE to represent collective algorithm

```
def allgather_ring(size, channels, instances, protocol):
    (...)
    for step in range(0, size-1):
        for index in range(0, size):
                rank = (index + step) % size
                next_rank = (index + step + 1) % size
                c = chunk(rank, Buffer.output, index)
                c = c.copy(next_rank, Buffer.output, index, sendtb=channel, recvtb=channel, ch=channel)

    XML()
```

← **Represent arbitrary collective algorithm and store it in standardized XML format**

- ## Custom collective algorithms represented by MSCCL-xml

```xml
<algo name="allreduce_ring_1channelsperring" (...) >
<gpu id="0" i_chunks="4" o_chunks="0" s_chunks="0">

<tb id="0" send="1" recv="3" chan="0">
<step s="0" type="s" srcbuf="i" srcoff="0" dstbuf="i" dstoff="0" cnt="1" depid="-1" deps="-1" hasdep="0"/>
<step s="1" type="rrc" srcbuf="i" srcoff="3" dstbuf="i" dstoff="3" cnt="1" depid="-1" deps="-1" hasdep="0"/>
<step s="2" type="s" srcbuf="i" srcoff="3" dstbuf="i" dstoff="3" cnt="1" depid="-1" deps="-1" hasdep="0"/>
<step s="3" type="rrc" srcbuf="i" srcoff="2" dstbuf="i" dstoff="2" cnt="1" depid="-1" deps="-1" hasdep="0"/>

(...)
```
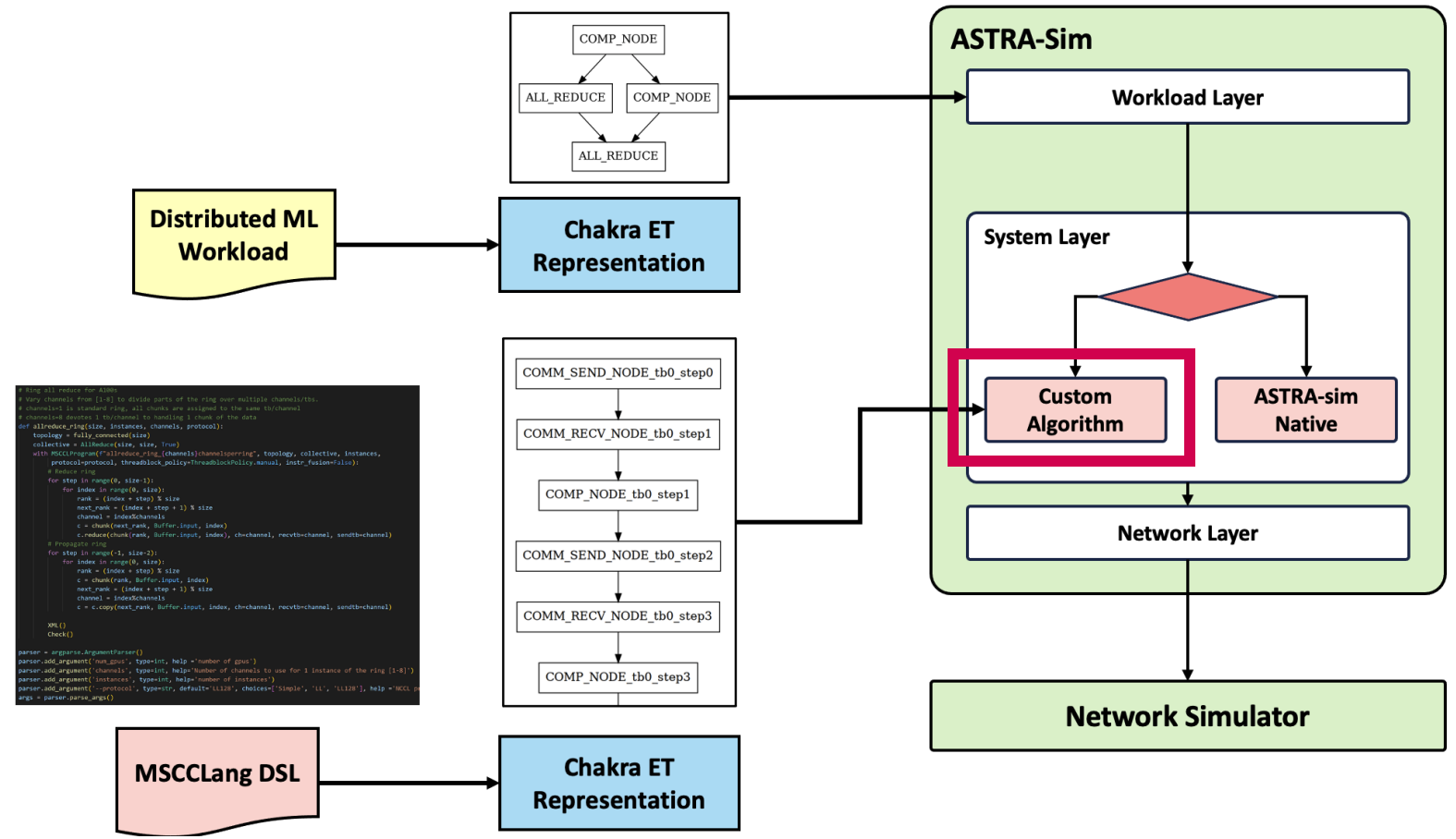
# Sneak Peek: Collective API

**Work in Progress!**

- HotI paper: *https://arxiv.org/abs/2408.11008*



*Slide courtesy: Jinsun Yoo <jinsun@gatech.edu>*

# Sneak Peek: TACOS

**To appear in MICRO '24!**
- https://arxiv.org/abs/2304.05301

- A mechanism to generate **topology-aware custom collective plan**



CollectiveAPI

TACOS custom collective plan

*Slide courtesy: Jinsun Yoo <jinsun@gatech.edu>*