# Themis: A Network Bandwidth-Aware Collective Scheduling Policy for Distributed Training of DL Models

Saeed Rashidi[1], William Won[1], Sudarshan Srinivasan[2], Srinivas Sridharan[3], and Tushar Krishna[1]

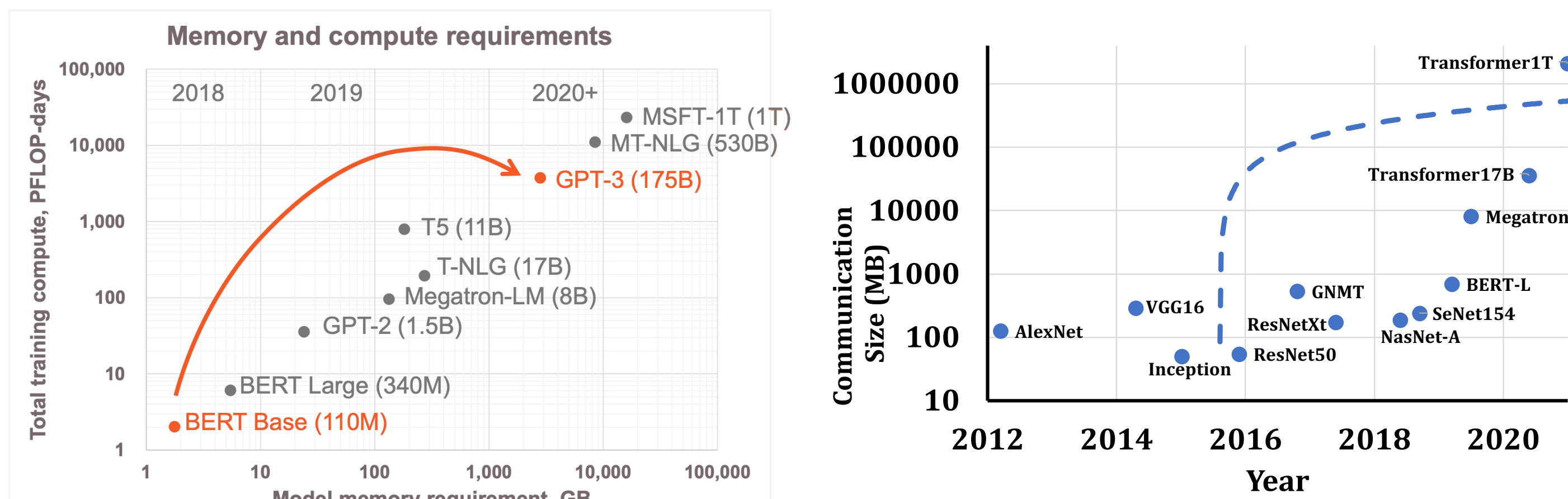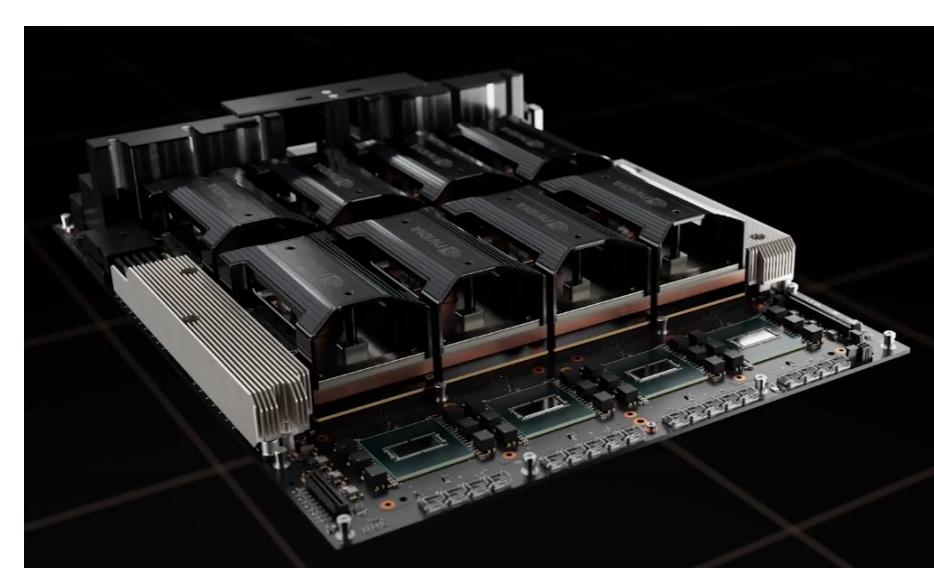[1]Georgia Institute of Technology, [2]Intel, [3]Meta

## Large AI Models

- Deep Learning: **Model size is increasing** (2× / 3.4 months)
  - GPT-3: 355 GPU years to train
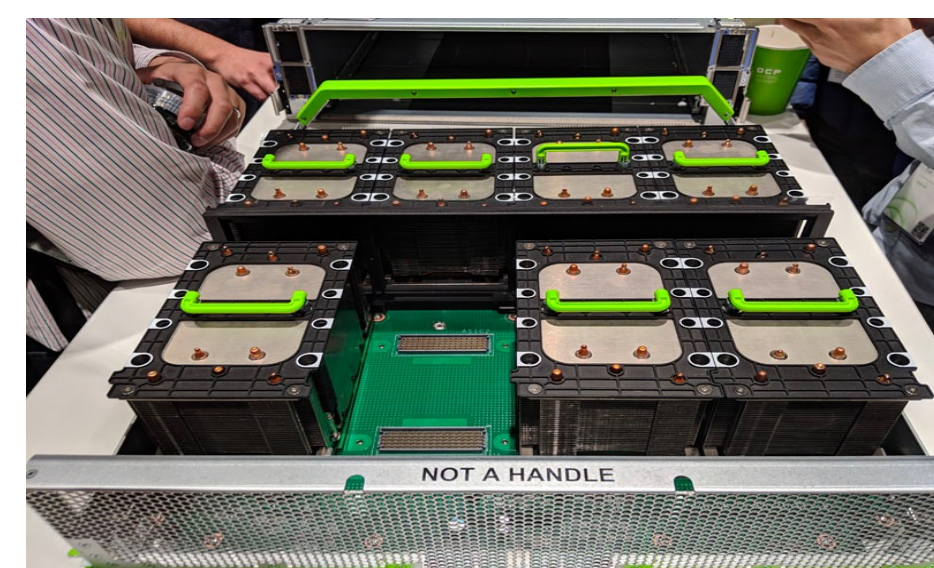- Necessitates **distributed training platforms**



## DL Training Platforms

- Futuristic training networks: **multi-dimension + heterogeneous BW**
  - Ring, Switch (SW), FullyConnected (FC)
  - NVIDIA HGX-H100: SW_SW_SW
  - Facebook Zion: FC_SW
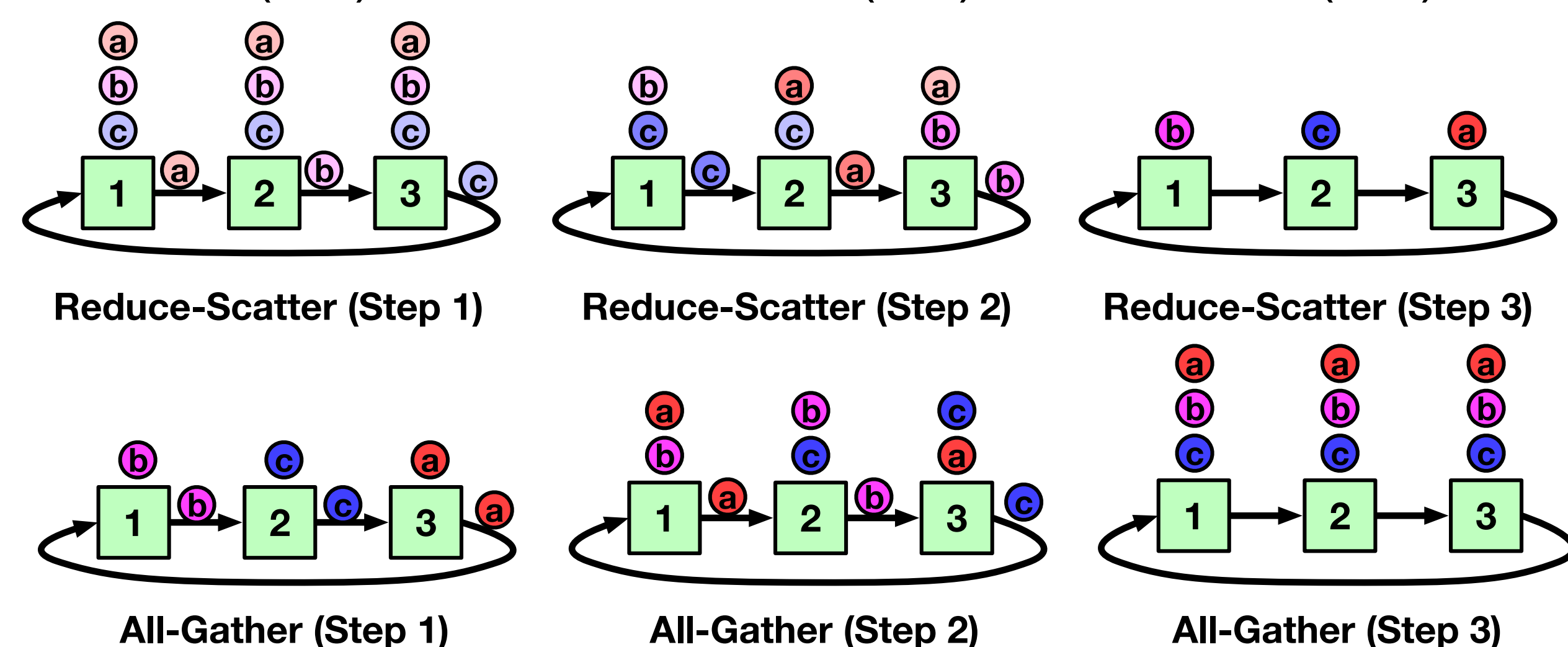  - Google TPUv4: Ring_Ring_Ring
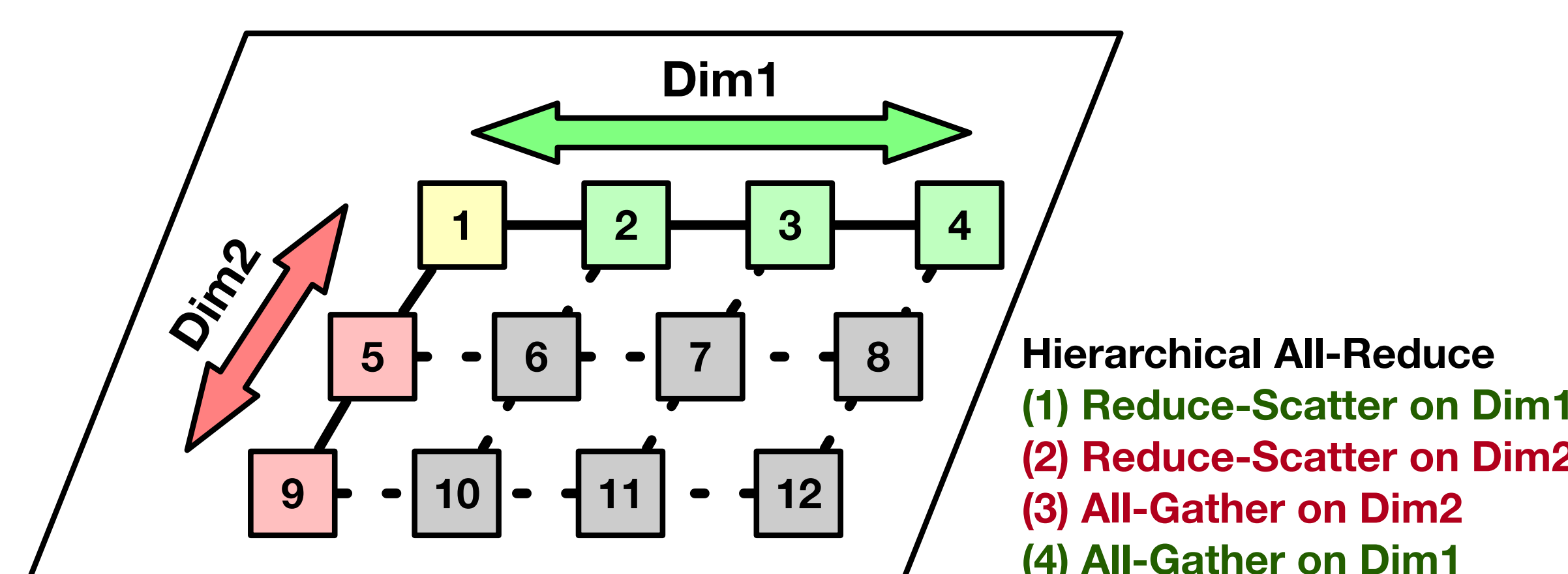


NVIDIA HGX-H100    Facebook Zion    Google TPUv4

## All-Reduce Algorithm

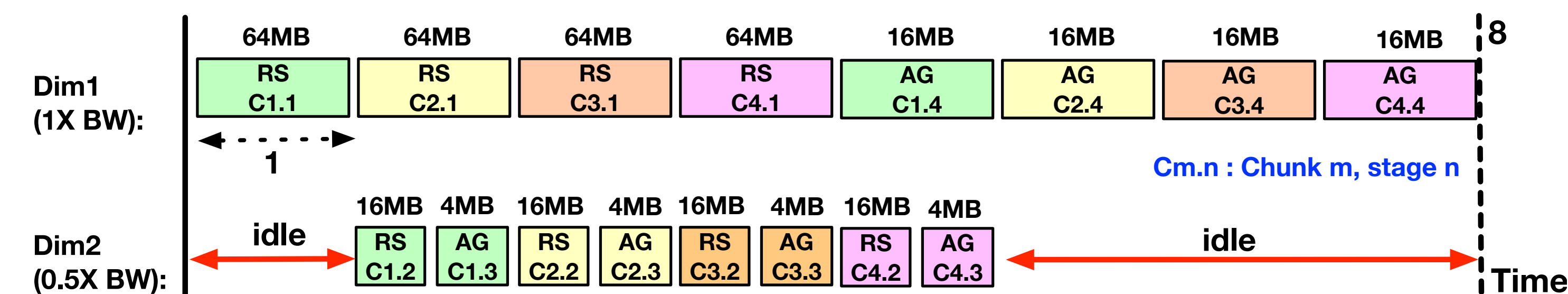- All-Reduce (AR) = Reduce-Scatter (RS) + All-Gather (AG)



Reduce-Scatter (Step 1)   Reduce-Scatter (Step 2)   Reduce-Scatter (Step 3)

All-Gather (Step 1)   All-Gather (Step 2)   All-Gather (Step 3)

- **Hierarchical** All-Reduce: Traverses dimensions **in-order**



Hierarchical All-Reduce
(1) Reduce-Scatter on Dim1
(2) Reduce-Scatter on Dim2
(3) All-Gather on Dim2
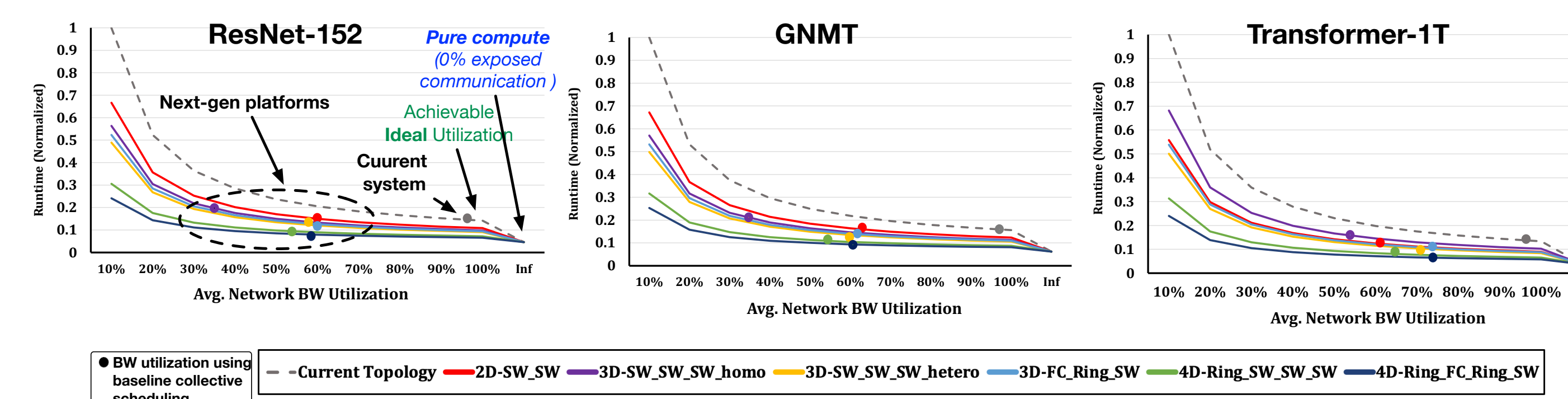(4) All-Gather on Dim1

## Motivation

- **Message size decreases** as traversing dimensions (Hierarchical AR)
- Each network dimension has **different BW** (Networking Technology)
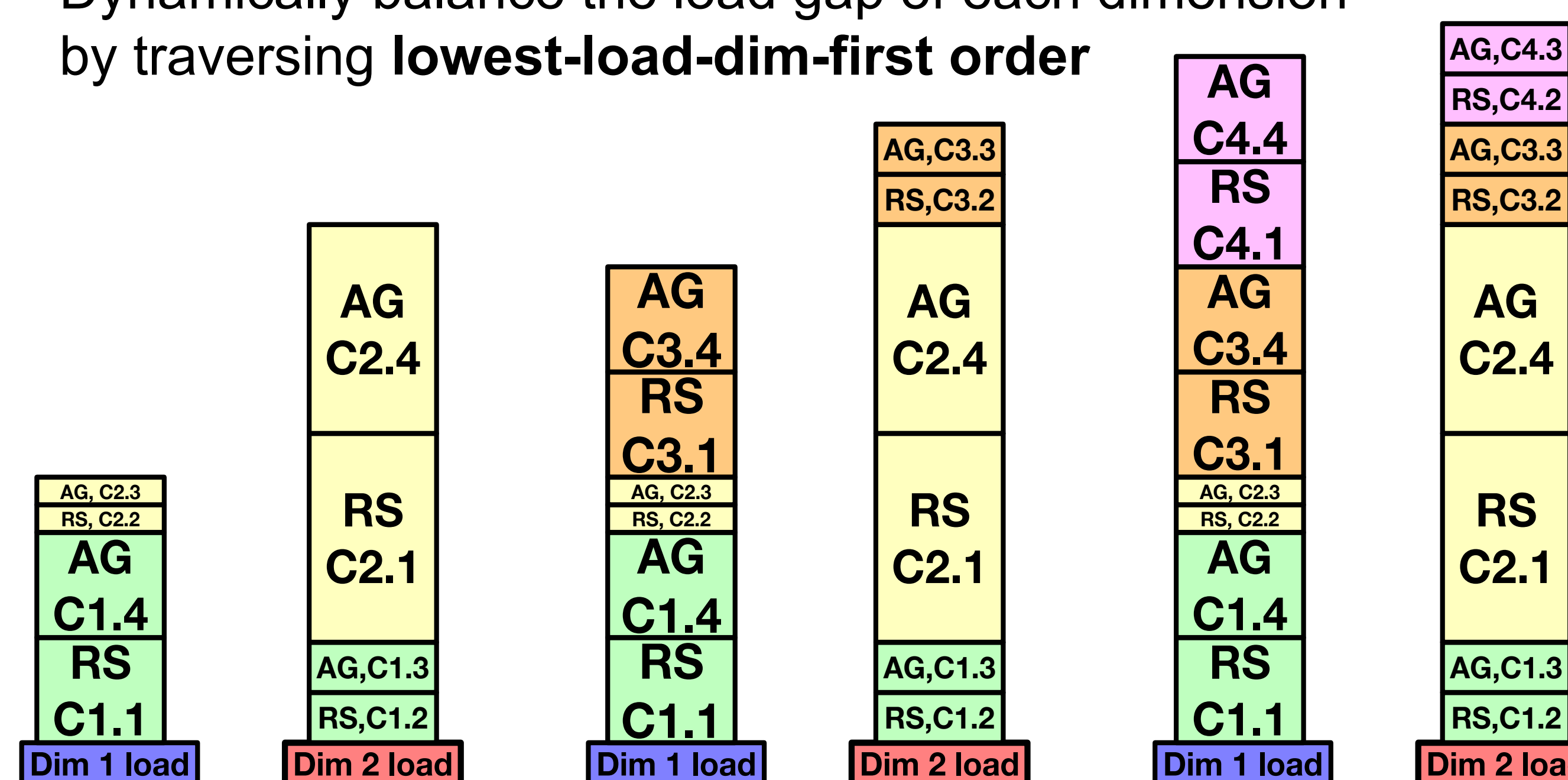  - Network BW and chunk size **mismatch across dimensions**



Cm.n : Chunk m, stage n

- **Underutilization** of Network BW Resource
  - Baseline: ~59.7% network BW utilization for next-gen topologies
  - ~1.37× (2.34× max) speedup if 100% BW utilization is achieved
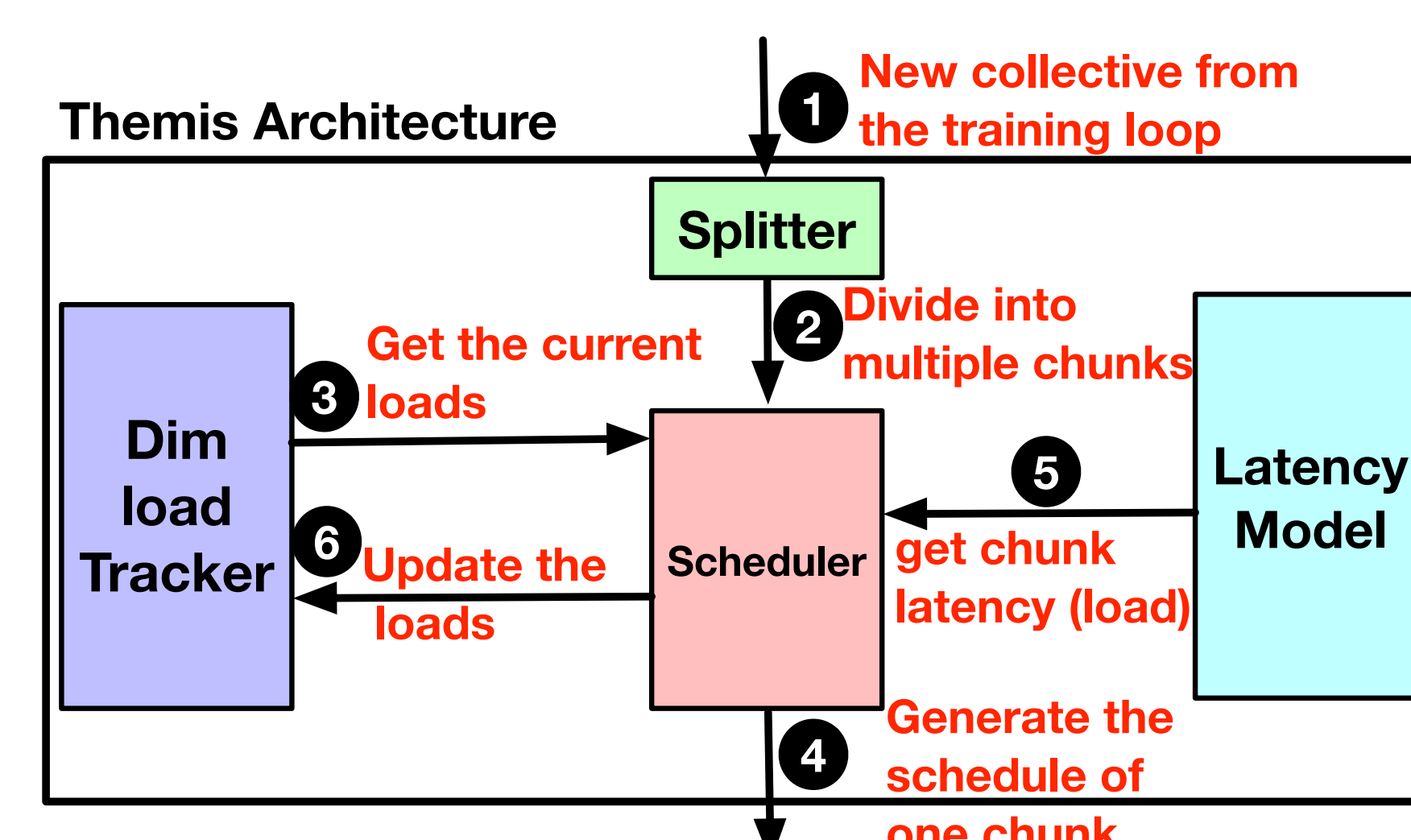


## Themis

- Observations:
  - Each chunk can traverse network dimensions **out-of-order**
  - Individual chunks can have **separate scheduling policy**
- **Themis**
  - **Track the current load** of each dimension
  - Dynamically balance the load gap of each dimension by traversing **lowest-load-dim-first order**



- Themis Components: **Scheduler, Latency Model, Dim Load Tracker**
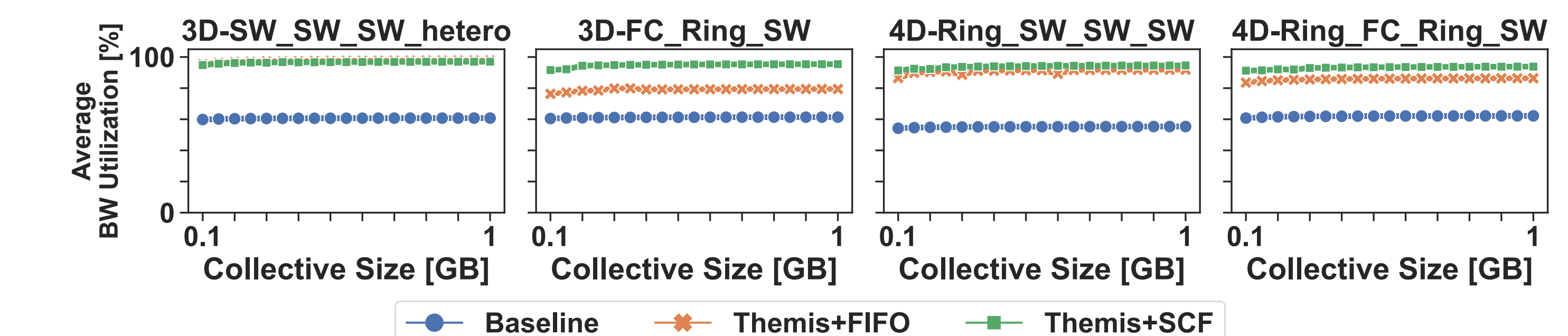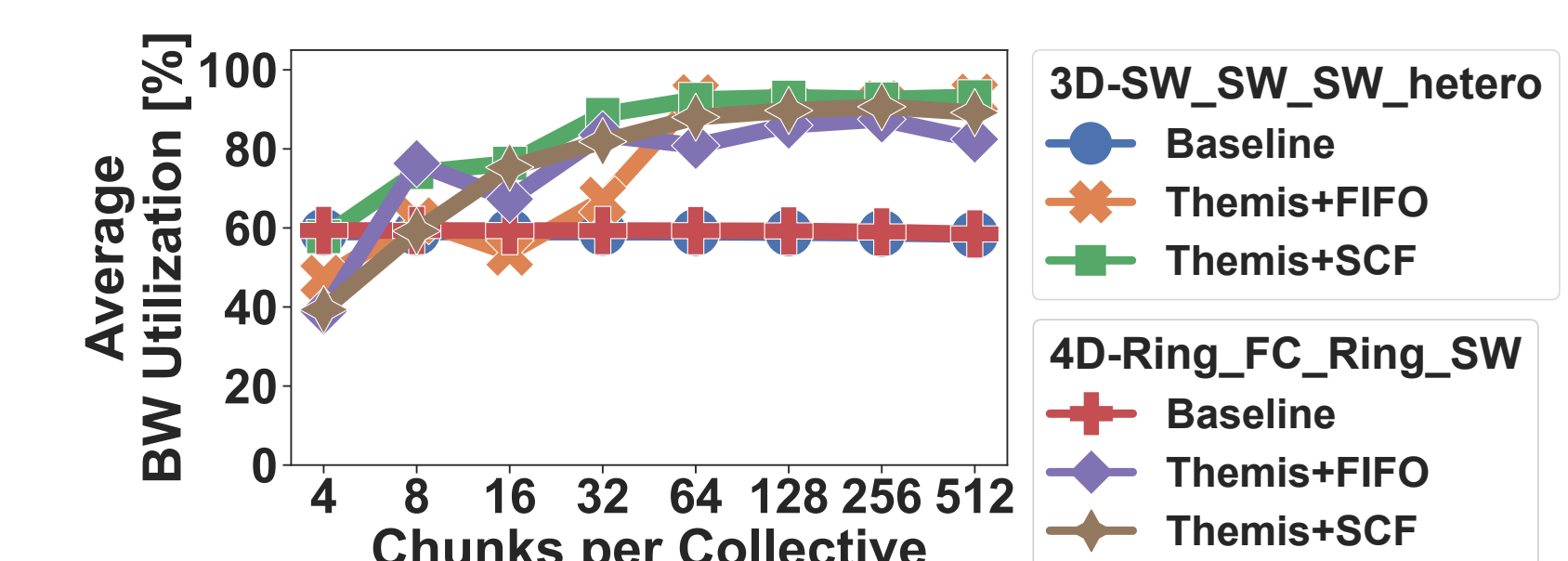


Themis Architecture
① New collective from the training loop
② Divide into multiple chunks
③ Get the current loads
④ Generate the schedule of one chunk
⑤ get chunk latency (load)
⑥ Update the loads

## Results

- Simulation Infrastructure: **ASTRA-sim**
  - https://astra-sim.github.io

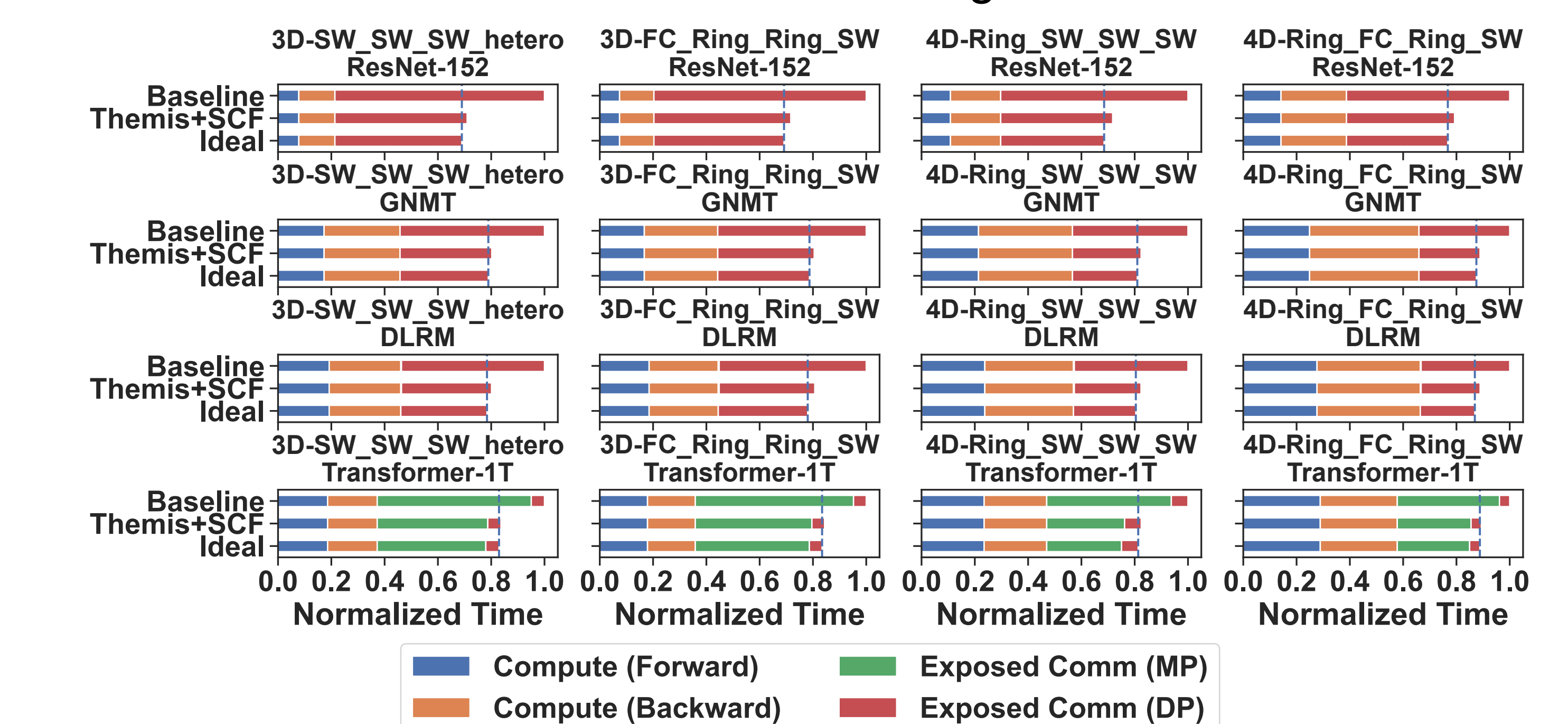| Name | Shape | BW (GB/s / Dim) | Latency (ns / Dim) |
|---|---|---|---|
| 3D-SW_SW_SW_hetero | 16 × 8 × 8 | 200, 100, 50 | 700, 1700 |
| 3D-FC_Ring_SW | 8 × 16 × 8 | 175, 100, 50 | 700, 700, 1700 |
| 4D-Ring_SW_SW_SW | 4 × 4 × 8 × 8 | 250, 200, 100, 50 | 20, 700, 700, 1700 |
| 4D-Ring_FC_Ring_SW | 4 × 8 × 4 × 8 | 375, 175, 150, 100 | 20, 700, 700, 1700 |

- Single All-Reduce: Themis achieves ~**95.14% BW Utilization**
  - Baseline: ~56.31% (1.72× speedup)



- **More #chunks = better load balancing** capabilities



- Workloads: Themis reaches **near-ideal training performance**
  - ~1.49× (ResNet), ~1.30× (GNMT/DLRM), ~1.25× (T-1T) speedup over baseline hierarchical collective algorithm



Compute (Forward)   Compute (Backward)   Exposed Comm (MP)   Exposed Comm (DP)

## Conclusion

- Understanding futuristic training platforms
  - **Multi-dimensional** network with **heterogeneous BW**
- **Huge network BW underutilization** is observed
  - Due to chunk size and network BW mismatch across dimensions
- Themis: **Dynamic chunk scheduler** to improve BW utilization
  - By monitoring and **balancing loads** of each dimension
  - **95.14% network BW** utilization (Single All-Reduce)
  - **1.49×** (ResNet-152), **1.25×** (Transformer-1T) speedup