

Swing: Short-cutting Rings for Higher Bandwidth Allreduce

Daniele De Sensi¹, Tommaso Bonato², David Saam³, Torsten Hoefler²

2

ETH zürich

1



3

RWTHAACHEN
UNIVERSITY

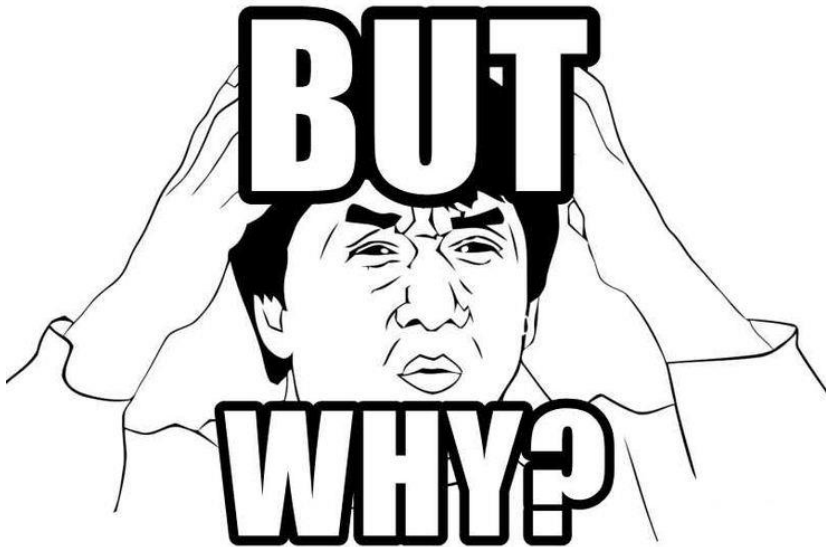
NSDI

Santa Clara – April 16-18, 2024



Swing

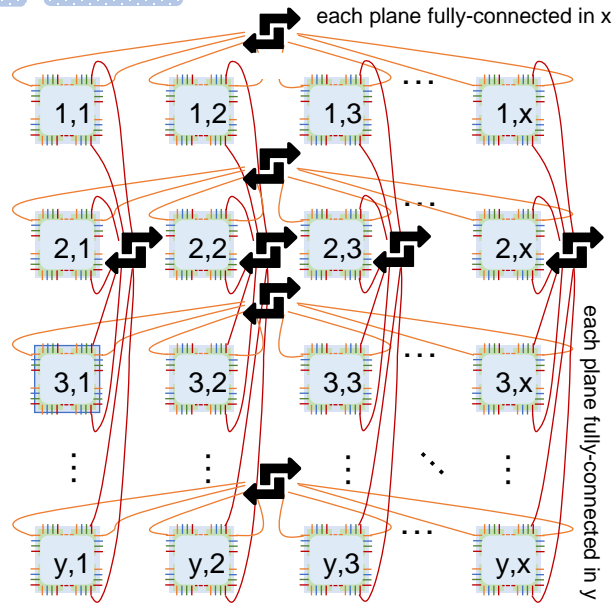
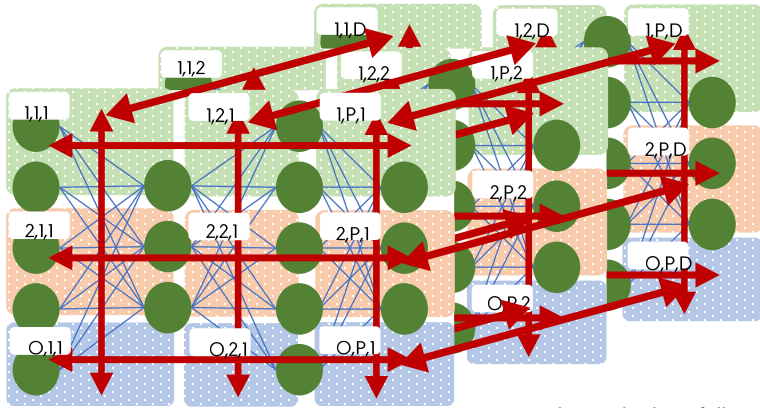
A novel **allreduce/allgather/reduce-scatter** algorithm
optimized for multi-dimensional **torus networks**



(expected advantages on any
blocking network)

Why torus and why allreduce?

3D - Data, Pipeline, and Operator Parallelism



AI & Machine Learning

Enabling next-generation AI workloads: Announcing TPU v5p and AI Hypercomputer

December 7, 2023

Google's TPU v5p Pod
(> 9,000 chips on a 3D torus)

Amazon EC2 Trn1 Instances

High-performance, cost-effective training of generative AI models

Get started with Trn1 instances using AWS Neuron

AWS Trainium Instances
(16 chips on a 2D torus)

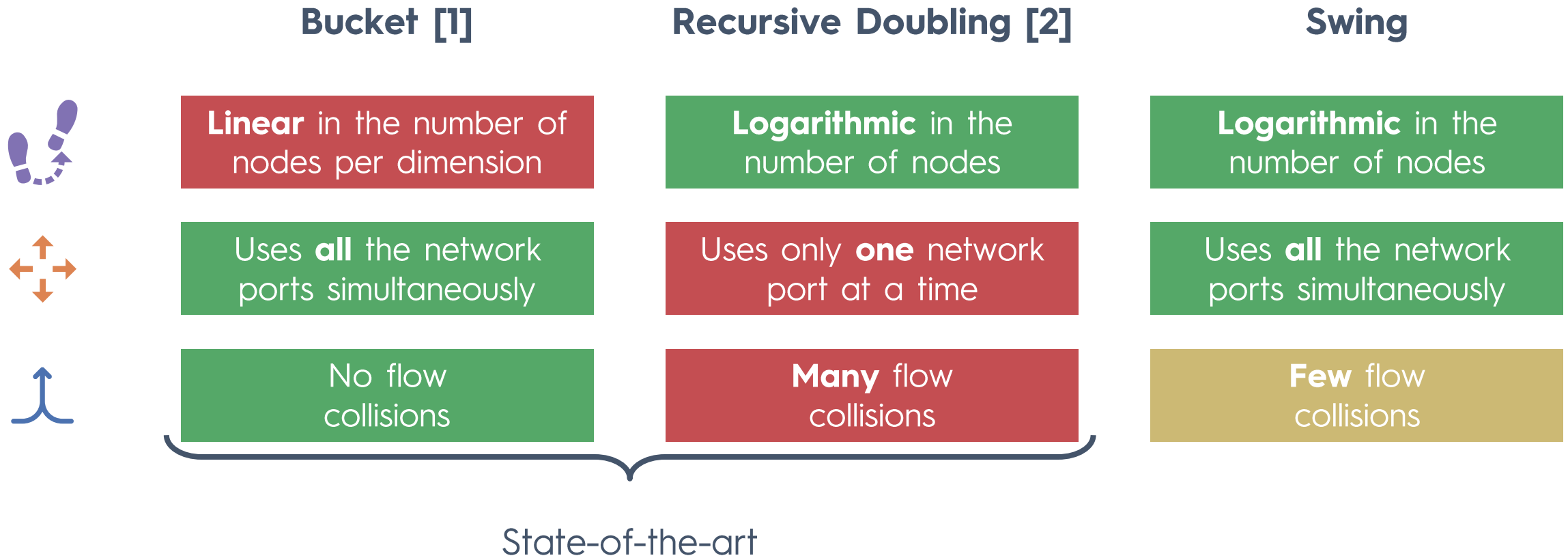
GRAPHCORE

BUILD: IPU-POD₆₄

Graphcore IPU-POD
(64 chips on a 2D torus)

HammingMesh: A Network Topology for Large-Scale Deep Learning (2022)
Torsten Hoefler, Tommaso Bonato, Daniele De Sensi, Salvatore Di Girolamo, Shigang Li, Marco Heddes, Jon Belk, Deepak Goel, Miguel Castro, Steve Scott

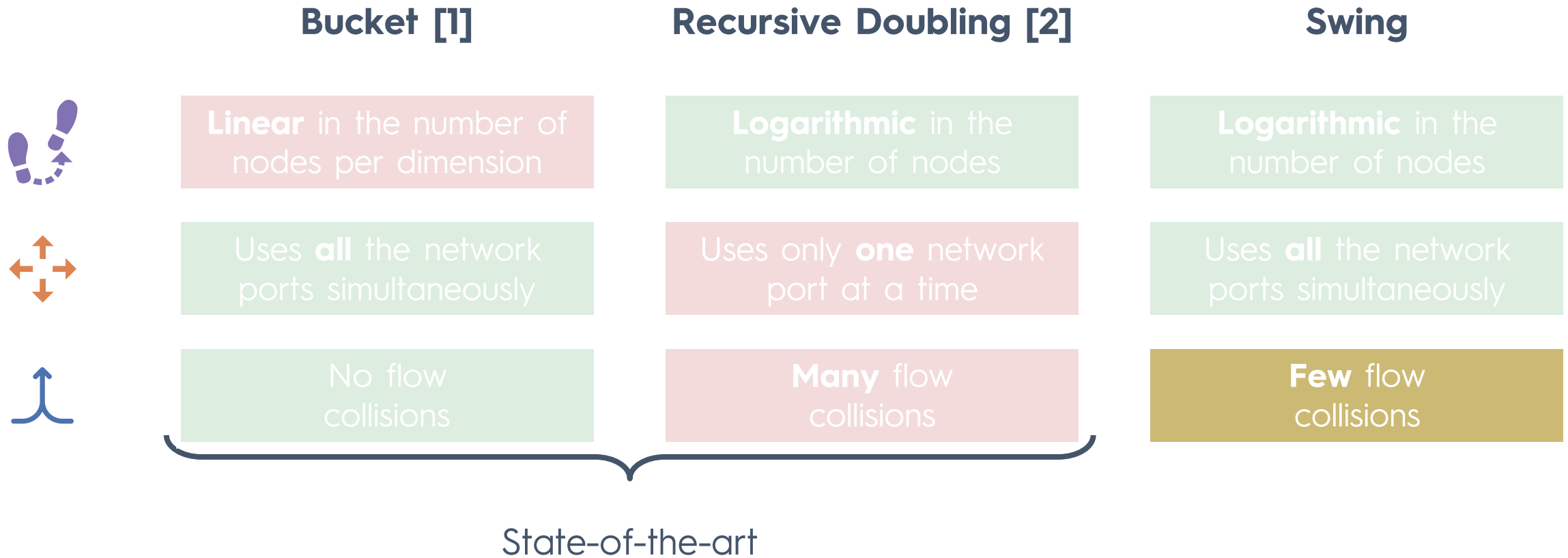
Allreduce Algorithms



[1] N. Jain and Y. Sabharwal. *Optimal bucket algorithms for large MPI collectives on torus interconnects* (2010)

[2] P. Sack and W. Gropp. *Collective algorithms for multiported torus networks* (2015)

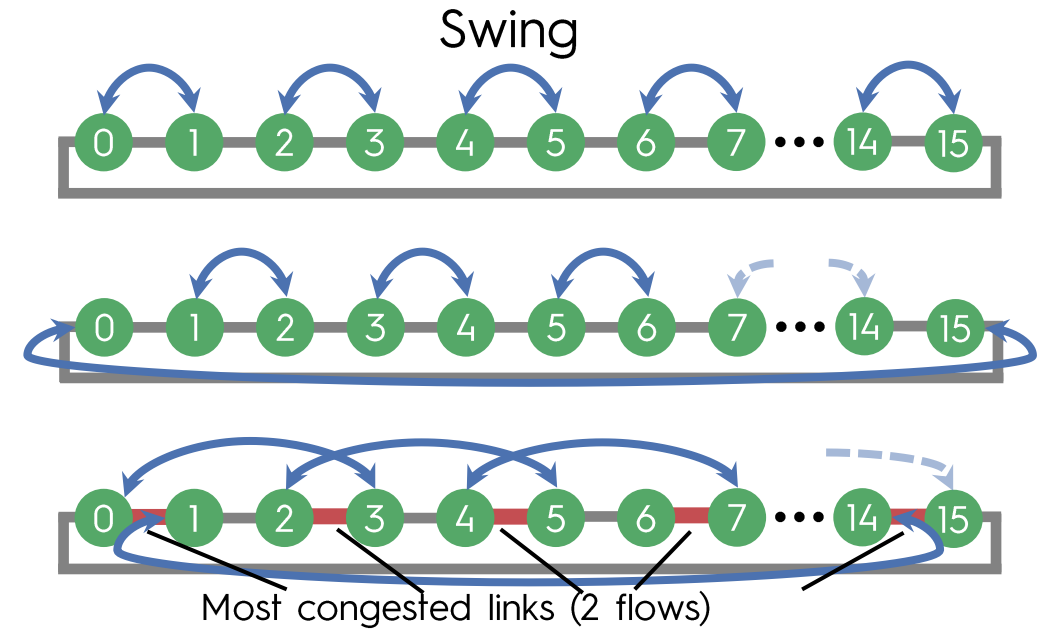
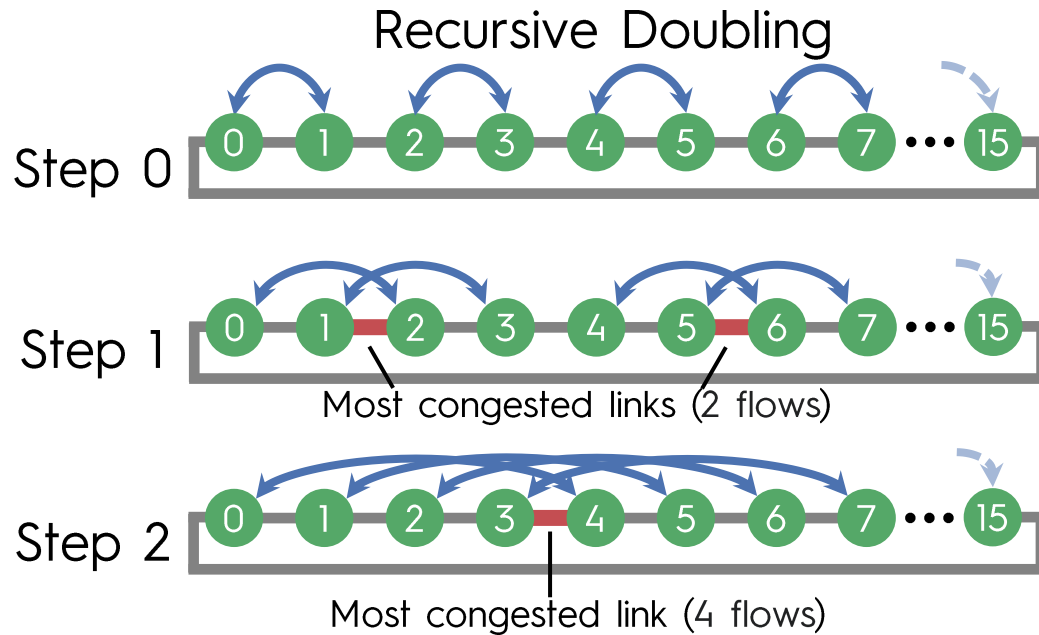
Allreduce Algorithms



[1] N. Jain and Y. Sabharwal. *Optimal bucket algorithms for large MPI collectives on torus interconnects* (2010)

[2] P. Sack and W. Gropp. *Collective algorithms for multiported torus networks* (2015)

Swing Allreduce



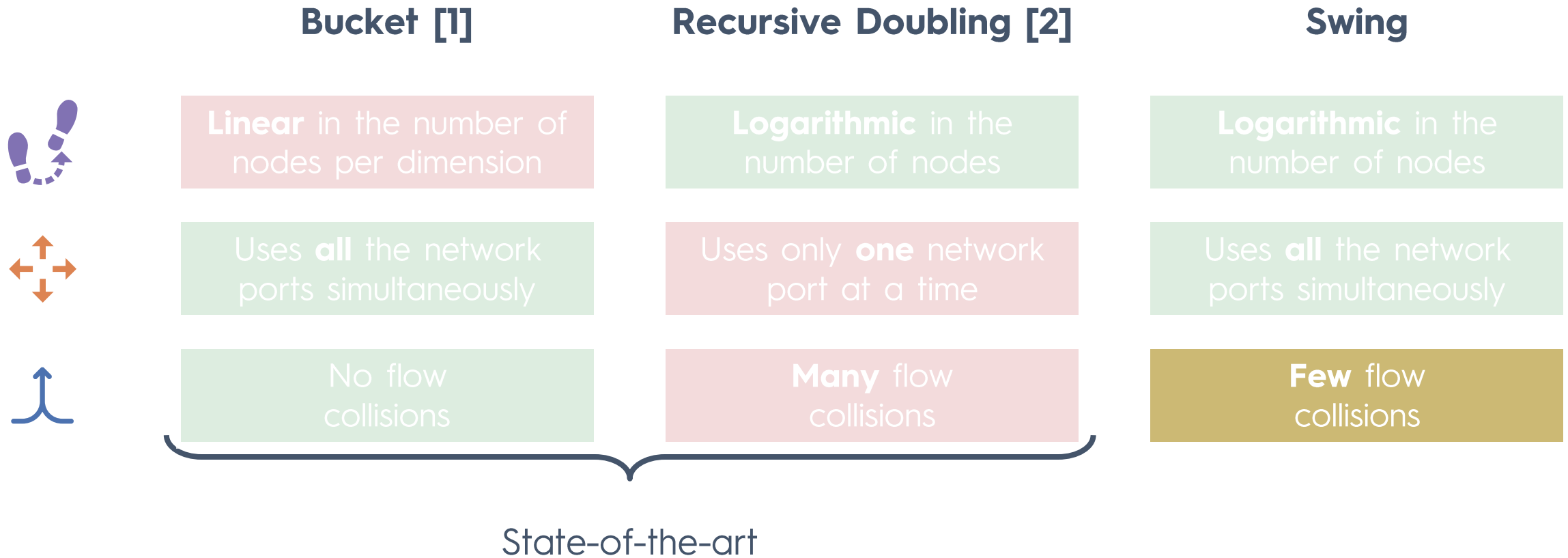
Swing Allreduce

Details, generalization to D dimensions, and correctness proof in the paper

Recursive Doubling



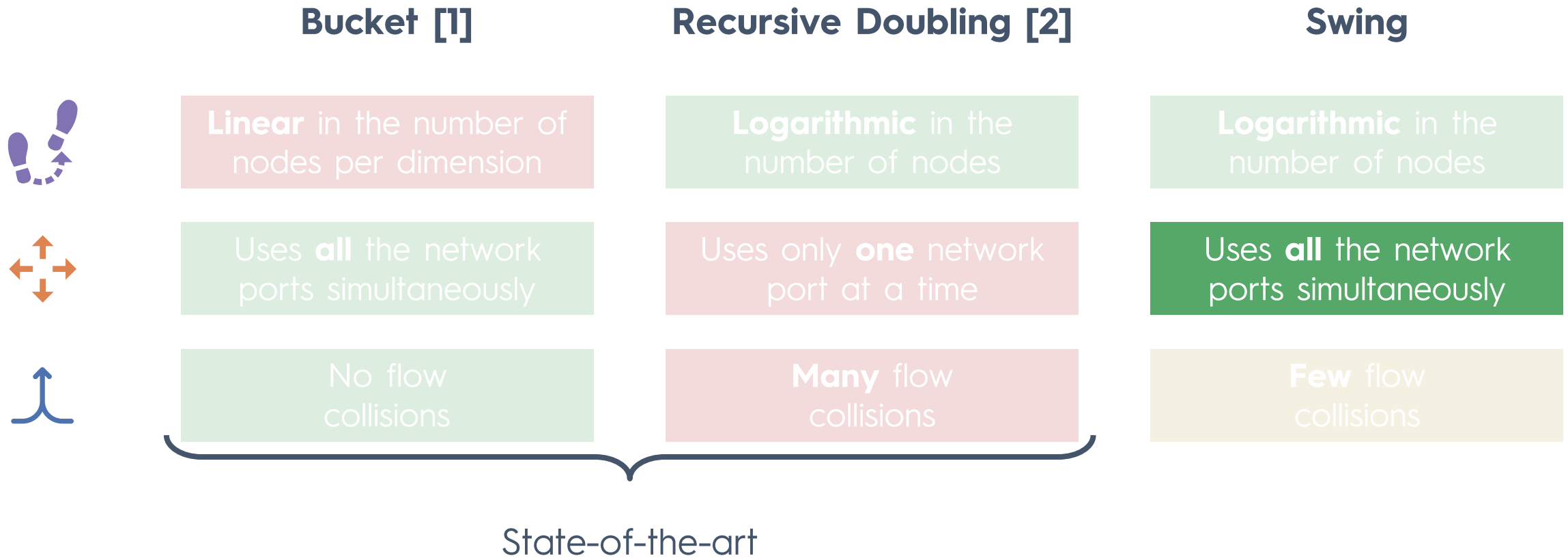
Allreduce Algorithms



[1] N. Jain and Y. Sabharwal. *Optimal bucket algorithms for large MPI collectives on torus interconnects* (2010)

[2] P. Sack and W. Gropp. *Collective algorithms for multiported torus networks* (2015)

Allreduce Algorithms

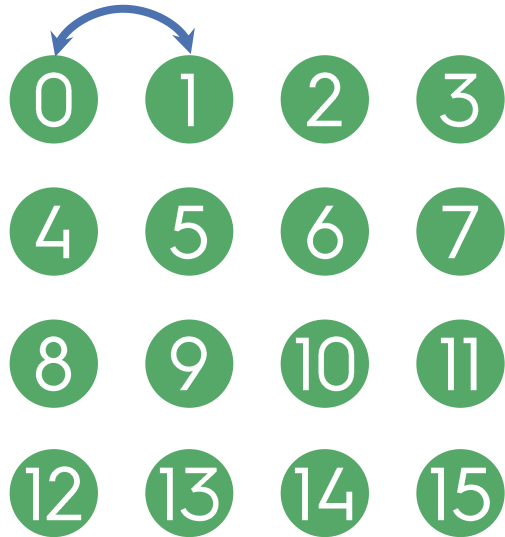


[1] N. Jain and Y. Sabharwal. *Optimal bucket algorithms for large MPI collectives on torus interconnects* (2010)

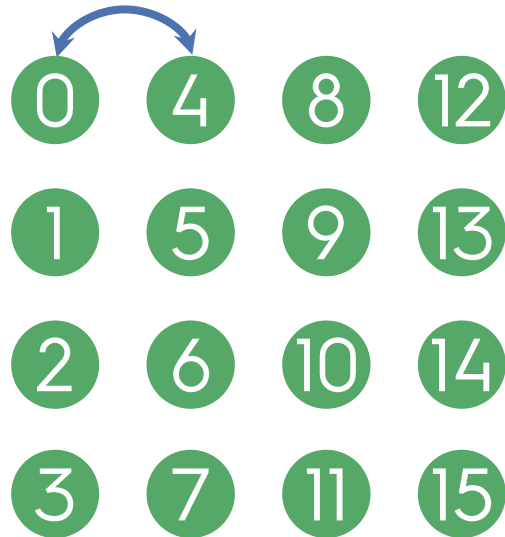
[2] P. Sack and W. Gropp. *Collective algorithms for multiported torus networks* (2015)

Multiport Swing

Step 0, port 0



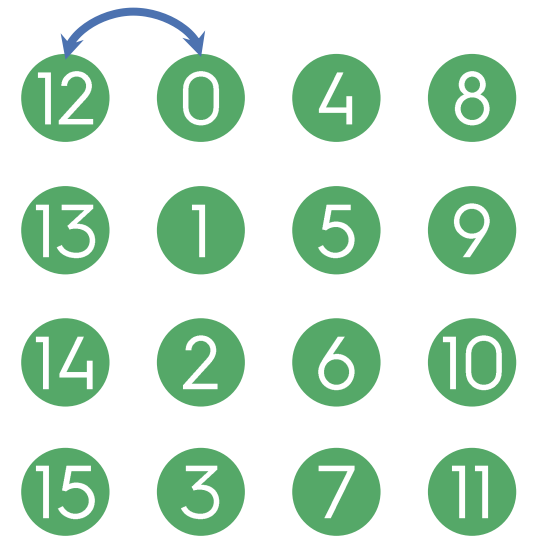
Step 0, port 1



Step 0, port 2



Step 0, port 3



Transpose

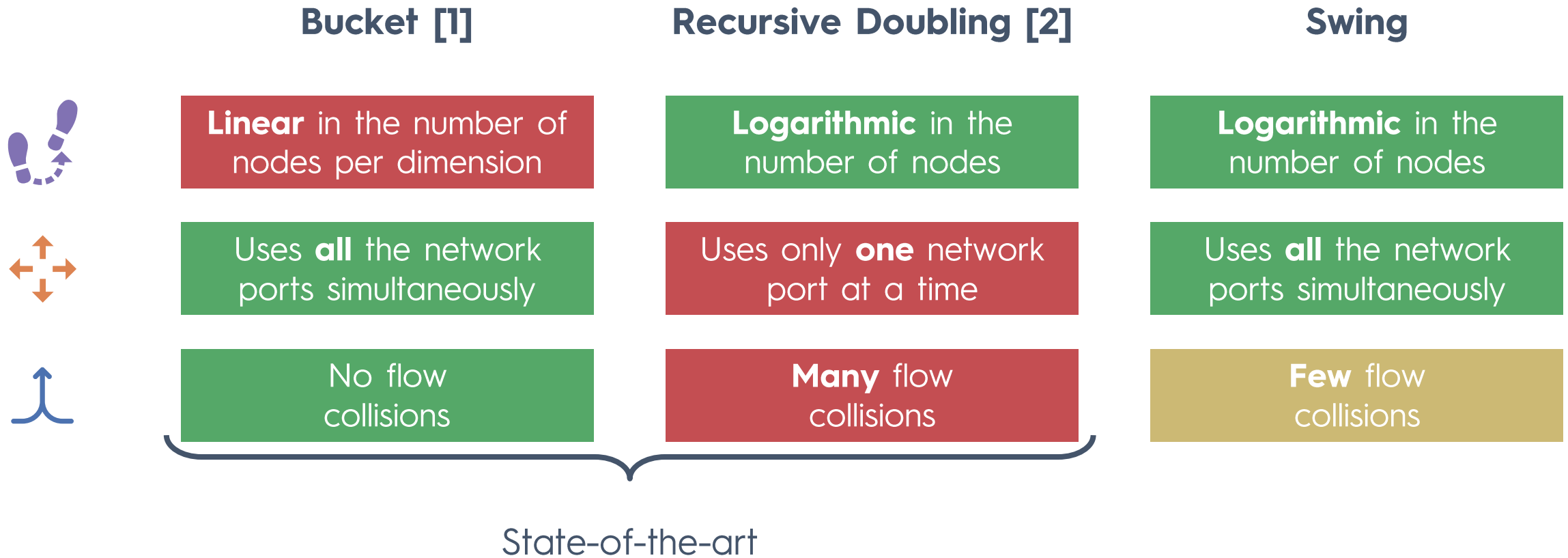


Shift columns by one



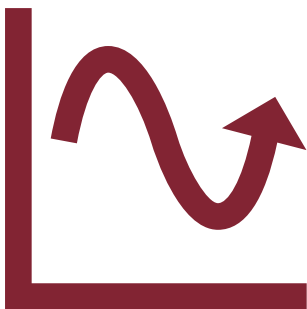
Transpose
+
shift columns by one

Allreduce Algorithms



[1] N. Jain and Y. Sabharwal. *Optimal bucket algorithms for large MPI collectives on torus interconnects* (2010)

[2] P. Sack and W. Gropp. *Collective algorithms for multiported torus networks* (2015)



Evaluation

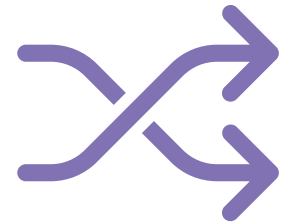
Setup



SST packet-level
network simulator

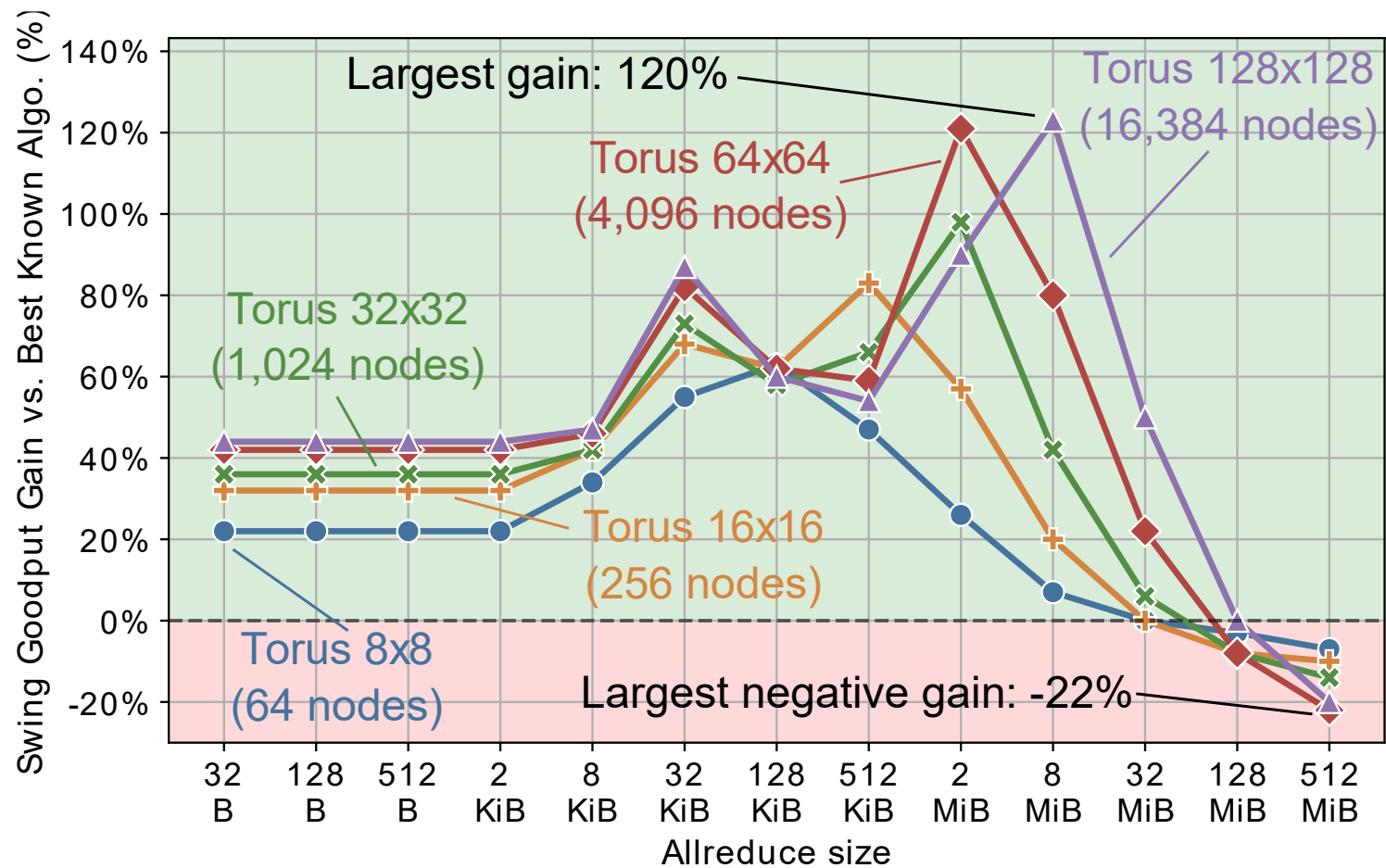


400 Gb/s links
100 ns latency

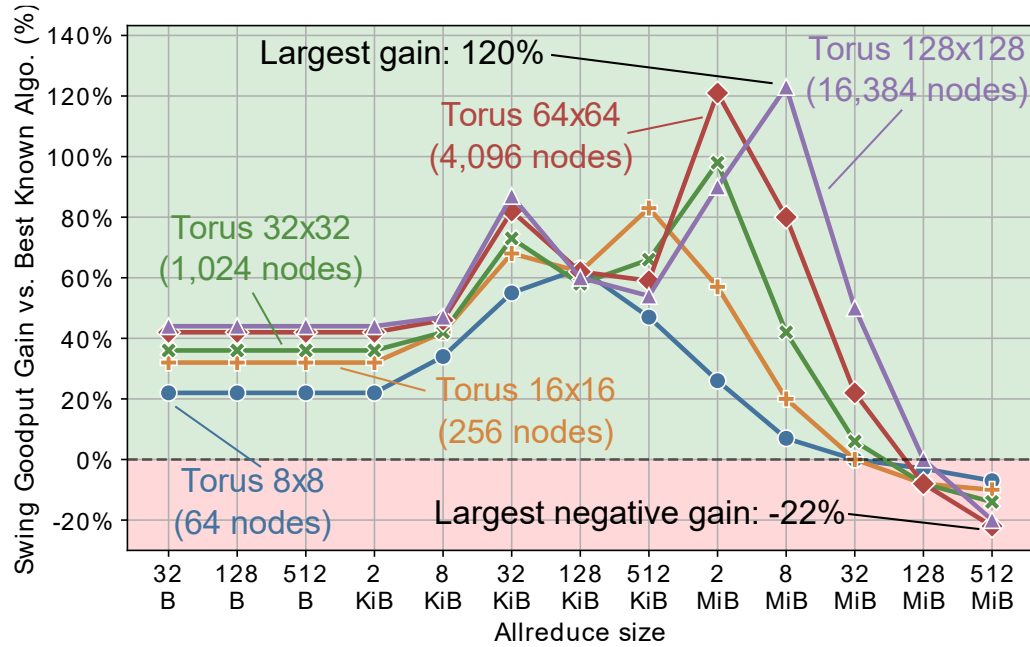


300 ns
per-hop latency

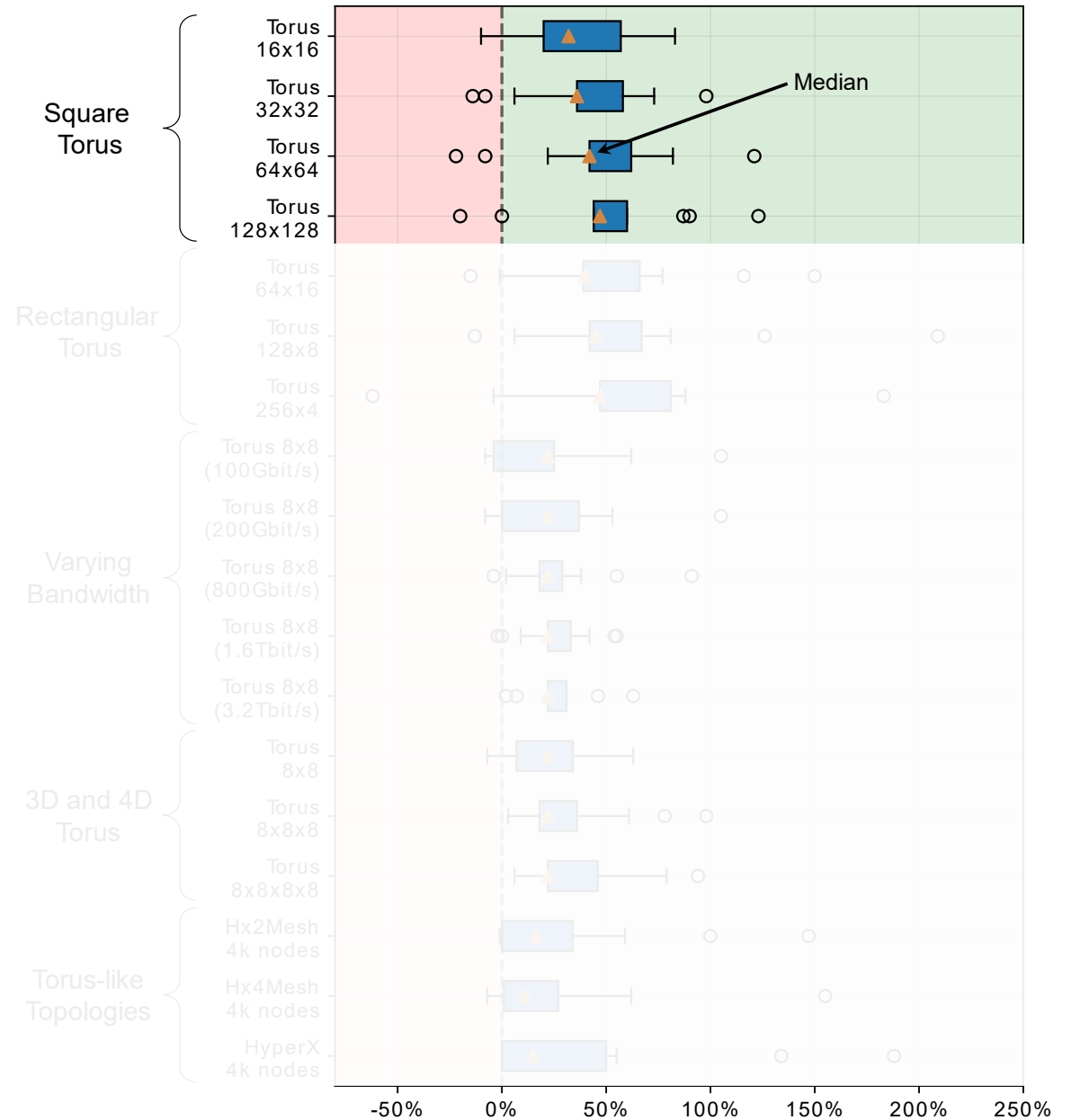
Results - Performance Gain



Results Summary



Goodput Gain vs. Best Known Algo. for Allreduce <= 512MiB

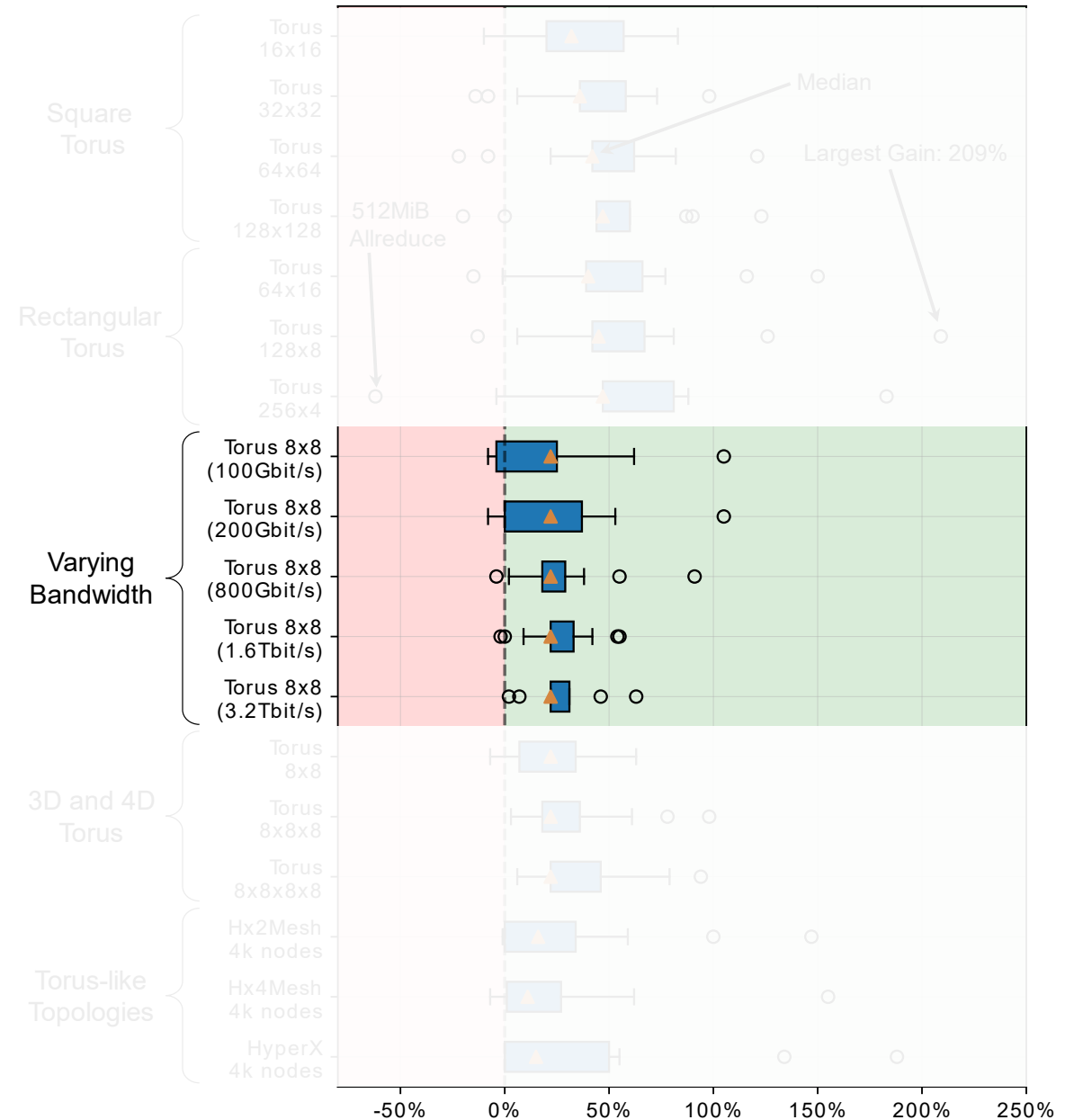


Results Summary

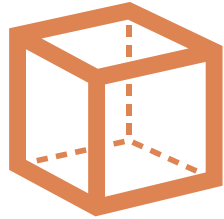


At higher bandwidth, the number of steps has a higher relative impact on performance

Goodput Gain vs. Best Known Algo. for Allreduce $\leq 512\text{MiB}$

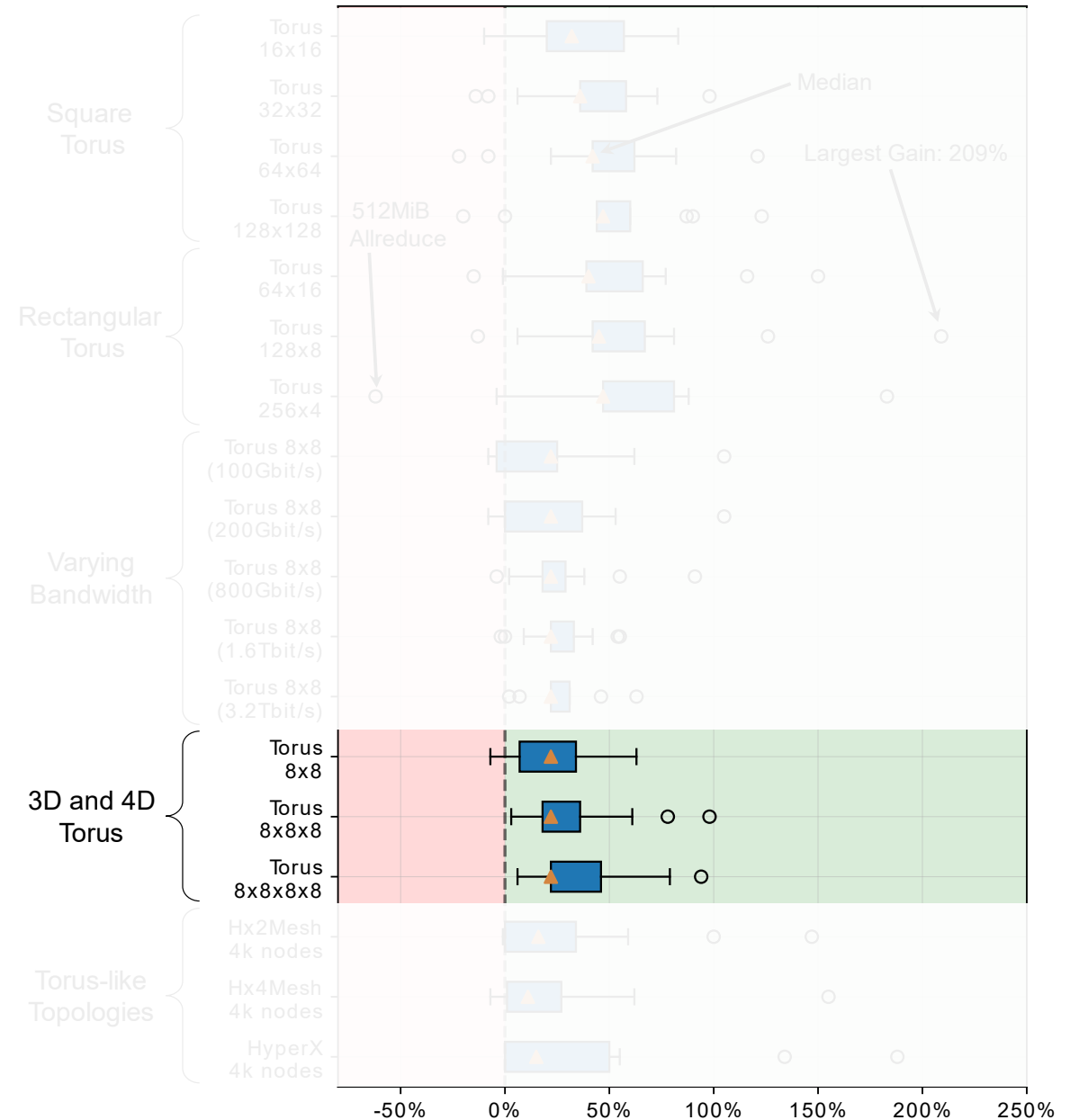


Results Summary

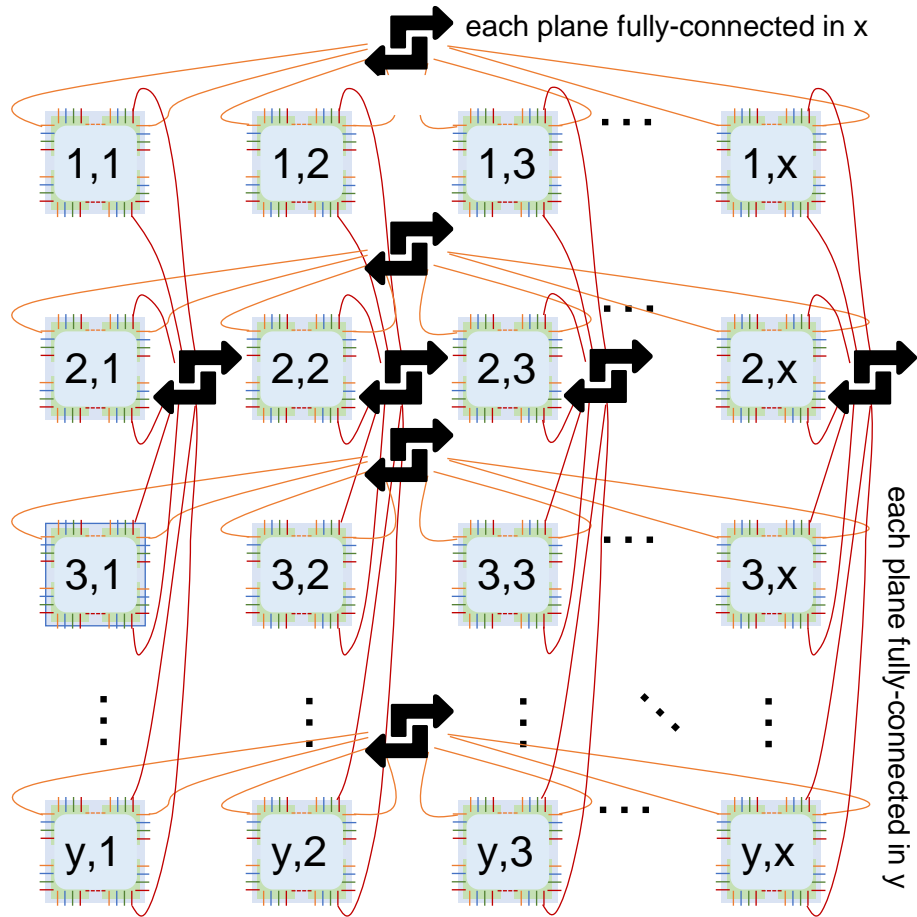


More dimensions imply more communications with close nodes

Goodput Gain vs. Best Known Algo. for Allreduce $\leq 512\text{MiB}$

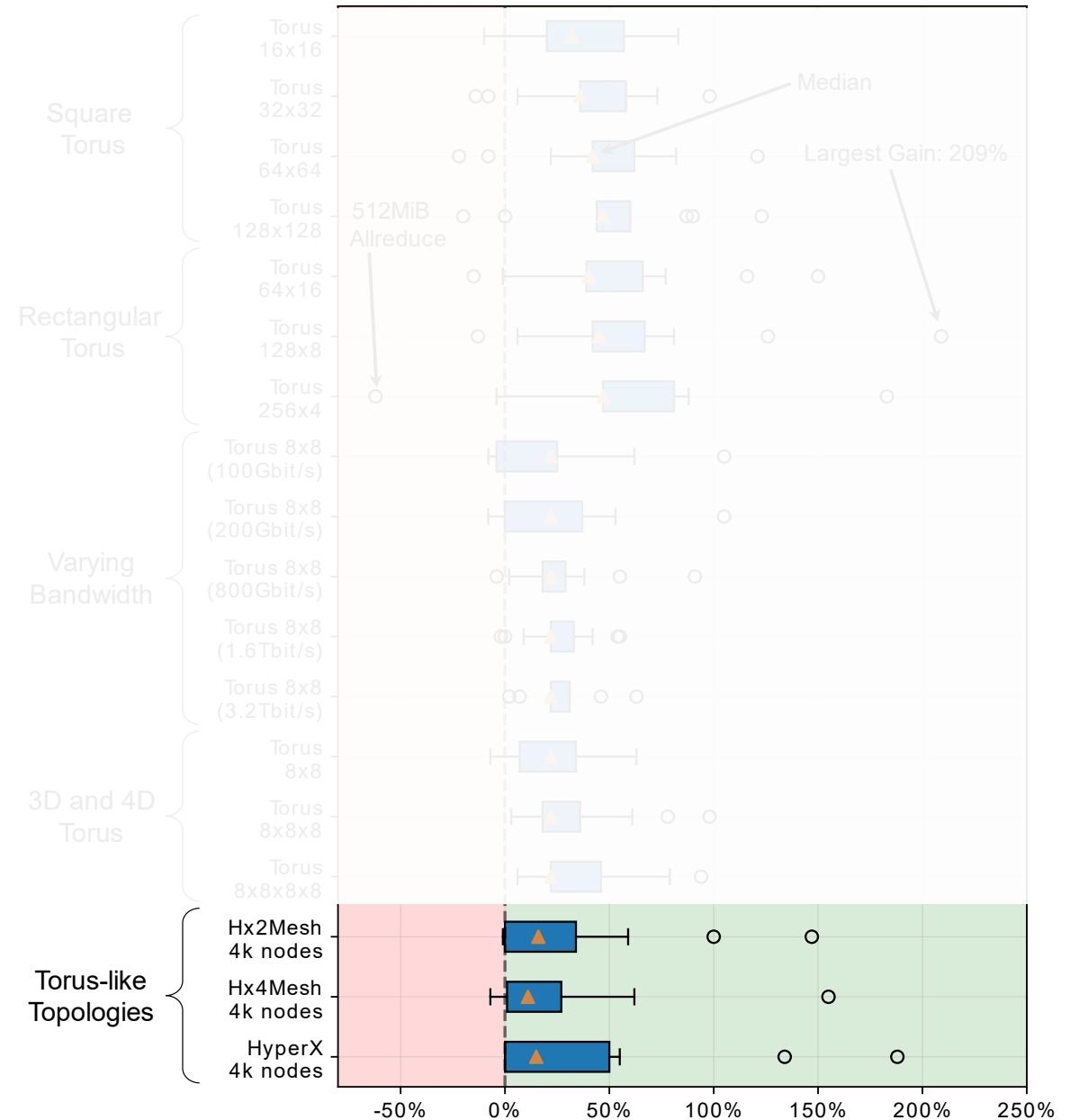


Results Summary

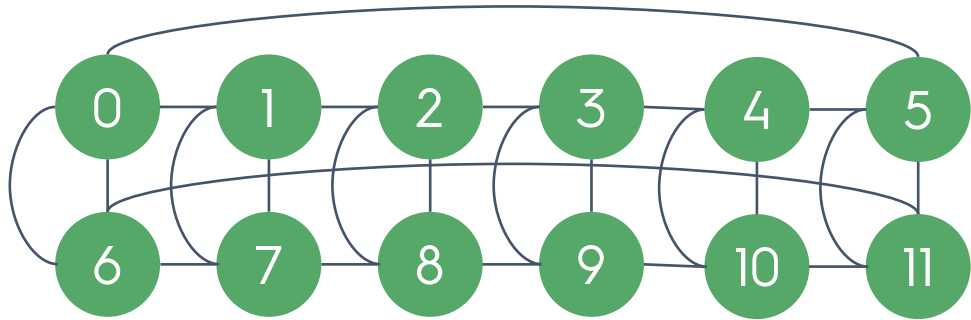


HammingMesh: A Network Topology for Large-Scale Deep Learning (2022)
 Torsten Hoefler, Tommaso Bonato, Daniele De Sensi, Salvatore Di Girolamo, Shigang Li, Marco Heddes, Jon Belk, Deepak Goel, Miguel Castro, Steve Scott

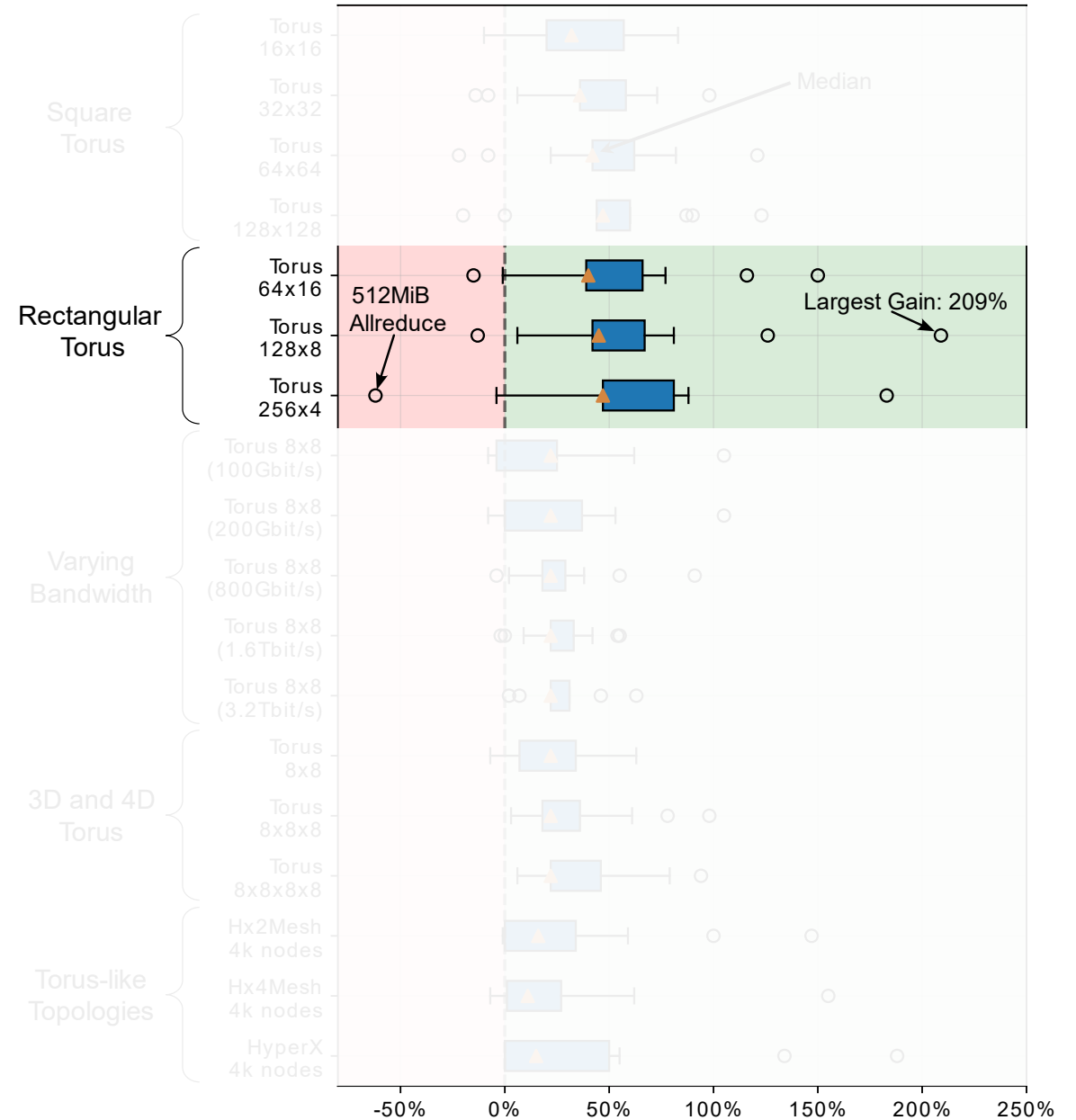
Goodput Gain vs. Best Known Algo. for Allreduce $\leq 512\text{MiB}$



Results Summary



Goodput Gain vs. Best Known Algo. for Allreduce <= 512MiB



Results Summary

5.3 Performance for 3D and 4D Torus

As discussed in Sec. 4 and summarized in Table 2, the performance of the allreduce algorithm for multidimensional torus also depends on the number of dimensions. Thus, we evaluate the performance of the different allreduce algorithms on 8^2 , 8^3 , and 8^4 torus networks.

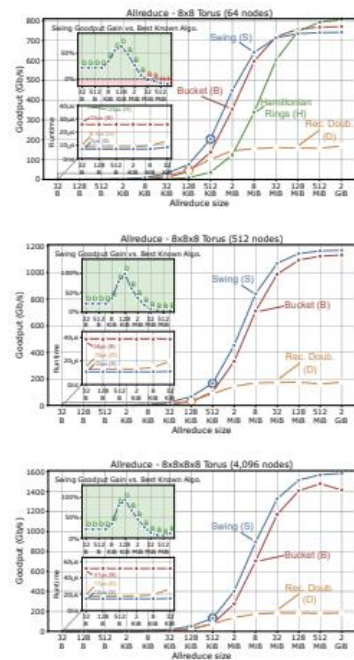


Figure 11: Goodput on higher-dimensional torus networks: 2D 8x8, 3D (8x8x8), and 4D (8x8x8x8).

We report the evaluation result in Fig. 11. We do not include the Hamiltonian ring algorithm in the 3D and 4D torus results since it only works for 2D torus networks. When increasing the number of dimensions, the goodput gain of Swing in-

creases because, as shown in Table 2 and discussed in Sec. 4, the congestion deficiency drops to 3% on 3D torus and to 0.8% on 4D torus. Consequently, for 3D and 4D torus networks, Swing outperforms by up to 2x all existing algorithms on allreduce ranging from 32B to 2GB.

5.4 Performance on Torus-Like Topologies

Some topologies like HammingMesh [26] and HyperX [3,20] extend torus by adding additional links, thus increasing the network bisection bandwidth. Seen from a different perspective, those extra links allow distant nodes to communicate crossing fewer hops, decreasing Swing congestion deficiency.

5.4.1 Performance on HammingMesh

HammingMesh [26] groups nodes into square boards. Each board is a 2D mesh, and nodes on the same column (or row) located at the edge of the boards are connected together using fat trees. Due to its higher performance and flexibility compared to a torus a similar topology is used, for example, to interconnect TPUv4 devices [31]. Because of the extra links, the congestion deficiency of Swing on a HammingMesh is lower than that on a 2D torus. Moreover, for a fixed number of nodes, having smaller boards increases the number of extra (fat tree) links and, thus, decreases the congestion deficiency.

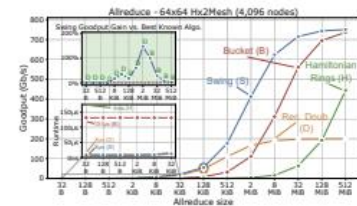
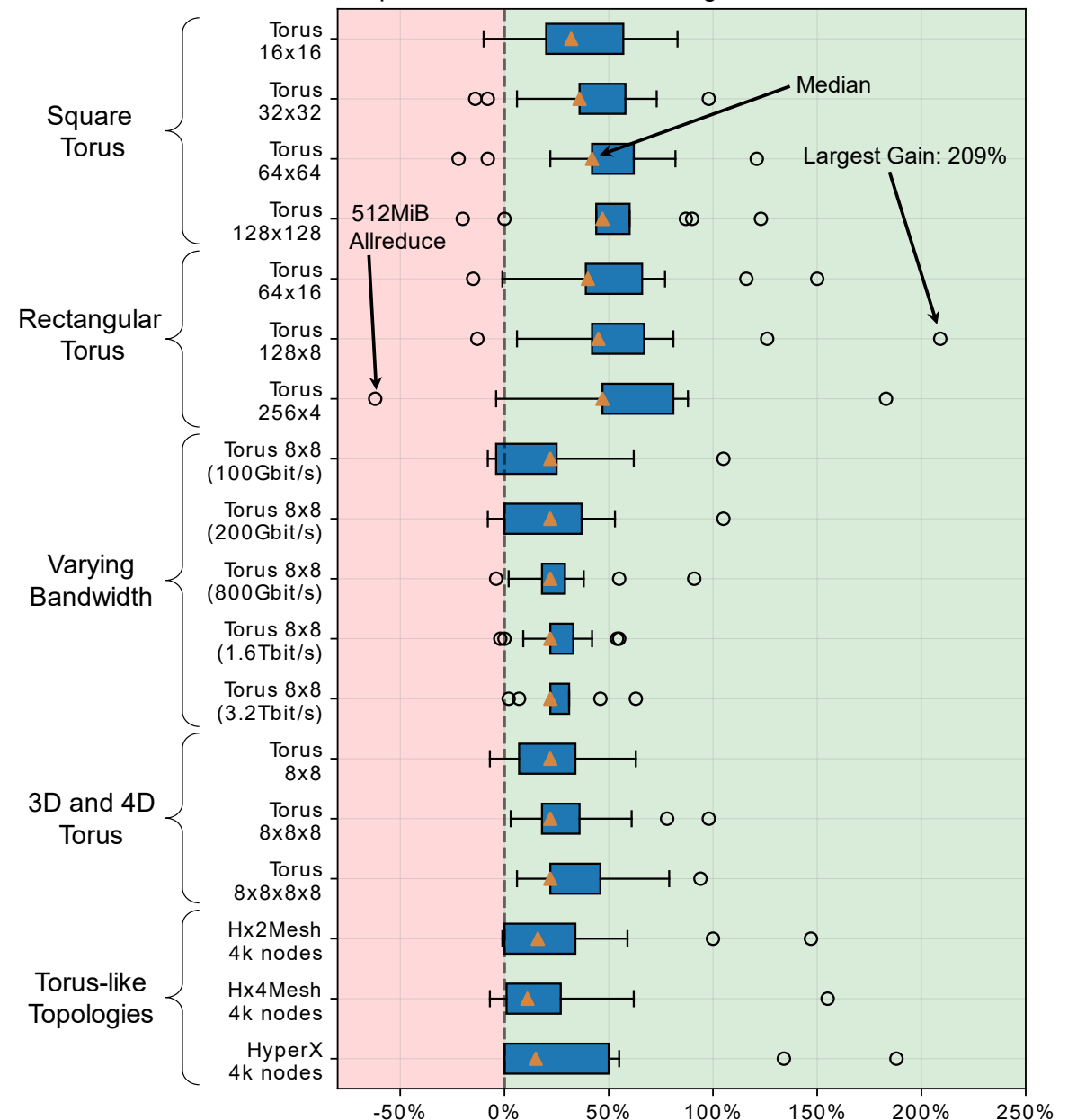


Figure 12: Goodput on a 4,096 nodes Hx2Mesh.

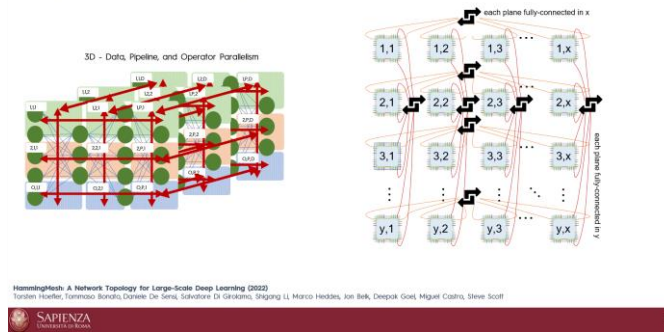
We show in Fig. 12 the performance of the different algorithms for a Hx2Mesh network with 4,096 nodes (2x2 boards arranged in a 32x32 configuration). For such configuration, Swing outperforms the state-of-the-art algorithms at any size, up to 2.5x for 2MiB allreduce. Moreover, because of the lower congestion deficiency, we observe how the peak Swing performance is higher compared to a 2D torus with the same number of nodes (Fig. 6). Last, we also observe a runtime reduction for all the algorithms for small vectors, since nodes on the same board on HammingMesh are connected through PCB traces, with lower latency than optical network cables.

Goodput Gain vs. Best Known Algo. for Allreduce <= 512MiB

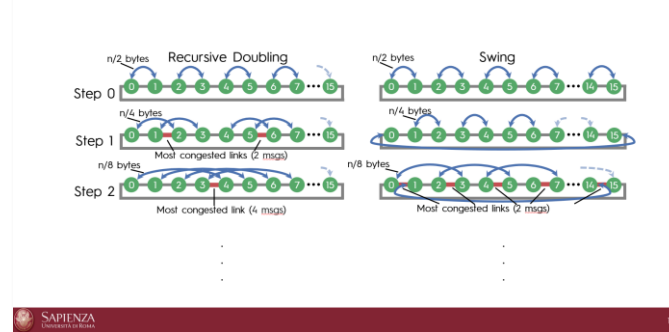


Conclusions

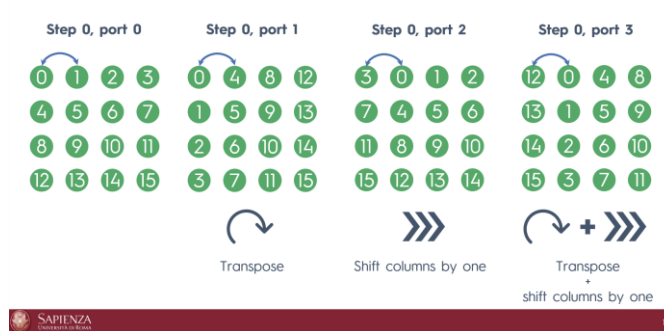
Why torus and why allreduce?



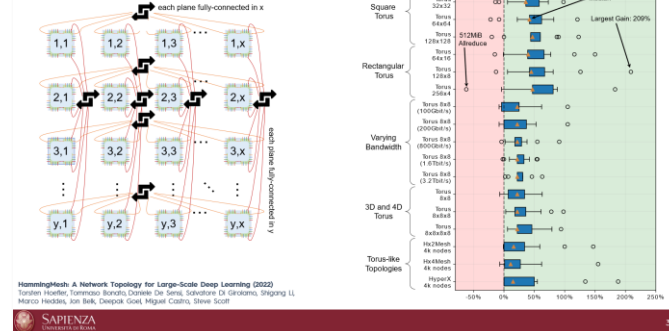
Swing Allreduce



Multiport Swing



Results Summary



Goodput Gain vs. Best Known Algo. for Allreduce <= 512MiB

