

Generating High-Performance Communication Kernels

Madan Musuvathi

Peng Cheng, Changho Hwang, Abhinav Jangda, Binyang Li, Saeed Maleki, Youshan Miao, Jacob Nelson,
Olli Saarikivi, Aashaka Shah, Yifan Xiong, Yongqiang Xiong, Fan Yang, Ziyue Yang

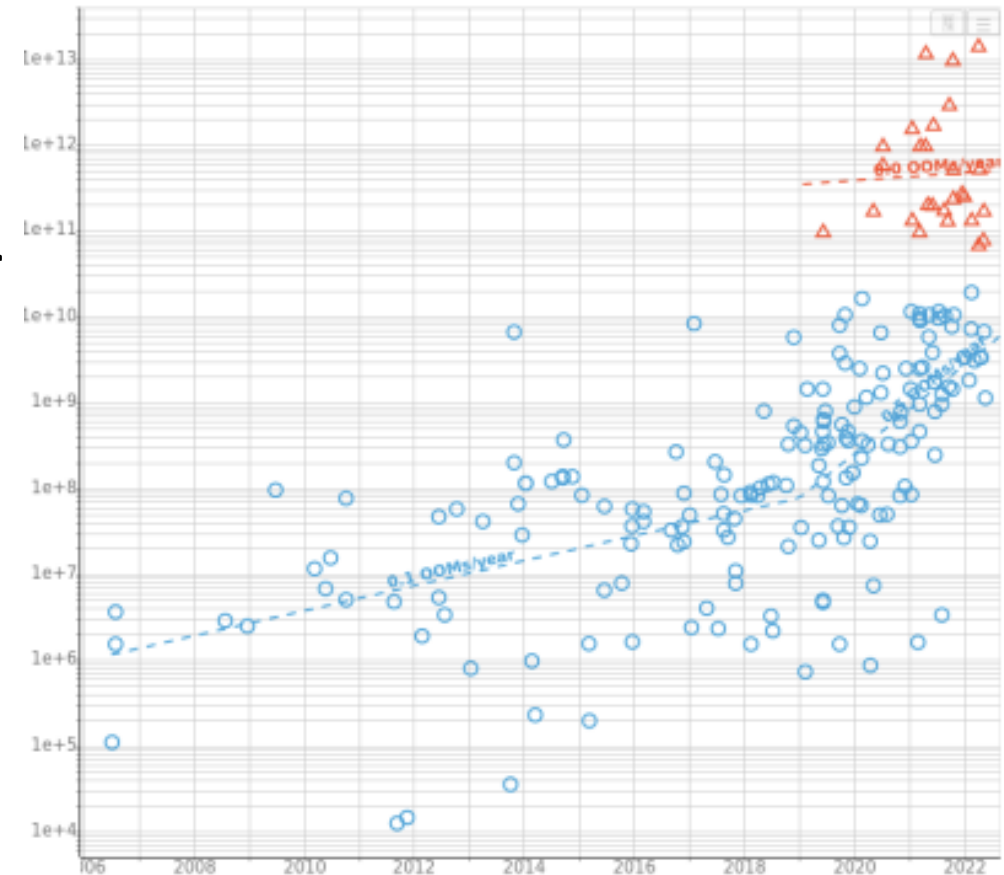
[Microsoft Research]

Meghan Cowan, Todd Mytkowicz [Google]

Vijay Chidambaram [UT Austin], Rachee Singh [Cornell], Zixian Cai [ANU], Zhengyang Liu [Utah]

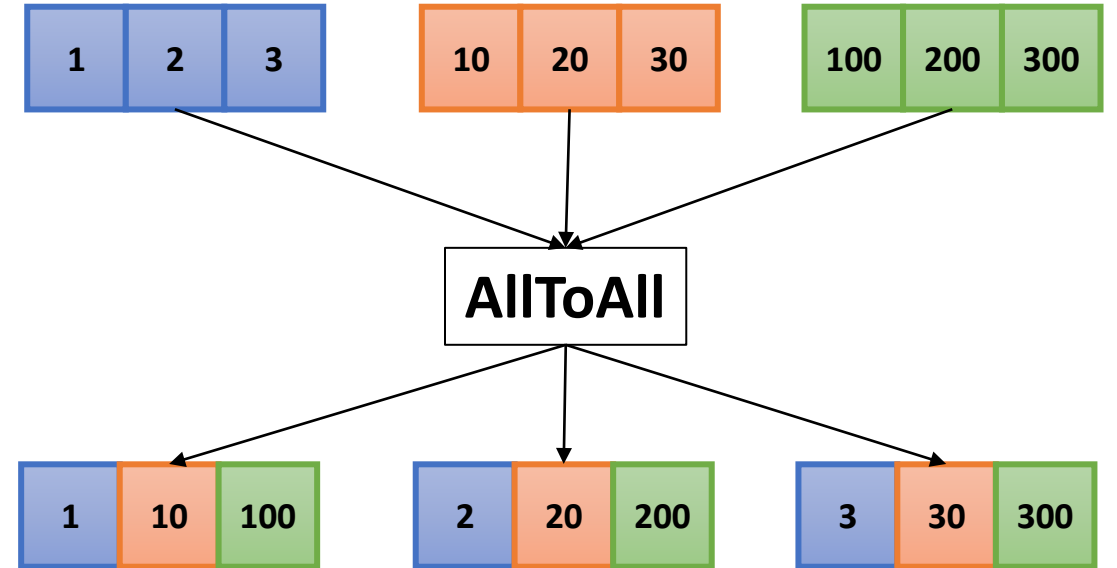
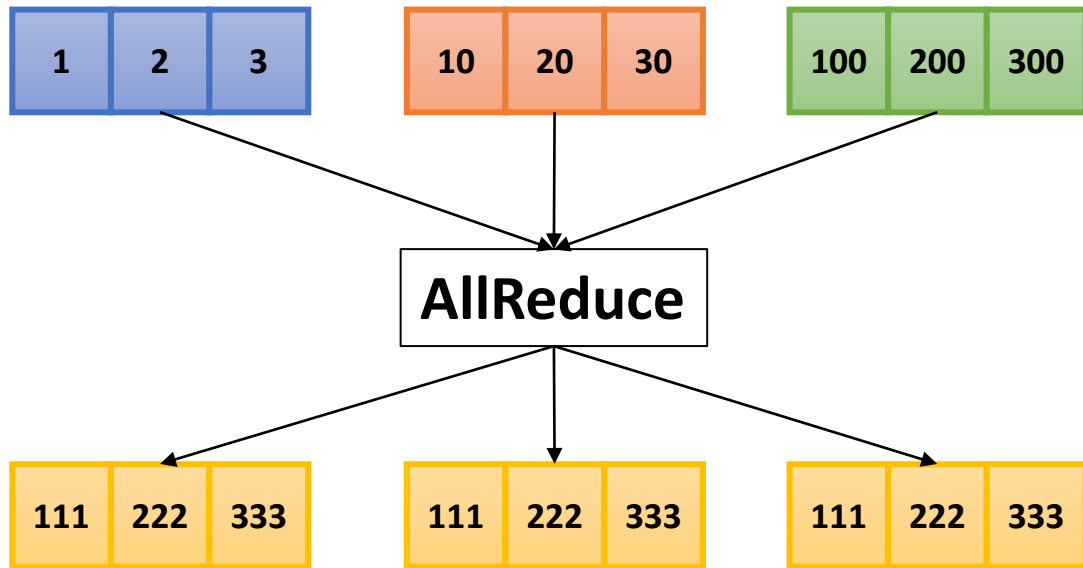
ML computations are increasingly communication bound

- Model sizes increasing $\sim 10\times$ / year
- Need to distribute computations for both inference and training

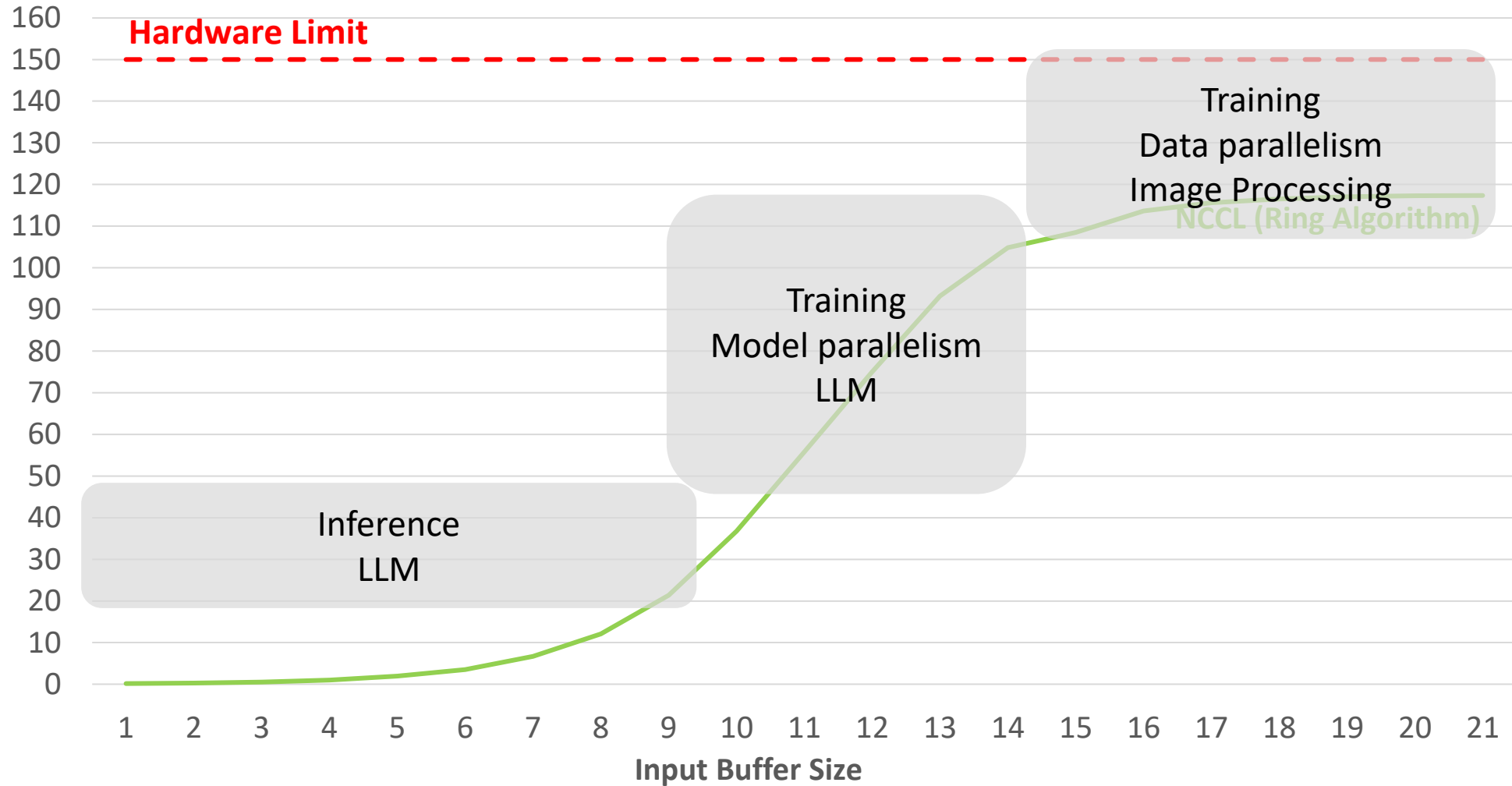


Villalobos et al. '22

Collective Communication



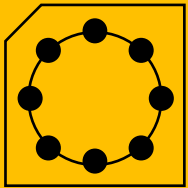
AllGather Throughput (GBps) on DGX1 (8x V100)



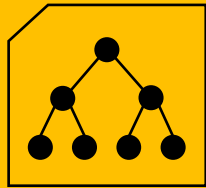
Implementing Communication Kernels

Algorithmic choices

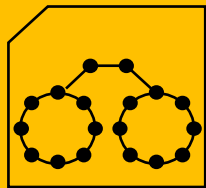
Ring



Tree



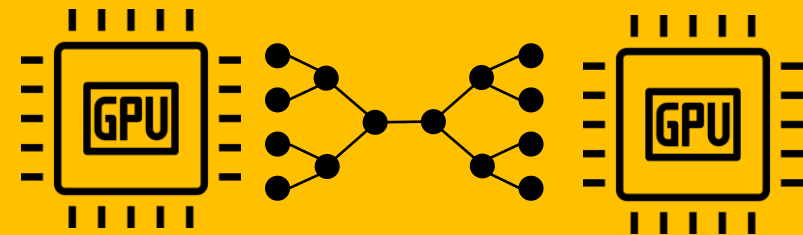
Hierarchical



...

Route data
Distribute compute
Latency vs bandwidth
Intra- vs inter- node

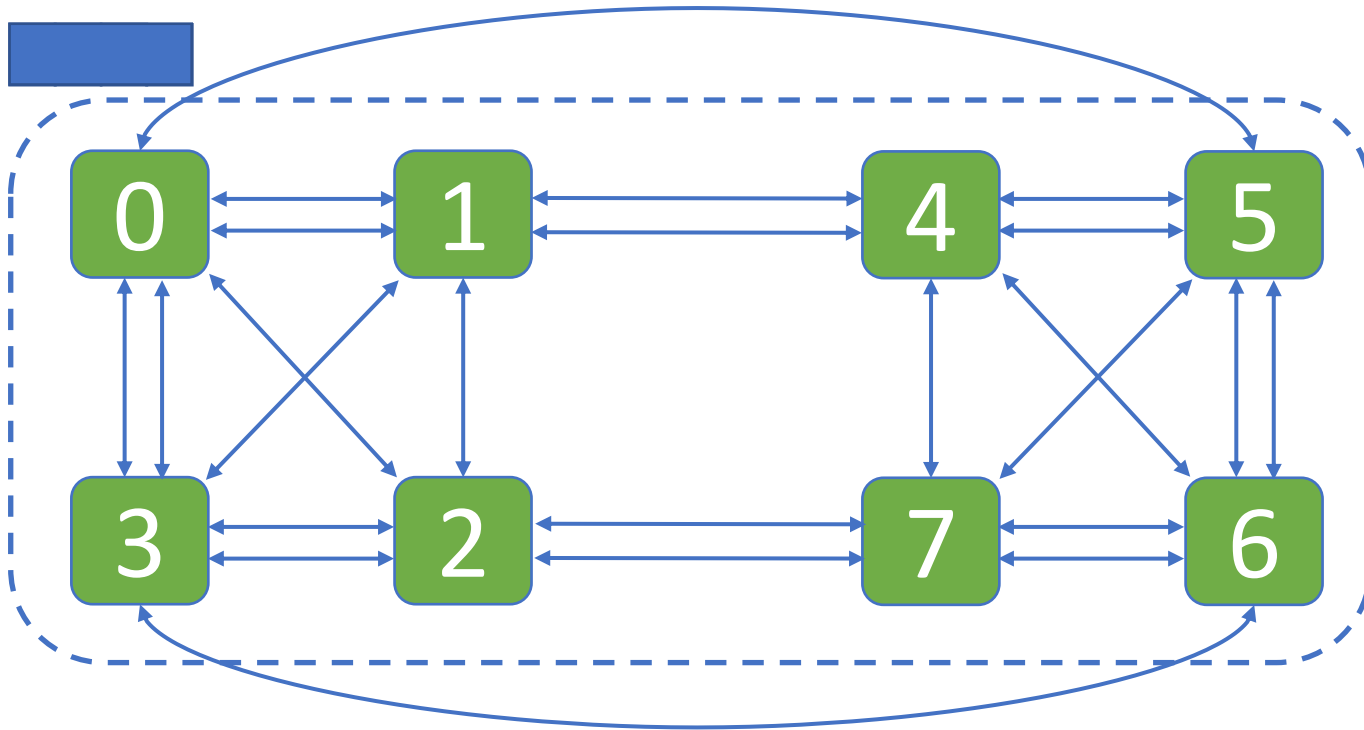
Implementation choices



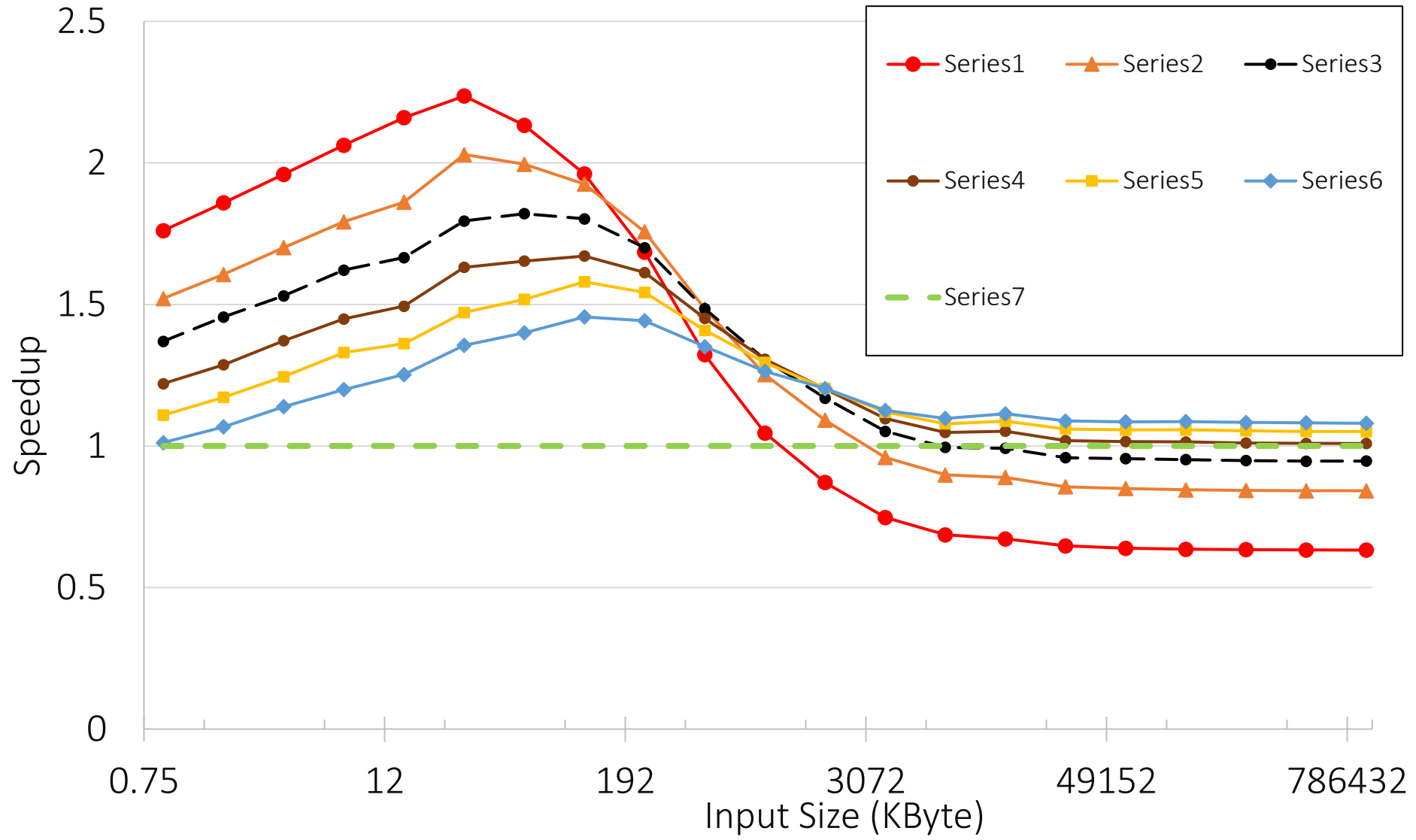
Maximize parallelism
Maximize link utilization
Minimize launch + synch overheads
Deal with relaxed consistency, deadlocks, ...

Synthesizing Optimal Collective Algorithms [PPoPP '21]

Pareto optimality: minimize latency (bandwidth) for a given bandwidth (latency) cost

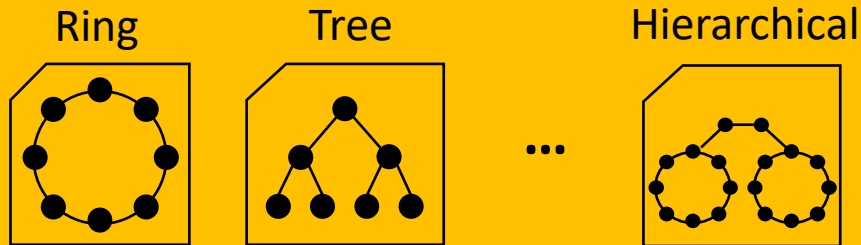


- Split input into c chunks
- Bound routing steps s
- Search (with z3) for algorithms
- Pareto frontier can be generated by minimizing s for various c



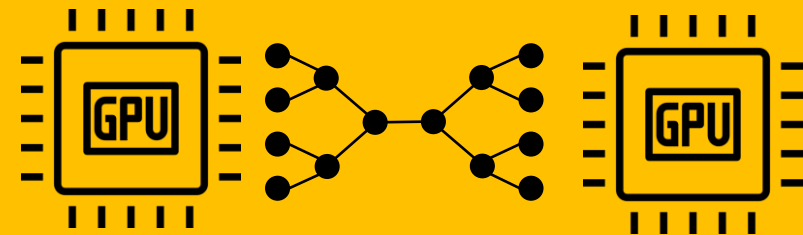
Implementing Communication Kernels

Algorithmic choices



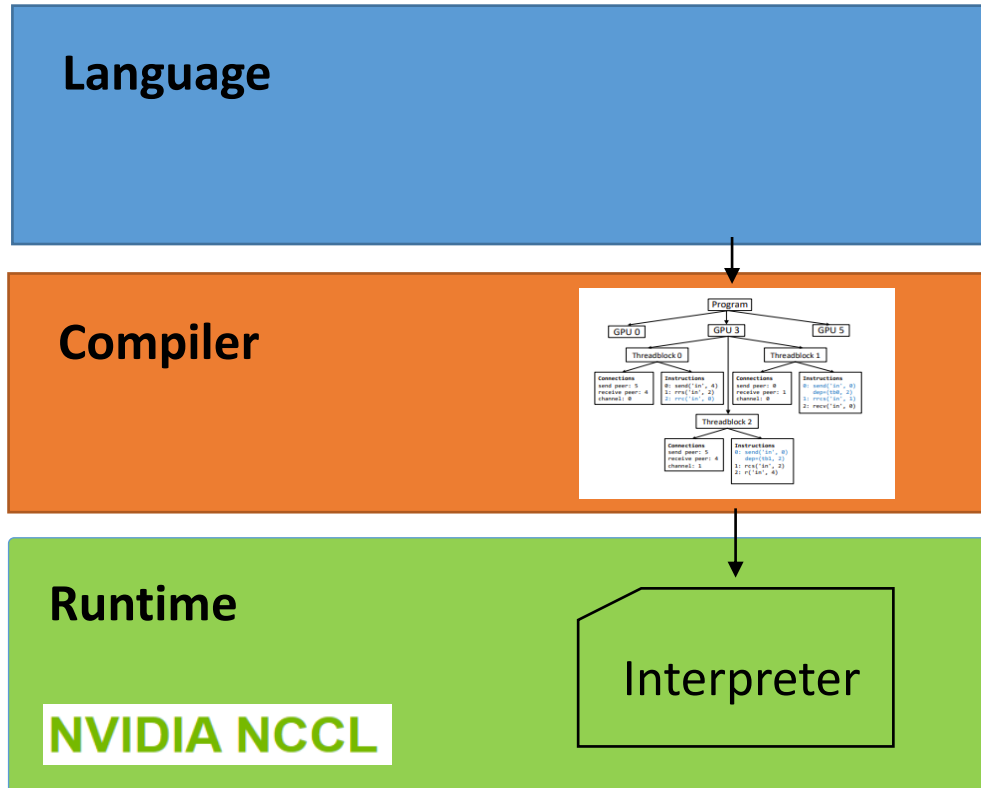
Route data
Distribute compute
Latency vs bandwidth
Intra- vs inter- node

Implementation choices



Maximize parallelism
Maximize link utilization
Minimize launch + synch overheads
Deal with relaxed consistency, deadlocks, ...

MSCCL [ASPLOS '23]

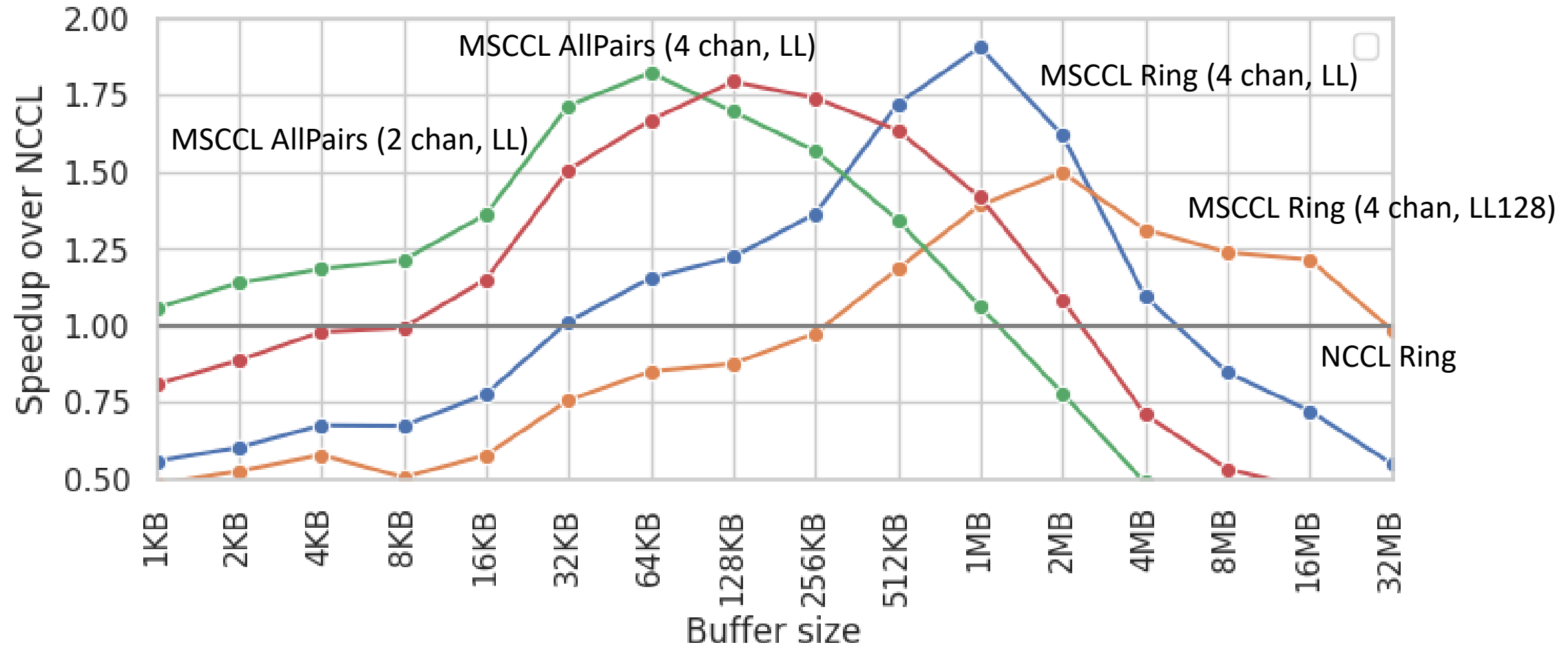


Pythonic chunk-oriented DSL

```
1 # N: number of nodes
2 # n = 0 to G-1: GPU id
3 def alltoall(N, G):
4     for n in range(N):
5         for g in range(G):
6             for m in range(N):
7                 for i in range(G):
8                     c = chunk((m,i), 'in', (n,g))
9                     c.copy((n,g), 'out', (m,i))
10                    c.copy((n,g), 'out', (m,i))
11                else:
12                    c.copy((m,g), 'sc', (n,i))
13
14 # Coalesced IB send
15 c = chunk((m,g), 'sc', n*G, sz=G)
16 c.copy((n,g), 'out', m*G)
```

Map chunk-oriented comp. to threads blocks
Parallelization/tiling/scheduling optimizations
Ensure correctness (no deadlocks/data races)
Embedded interpreter in NCCL
Reuse NCCL transport mechanisms
API compatible with NCCL

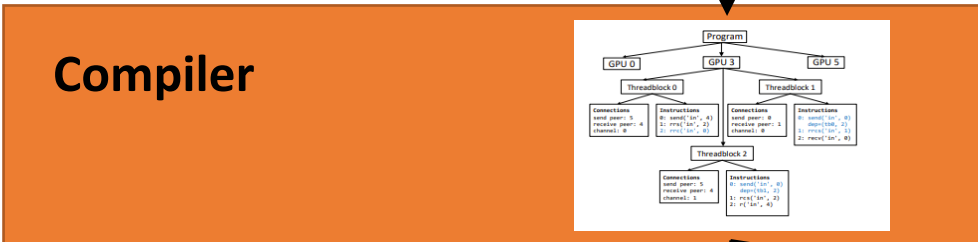
All-Reduce on 8xA100



MSCCL

Language

```
1 # N: number of nodes
2 # G: number of GPUs per node
3 def alltoall(N, G):
4     for n in range(N):
5         for g in range(G):
6             for i in range(G):
7                 c = chunk((n,i), 'in', (n,g))
8                 if 0 == 0:
9                     c.copy((n,g), 'out', (n,i))
10                else:
11                    c.copy((n,g), 'sc', (n,i))
12            # Coalesced IB send
13            c = chunk((n,g), 'sc', n*G, sz=G)
14            c.copy((n,g), 'out', n*G)
```



Runtime

NVIDIA NCCL

Interpreter



Interpreter

AMD RCCL