MSCCL++: Rethinking GPU Communication Abstractions for Cutting-Edge Al Applications

Presented by Yang Zhou Thur 10/23

Abstract

Big problem:

- Modern cutting-edge AI applications are being developed over fast-evolving, heterogeneous, nascent hardware devices.
- This requires frequent reworking of the AI software stack to adopt bottom-up changes from new hardware, which takes time for general-purpose software libraries (NCCL)

MSCCL++ provides both portability and performance by essentially a reusable lib

- a primitive interface provides a minimal hardware abstraction
- higher-level portable interfaces and specialized implementations

Background

Various communication channels:

- CPU-GPU: PCIe, NVLink-C2C
- Intra-node: NVLink, xGMI,
- Inter-node: RDMA (IB, RoCE)

Key challenge: people want high communication performance, but this requires writing custom communication code, often from scratch and taking long time.

ML workloads move fast

Motivation examples

The custom communication of TensorRT-LLM outperforms NCCL in a wide range of LLM scenarios, especially when the data size is relatively small, while TensorRTLLM still uses NCCL for larger data sizes.

What is 1-shot allreduce?

Key idea

Separate high-level abstractions and optimizations from primitive hardware abstractions

- Application developers to optimize and fine-tune by enabling precise control of the hardware
- Library developers to support new hardware features by only providing a shallow layer of abstraction over it

Basically, they are building a reusable communication library with an abstraction for various communication channels.

NCCL examples

Hard-coded GPU kernels

Send, recv not flexible enough

Inflexible synchronization

```
1 global ringRS(in, nelem, ring)
    send = ring.sendbuff; recv = ring.recvbuff;
    sz = ring.buffSz;
    prim = Primitives<half>
      (tid, ring.buffSz, ring.prev, ring.next);
    for (off = 0; off < nelem; off += sz)</pre>
      //step 0: push data to the next GPU
      idx = off + ring.ranks[ring.ranks-1]*sz;
10
      prims.copy(in+idx, send, sz);
      prims.send();
11
12
      //k-2 steps: reduce and copy to next GPU
13
14
      for (j = 2; j < ring.ranks; ++j)
15
        rankDest = ring.ranks[ring.ranks-j];
16
        idx = off + rankDest * buffSz;
17
        prims.recv();
        prim.reduce(in+idx, recv, send, sz, "+");
18
19
        prims.send();
20
21
      //step k-1: write the result for this rank
22
      idx = off + ring.rank * buffSz;
23
      prims.recv();
24
      prims.reduce(recv, in+idx, sz, "+");
```

Figure 1. Ring ReduceScatter (simplified) in NCCL.

MSCCL++ overview

High-level APIs

+

Low-level interfaces

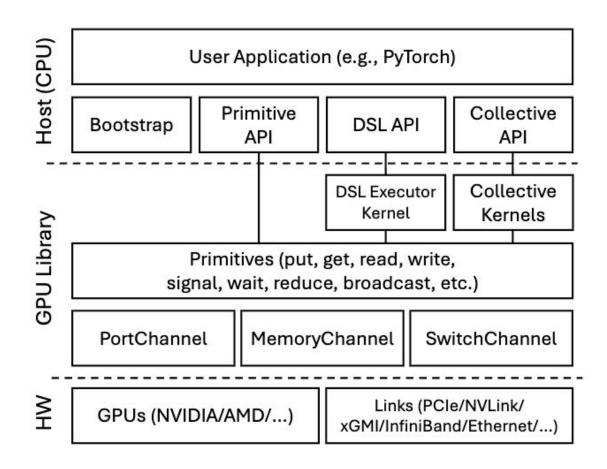


Figure 3. Overview of the MSCCL++ stack.

Communication channels

PortChannel: DMA engines on GPUs or RDMA NICs

MemoryChannel: peer-to-peer memory access over NVLink/xGMI

SwitchChannel: interconnection-switch-enabled multimem memory

Communication Primitives

All MSCCL++ primitives are defined as a method of a channel.

```
//async
put(src0, dst1, size)
//unsafe to reuse src0
signal() //async
flush() //sync
//safe to reuse src0
//unsafe to read dst1
wait() //sync
//safe to read dst1
//safe to read dst1
//safe to read dst1
//safe to read dst1
```

Figure 4. MSCCL++ data transfer abstractions. put asynchronously transfers data from one GPU to another. signal and wait synchronize data transfer between GPUs. flush ensures the completion of preceded data transfer.

An example

Implementing a different ReduceScatter

```
1 global allPairsRS(count, gpus, channels[gpus])
    sz = channel.scratchSz/gpus.num
    count = count/gpus.num
5
    for (off = 0; off < count; off += sz)</pre>
      //Send 1/Nth data to each GPU
      for (g = 0; g < gpus.num; g++)
        idx = off + g * count;
        channels[g].put(idx, sz*g, sz)
        channels[g].signal()
10
11
12
      //Reduce every pair of GPU
13
      for (g = 0; g < gpus.num-1; g++)
14
        channels[g].wait()
15
        reduce(in + off, channels[g].scratch)
16
17
      //barrier on all gpus
18
      multiDeviceBarrier();
```

Figure 5. All-pairs ReduceScatter kernel in MSCCL++. Channels are initialized with source as input and destination as scratch buffer.

Implementation deepdive

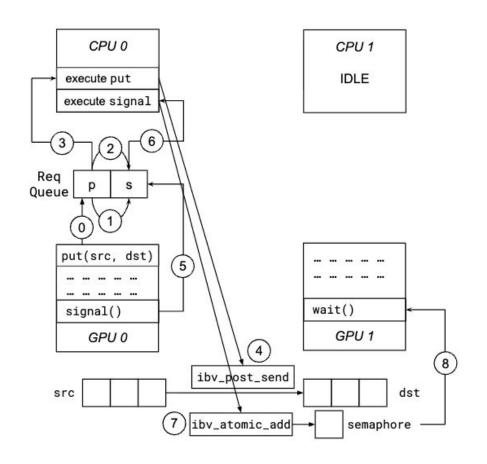


Figure 7. PortChannel workflow for IB from ① when GPU 0 calls put primitive to ⑧ when GPU 1 receives the data.

Testbed

Env. Name	GPU	Intra-node Link	Network
A100-40G	NVIDIA A100 (40G) (8x/node)	NVLink 3.0	Mellanox HDR InfiniBand (200 Gb/s, 1x NIC/GPU)
A100-80G	NVIDIA A100 (80G) (8x/node)	NVLink 3.0	Mellanox HDR InfiniBand (200 Gb/s, 1x NIC/GPU)
H100	NVIDIA H100 (8x/node)	NVLink 4.0	Quantum-2 CX7 InfiniBand (400 Gb/s, 1x NIC/GPU)
MI300x	AMD MI300x (8x/node)	Infinity Fabric Gen 4	Quantum-2 CX7 InfiniBand (400 Gb/s, 1x NIC/GPU)

Table 1. List of environments used for evaluation.

Some results

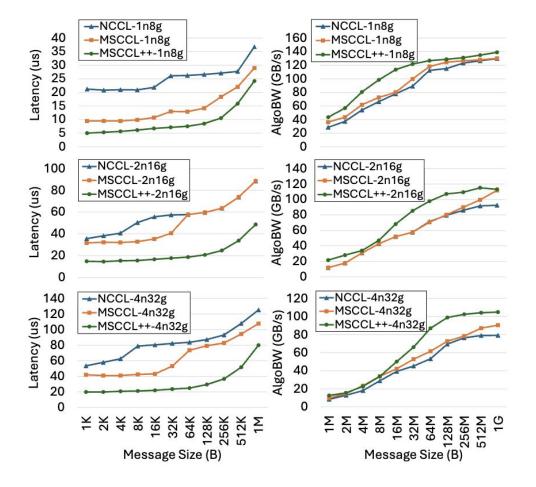


Figure 8. AllReduce, A100-40G. MnNg means M nodes and N GPUs. X-axis is message size (at the bottom of the figure).

Key takeaways

ML workloads is fast evolving, so do the system software (and so do hardware)

Communication is a cornerstone in ML workloads

We need a reusable library to develop custom communication strategies.

What so far we have covered

Datacenter networking:

- Microsecond and tail
- Topology design
- Congestion control
- How to use RDMA (efficiently)

Host networking:

- RDMA for GPU comm
- Efficient host TCP
- NCCL stack

What we will cover

LLM inference:

- Memory management
- Compute management
- Attention efficiency
- Holistic optimization

LLM training:

- Attention efficiency
- Training parallelism
- MoE (Mixture-of-Expert)
- Failure, auto-parallelism