

Collective Communication Optimizations (CCO)

An IETF118 Side Meeting, Prague
14:30-16:00 Thursday, Nov 9, 2023

Kehan Yao (China Mobile), Yizhou Li (Huawei)

Note Well

<https://www.ietf.org/about/note-well/>

This is a reminder of IETF policies in effect on various topics such as patents or code of conduct. It is only meant to point you in the right direction. Exceptions may apply. The IETF's patent policy and the definition of an IETF "contribution" and "participation" are set forth in BCP 79; please read it carefully.

As a reminder:

- By participating in the IETF, you agree to follow IETF processes and policies.
- If you are aware that any IETF contribution is covered by patents or patent applications that are owned or controlled by you or your sponsor, you must disclose that fact, or not participate in the discussion.
- As a participant in or attendee to any IETF activity you acknowledge that written, audio, video, and photographic records of meetings may be made public.
- Personal information that you provide to IETF will be handled in accordance with the IETF Privacy Statement.
- As a participant or attendee, you agree to work respectfully with other participants; please contact the ombudsteam (<https://www.ietf.org/contact/ombudsteam/>) if you have questions or concerns about this.

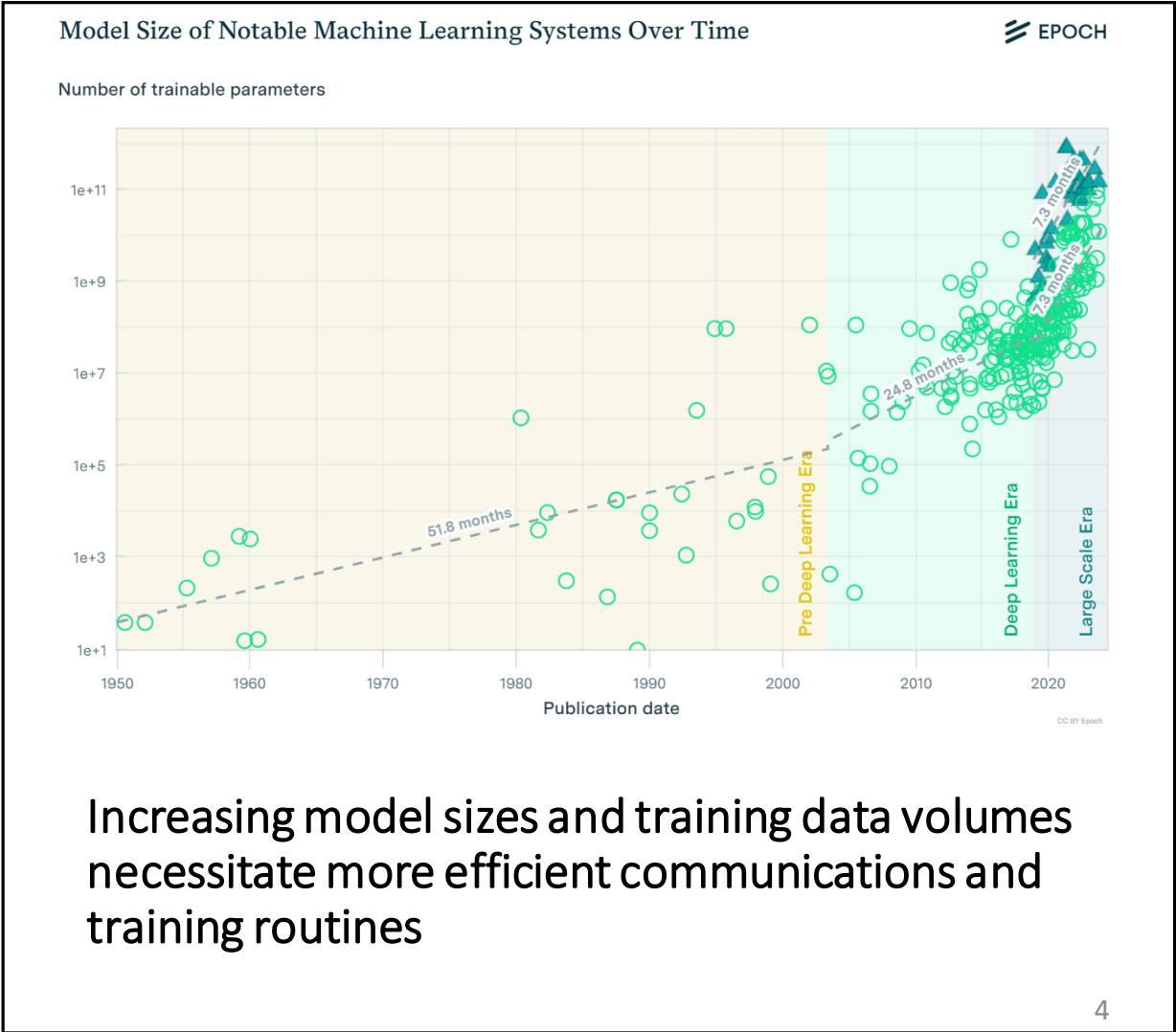
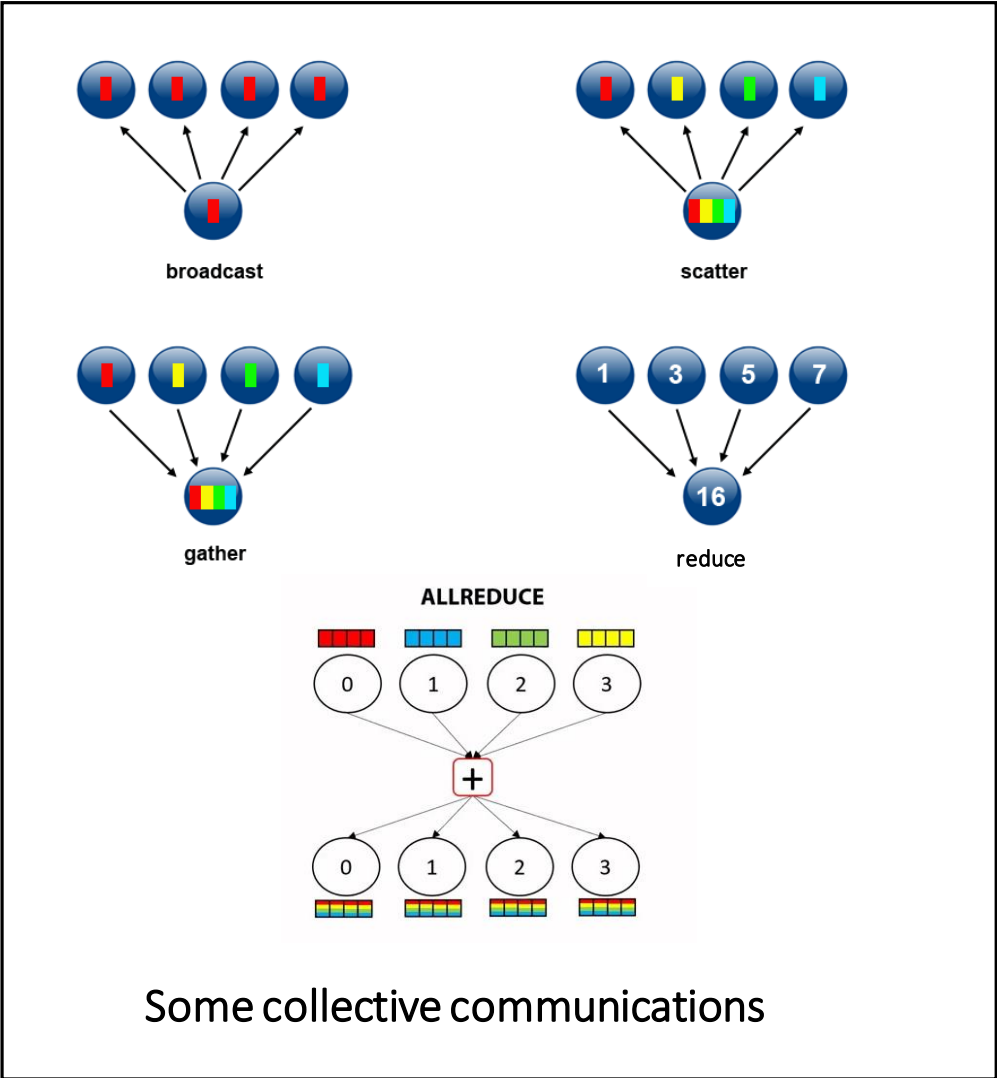
Definitive information is in the documents listed below and other IETF BCPs. For advice, please talk to WG chairs or ADs:

- [BCP 9](#) (Internet Standards Process)
- [BCP 25](#) (Working Group processes)
- [BCP 25](#) (Anti-Harassment Procedures)
- [BCP 54](#) (Code of Conduct)
- [BCP 78](#) (Copyright)
- [BCP 79](#) (Patents, Participation)
- <https://www.ietf.org/privacy-policy/> (Privacy Policy)

Agenda

- 1. Opening (5 mins)**
- 2. Use cases, problem space and requirements**
Kehan Yao (China Mobile) (20 mins)
- 3. Challenges in hardware offloading of collective operations**
Alex Margolin (Hebrew University of Jerusalem)
- 4. Signaling In-Network Computing Operations (SINC)**
David Lou (Huawei) (20 mins)
- 5. In Network Compute**
Surendra Anubolu (Broadcom Inc) (20 mins)
- 6. Open Discussions**

Collective communications are essential in the distributed training



The diagram shows the evolution of collective communication through two scenarios, (a) and (b), connected by an 'Evolve' arrow.

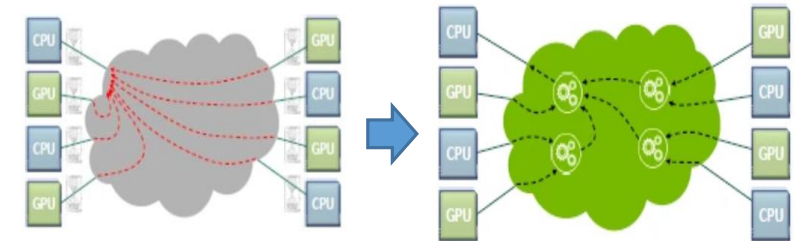
(a) Classic HPC scenarios:

- Programming Model:** MPI, OpenSHMEM, UCX, UCC.
- Collective:** Broadcast, Gather(v), Scatter(v), All-Gather, All-to-All, Reduce, All-Reduce, Reduce-Scatter.
- P2P Communication:** (Empty box)
- Hardware Topology:** Torus, Fat-Tree, Dragonfly, Hypercube.
- Hardware:** CPU, GPU, Accelerator.
- Highspeed Interconnect:** Ethernet, InfiniBand, RoCE.

(b) Emerging deep learning scenarios:

- xCCL:** NCCL, Gloo, MSCCL, ACCL, RCCL, oneCCL.
- Collective:** All-Reduce, All-Gather, All-to-All, Broadcast, Reduce-Scatter.
- P2P Communication:** (Empty box)
- Hardware Topology:** Ring, Torus, Fat-Tree.
- Hardware:** CPU, NPU(GPU/TPU).
- Highspeed Interconnect:** Ethernet, InfiniBand, RoCE, NVLink.

On the right side of the diagram, there is a vertical stack of labels: I E T F, ART, TSV, and INT/RTG, separated by a dashed blue line.



- xCCL: A survey of industry-led collective communication deep learning.
JOURNAL OF COMPUTER SCIENCE AND TECHNOLOGY 38(1): 166–1 libraries for 95 Jan. 2023.