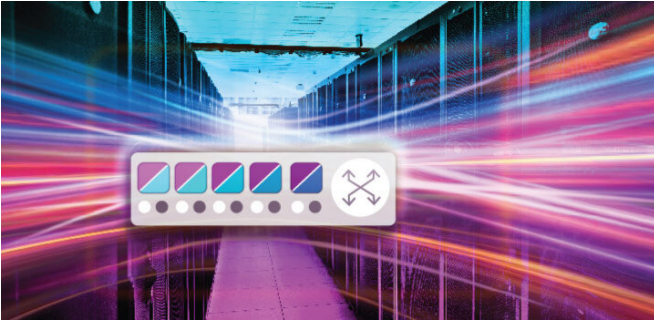


Machine learning at speed with in-network aggregation

12 April 2021



Technology developed through a KAUST-led collaboration with Intel, Microsoft and the University of Washington can dramatically increase the speed of machine learning on parallelized computing systems. Credit: KAUST; Anastasia Serin

Inserting lightweight optimization code in high-speed network devices has enabled a KAUST-led collaboration to increase the speed of machine learning on parallelized computing systems five-fold.

This "in-network aggregation" technology, developed with researchers and systems architects at Intel, Microsoft and the University of Washington, can provide dramatic speed improvements using readily available programmable network [hardware](#).

The fundamental benefit of artificial intelligence (AI) that gives it so much power to "understand" and interact with the world is the machine-learning step, in which the [model](#) is trained using large sets of labeled training data. The more data the AI is trained on, the better the model is likely to perform when exposed to new inputs.

The recent burst of AI applications is largely due to better machine learning and the use of larger models and more diverse datasets. Performing the

machine-learning computations, however, is an enormously taxing task that increasingly relies on large arrays of computers running the learning algorithm in parallel.

"How to train deep-learning models at a large scale is a very challenging problem," says Marco Canini from the KAUST research team. "The AI models can consist of billions of parameters, and we can use hundreds of processors that need to work efficiently in parallel. In such systems, communication among processors during incremental model updates easily becomes a major performance bottleneck."

The team found a potential solution in new network technology developed by Barefoot Networks, a division of Intel.

"We use Barefoot Networks' new programmable dataplane networking hardware to offload part of the work performed during distributed machine-learning training," explains Amedeo Sapio, a KAUST alumnus who has since joined the Barefoot Networks team at Intel. "Using this new programmable networking hardware, rather than just the network, to move data means that we can perform computations along the network paths."

The key innovation of the team's SwitchML platform is to allow the network hardware to perform the data aggregation task at each synchronization step during the model update phase of the machine-learning process. Not only does this offload part of the computational load, it also significantly reduces the amount of data transmission.

"Although the programmable switch dataplane can do operations very quickly, the operations it can do are limited," says Canini. "So our solution had to be simple enough for the hardware and yet flexible enough to solve challenges such as limited onboard memory capacity. SwitchML addresses this challenge by co-designing the communication

[network](#) and the distributed training algorithm, achieving an acceleration of up to 5.5 times compared to the state-of-the-art approach."

More information: Scaling Distributed Machine Learning with In-Network Aggregation.
arxiv.org/abs/1903.06701 arXiv:1903.06701v2
[cs.DC]

Provided by King Abdullah University of Science and Technology

APA citation: Machine learning at speed with in-network aggregation (2021, April 12) retrieved 5 December 2021 from <https://techxplore.com/news/2021-04-machine-in-network-aggregation.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.