

Xeon-class CPU," Shrivastava said.

Deep learning networks were inspired by biology, and their central feature, artificial neurons, are small pieces of computer code that can learn to perform a specific task. A deep learning network can contain millions or even billions of [artificial neurons](#), and working together they can learn to make human-level, expert decisions simply by studying large amounts of data. For example, if a deep neural network is trained to identify objects in photos, it will employ different neurons to recognize a photo of a cat than it will to recognize a school bus.

"You don't need to train all the neurons on every case," Medini said. "We thought, 'If we only want to pick the neurons that are relevant, then it's a search problem.' So, algorithmically, the idea was to use locality-sensitive hashing to get away from matrix multiplication."

Hashing is a data-indexing method invented for internet search in the 1990s. It uses numerical methods to encode large amounts of information, like entire webpages or chapters of a book, as a string of digits called a hash. Hash tables are lists of hashes that can be searched very quickly.

"It would have made no sense to implement our algorithm on TensorFlow or PyTorch because the first thing they want to do is convert whatever you're doing into a matrix multiplication problem," Chen said. "That is precisely what we wanted to get away from. So we wrote our own C++ code from scratch."

Shrivastava said SLIDE's biggest advantage over back-propagation is that it is data parallel.

"By data parallel I mean that if I have two data instances that I want to train on, let's say one is an image of a cat and the other of a bus, they will likely activate different neurons, and SLIDE can update, or train on these two independently," he said. "This is much a better utilization of parallelism for CPUs.

"The flipside, compared to GPU, is that we require a big memory," he said. "There is a cache hierarchy in main memory, and if you're not careful with it you

can run into a problem called cache thrashing, where you get a lot of cache misses."

Shrivastava said his group's first experiments with SLIDE produced significant cache thrashing, but their training times were still comparable to or faster than GPU training times. So he, Chen and Medini published the initial results on *arXiv* in March 2019 and uploaded their code to GitHub. A few weeks later, they were contacted by Intel.

"Our collaborators from Intel recognized the caching problem," he said. "They told us they could work with us to make it train even faster, and they were right. Our results improved by about 50% with their help."

Shrivastava said SLIDE hasn't yet come close to reaching its potential.

"We've just scratched the surface," he said. "There's a lot we can still do to optimize. We have not used vectorization, for example, or built-in accelerators in the CPU, like Intel Deep Learning Boost. There are a lot of other tricks we could still use to make this even faster."

Shrivastava said SLIDE is important because it shows there are other ways to implement [deep learning](#).

"The whole message is, 'Let's not be bottlenecked by multiplication matrix and GPU memory,'" Chen said. "Ours may be the first algorithmic approach to beat GPU, but I hope it's not the last. The field needs new ideas, and that is a big part of what MLSys is about."

More information: SLIDE : In Defense of Smart Algorithms over Hardware Acceleration for Large-Scale Deep Learning Systems, *arXiv*:1903.03129 [cs.DC] arxiv.org/abs/1903.03129

Provided by Rice University

APA citation: Deep learning rethink overcomes major obstacle in AI industry (2020, March 2) retrieved 30 November 2020 from <https://techxplore.com/news/2020-03-deep-rethink-major-obstacle-ai.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.