

Poseidon: A System Architecture for Efficient GPU-based Deep Learning on Multiple Machines

2016 USENIX
Annual Technical
Conference
JUNE 22-24, 2016
DENVER, CO
www.usenix.org/atc16

Hao Zhang, Zhiting Hu, Jinliang Wei, Pengtao Xie, Gunhee Kim, Qirong Ho, Eric Xing
{hao, zhitingh, jinlianw, pengtaox, gunhee, epxing}@cs.cmu.edu

Introduction

We propose a scalable open-source framework for large-scale distributed deep learning on GPU clusters. We build the framework upon the Caffe CNN libraries and the Petuum distributed ML framework as a starting point, but goes further by implementing three key contributions for efficient CNN training on clusters of GPU-equipped machines: (i) a three-level hybrid architecture to support both CPU-only clusters as well as GPU-equipped clusters, (ii) a distributed wait-free backpropagation (DWBP) algorithm to improve GPU utilization and to balance communication, and (iii) a dedicated structure-aware communication protocol (SACP) to minimize communication overheads. We empirically show that our framework converges to the same objective value as a single machine, and achieves state-of-art training speedup across multiple models and well-established datasets, using a commodity GPU cluster of 8 nodes.

Background: Iterative-Convergent Algorithms

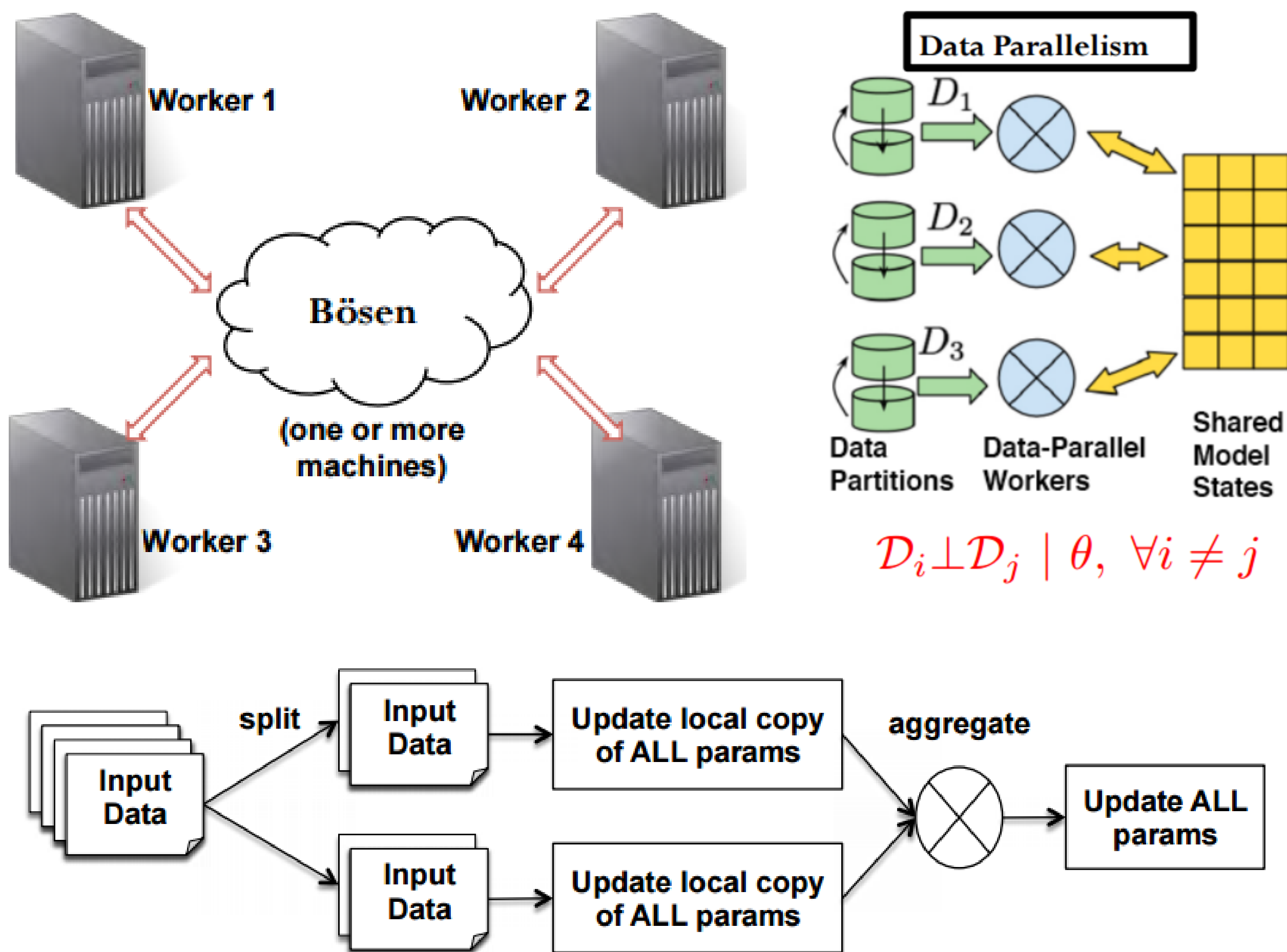
Iterative-Convergent Algorithm in General

The *iterative-convergent* algorithm can be represented as follows.

$$A^{(t)} = F(A^{(t-1)}, \Delta_\ell(A^{(t-1)}, D))$$

In large-scale machine learning, both data D and model A can be very large.

$$A^{(t)} = F(A^{(t-1)}, \sum_{p=1}^P \Delta_\ell(A^{(t-1)}, D_p))$$

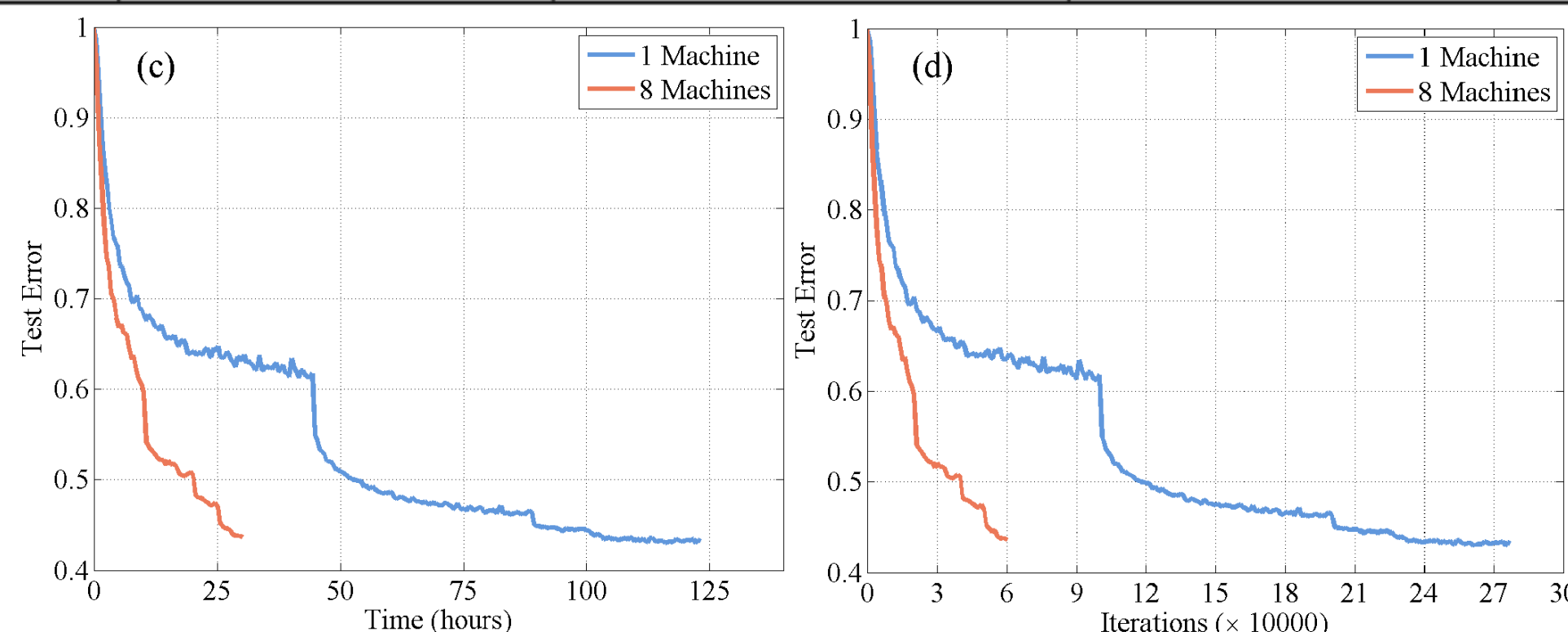


Experiments

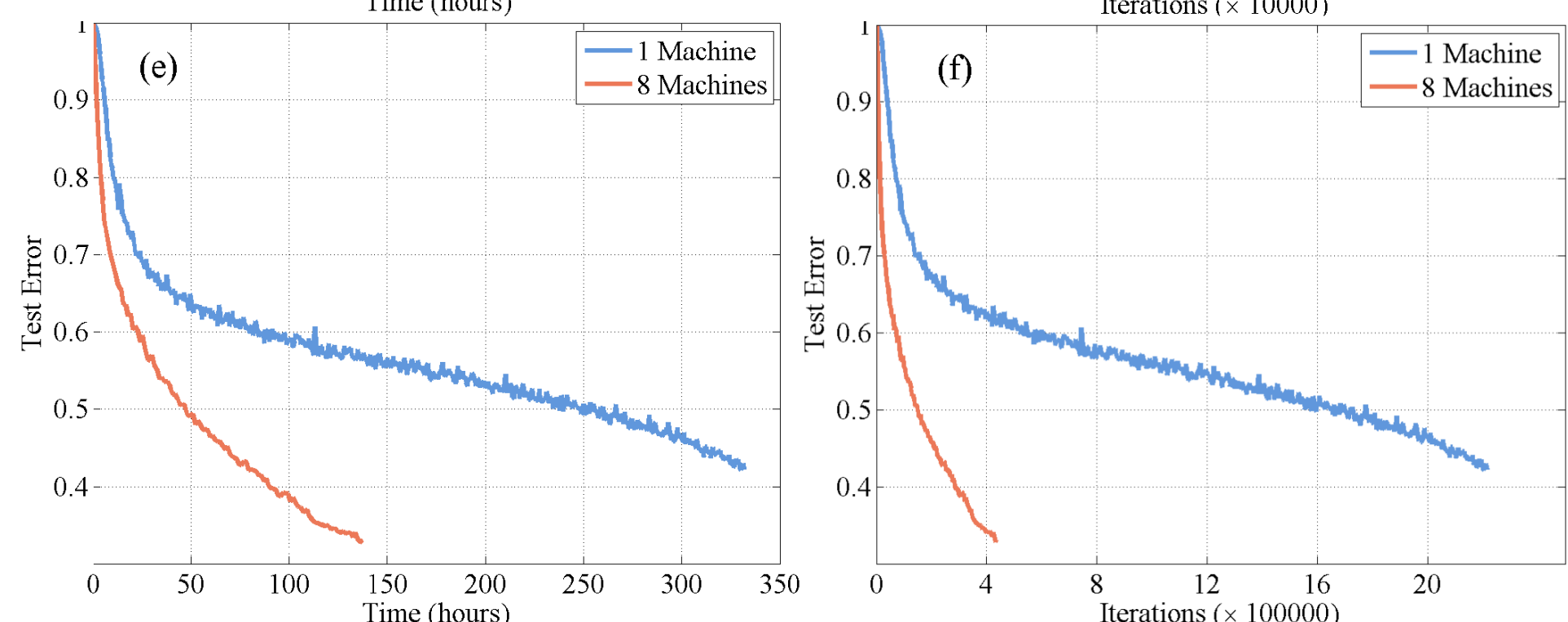
Accelerate the Training of Modern CNNs

Dataset	# of Images	Size of images	# of categories
ILSVRC2012	1.3M	$256 \times 256 \times 3$	1000
ImageNet22K	15M	$256 \times 256 \times 3$	21841

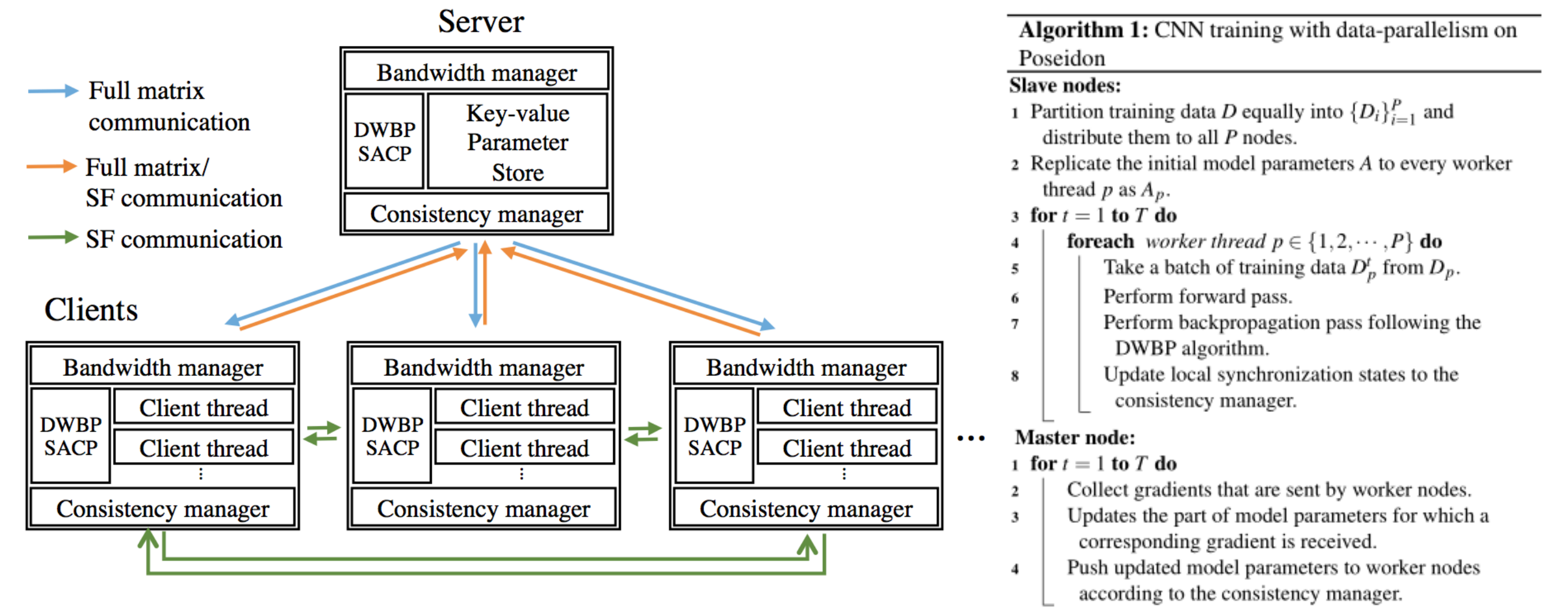
AlexNet Training
4.5x speedup 8 nodes
56.5% top-1 accuracy



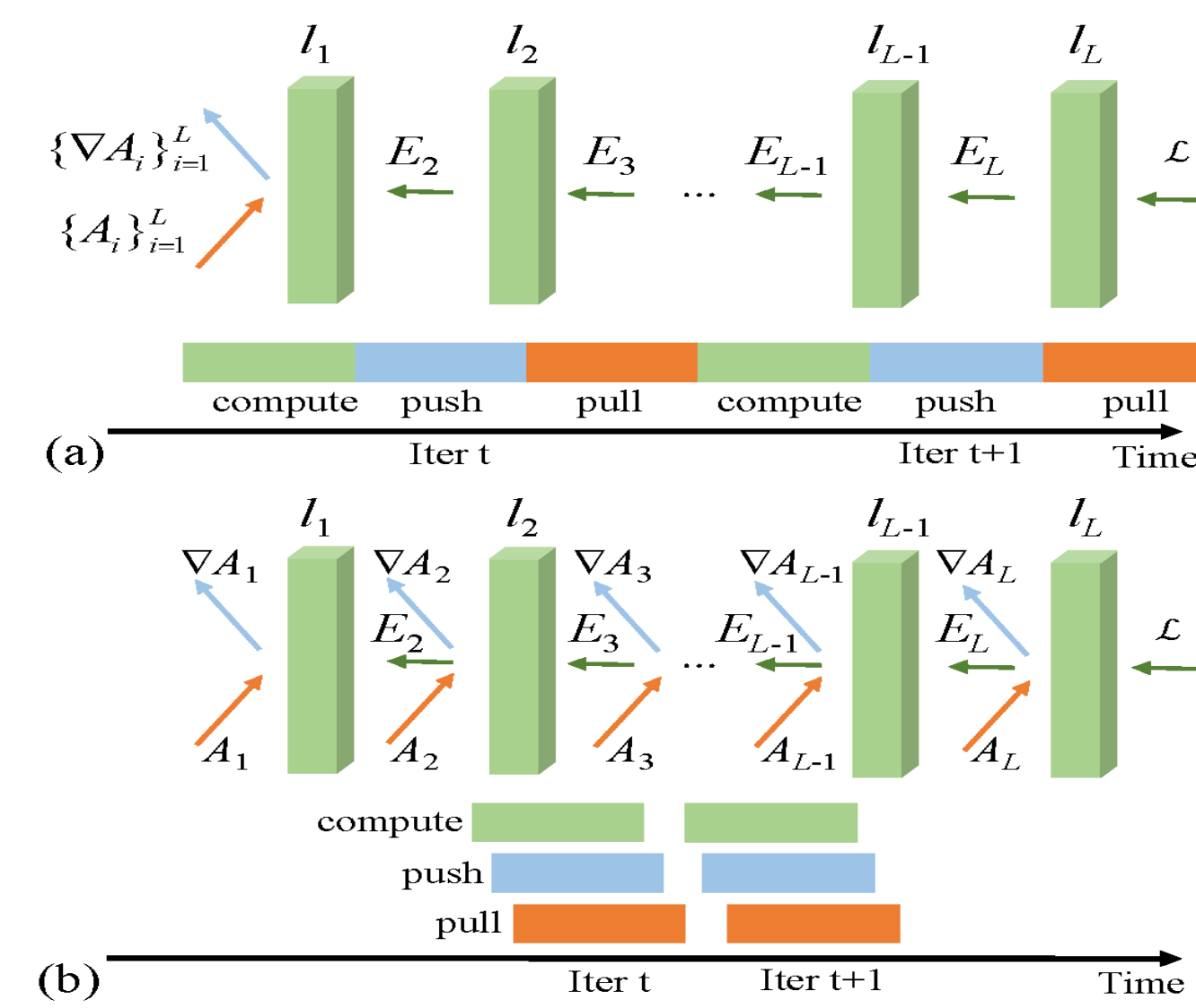
GoogLeNet Training
4x speedup 8 nodes
67.1% top-1 accuracy



System Architecture



Distributed Wait-free Backpropagation (DWBP)



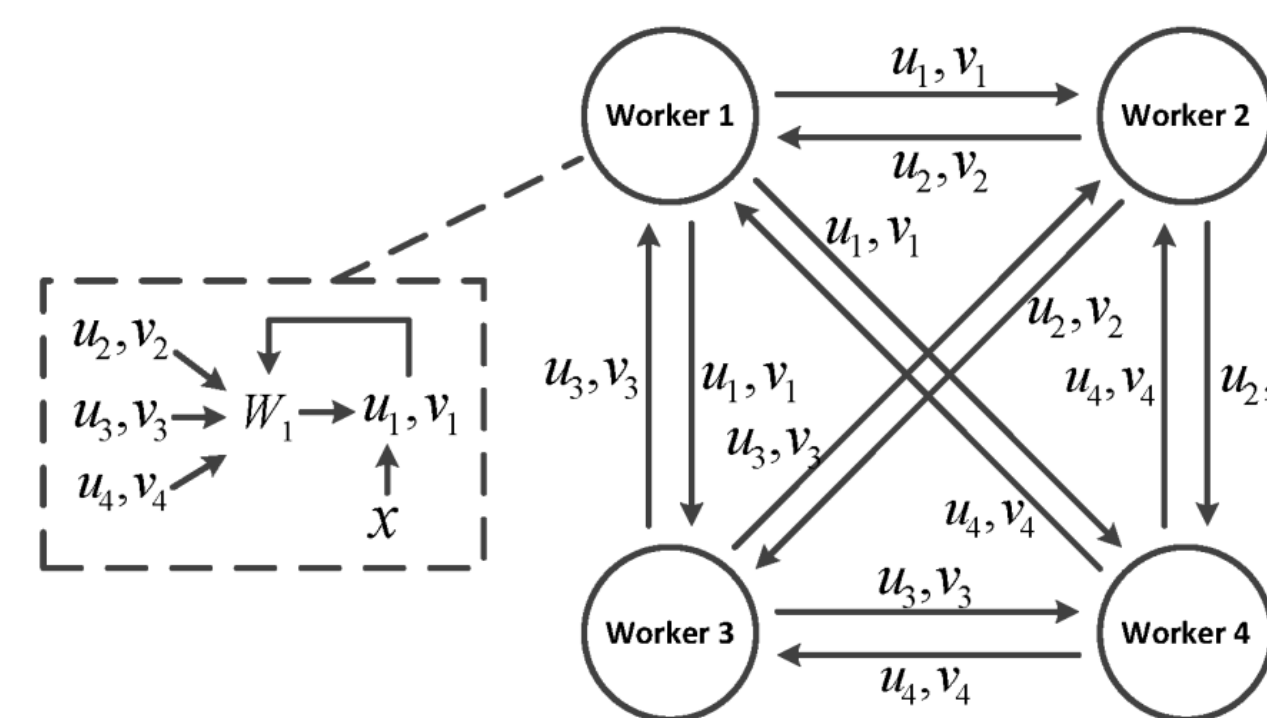
Algorithm 2: The Distributed Wait-free Backpropagation (DWBP) Algorithm

At iteration t on worker p :

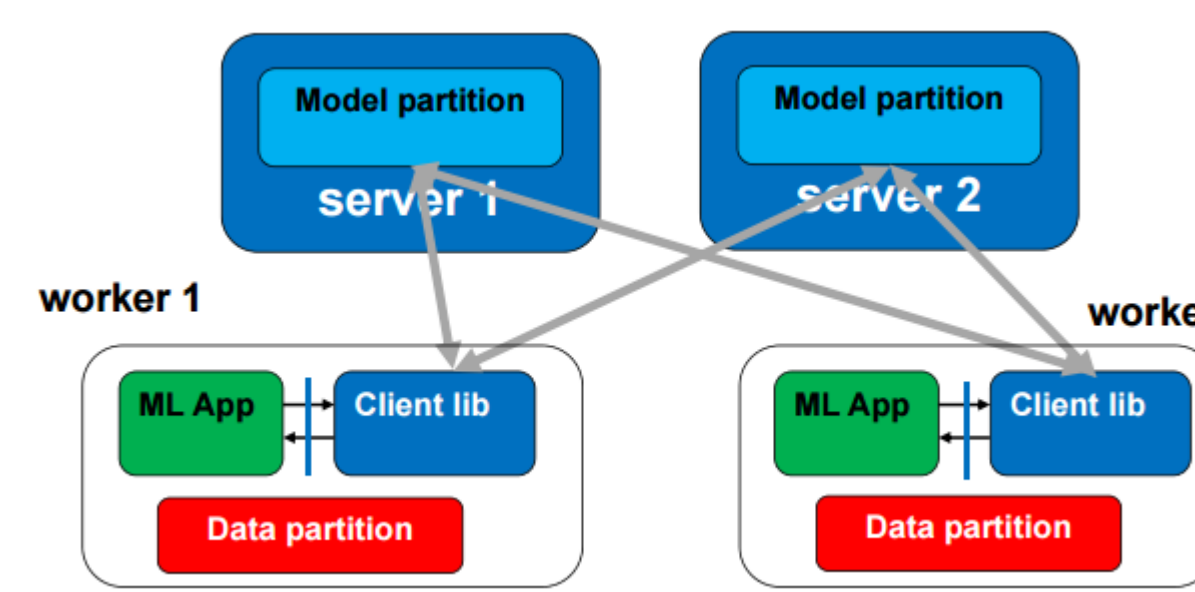
Input: Loss ℓ .

- for $i = L$ to 1 do
 - if $i == L$ then
 - Compute gradients $\nabla A_i = \frac{\partial \ell}{\partial A_i}$ using ℓ .
 - else
 - Receive error message E_{i+1} from layer $i+1$.
 - Compute gradients $\nabla A_i = \frac{\partial \ell}{\partial A_i}$ using E_{i+1} .
 - if $i \neq 1$ then
 - Compute error message E_i and pass to layer $i-1$.
 - Communicate: push out ∇A_i and pull in updated A_i following the SACP protocol;

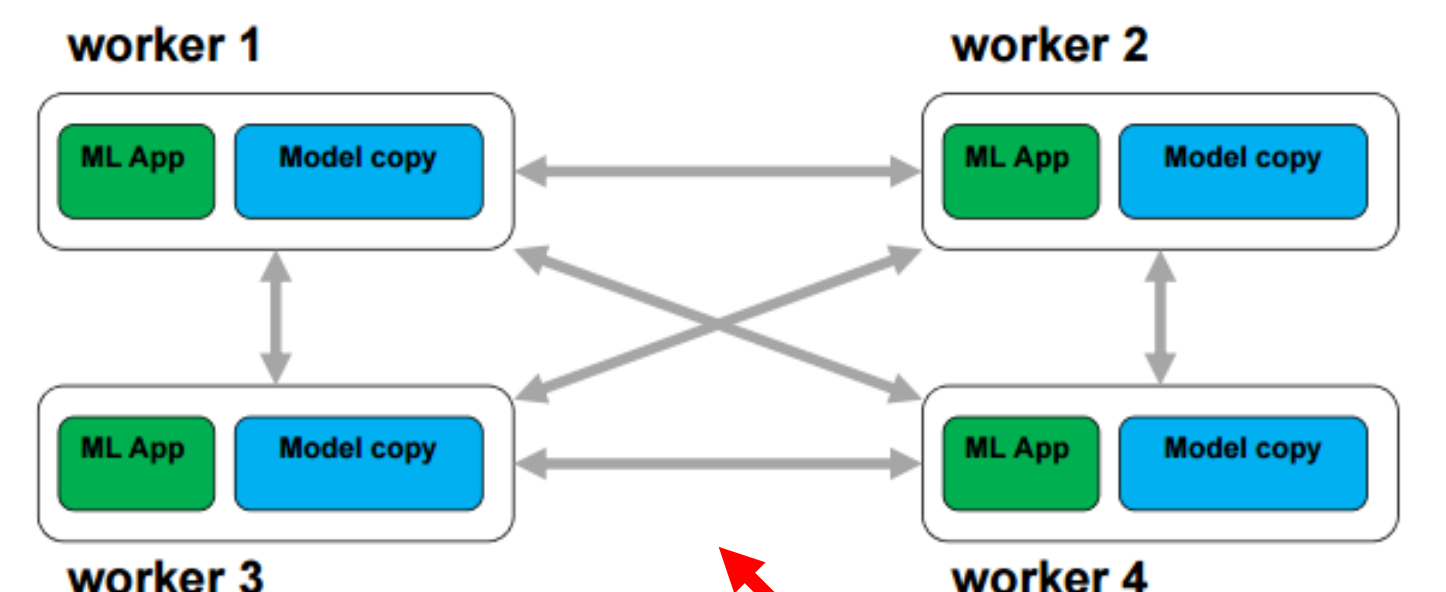
Structure-Aware Communication Protocol (SACP)



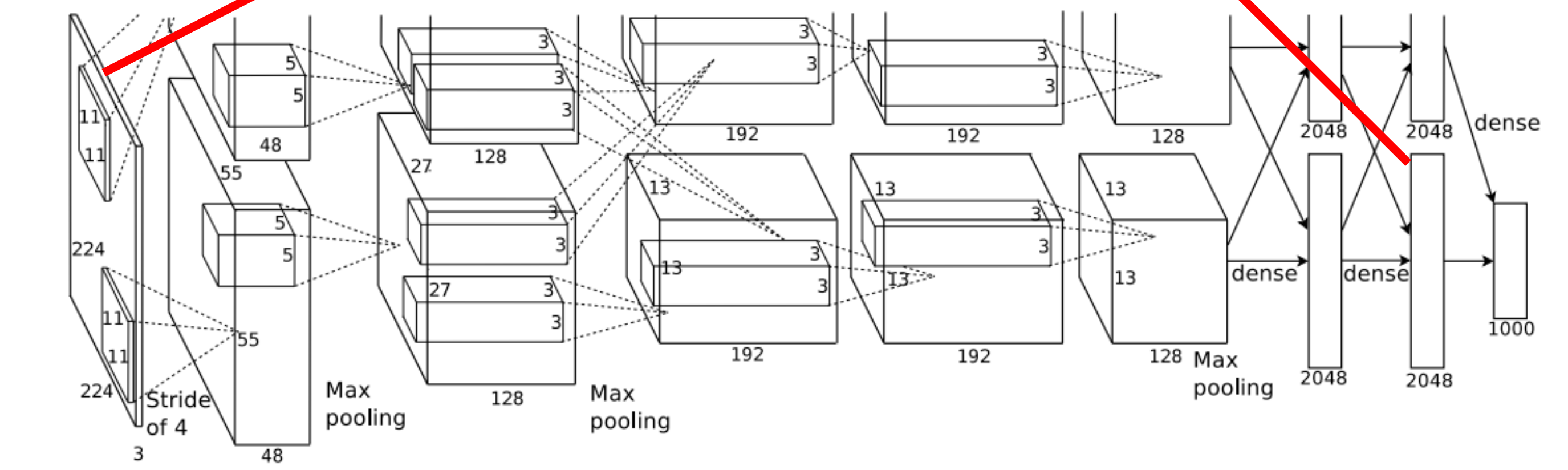
Topology: Master-Slave



Topology: Peer-to-Peer (P2P)

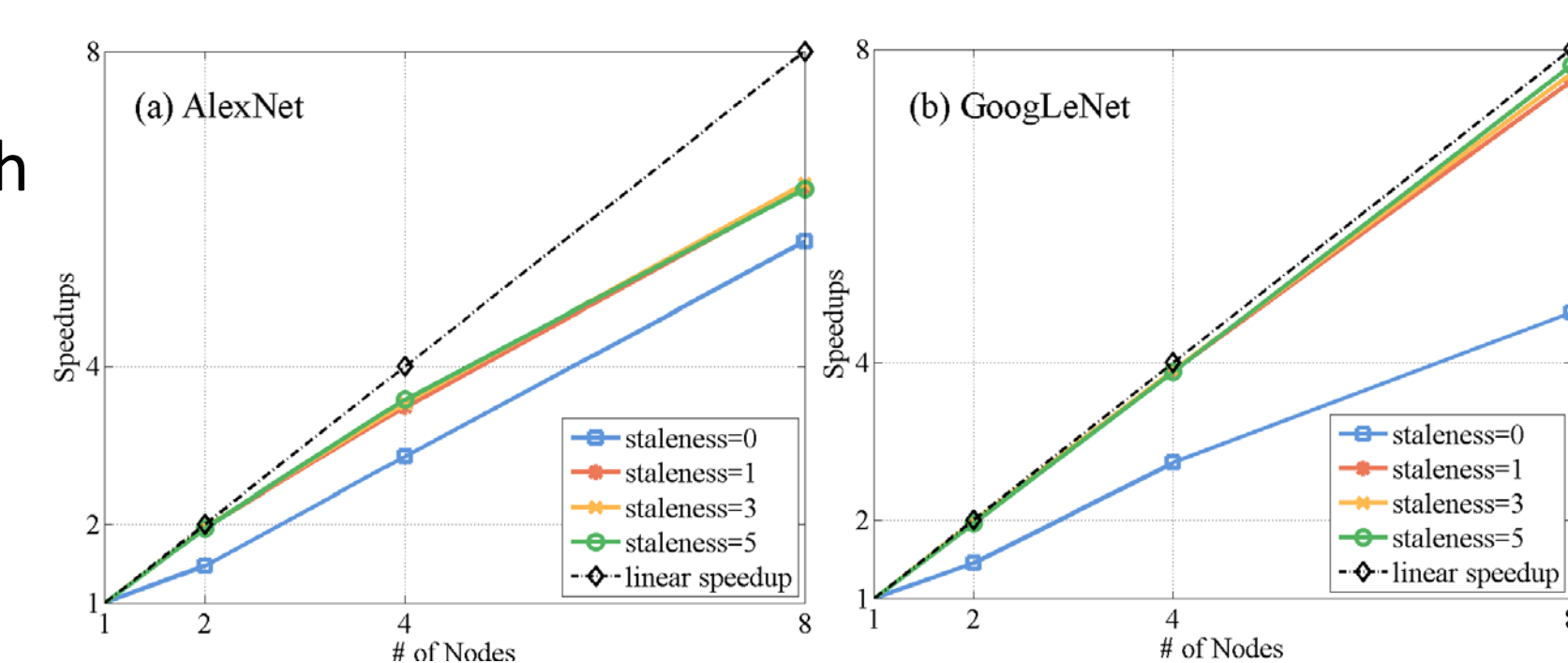


Small matrices



Big matrices

The speedups on throughput with different values of staleness, (a) Training AlexNet with batch size 256, and (b) Training GoogLeNet with batch size 32.



Training AlexNet and GoogLeNet with different number of GPU nodes and settings: (a) AlexNet with batch size 256 ; (b) GoogLeNet with batch size 32 . Compared to single machine Caffe

