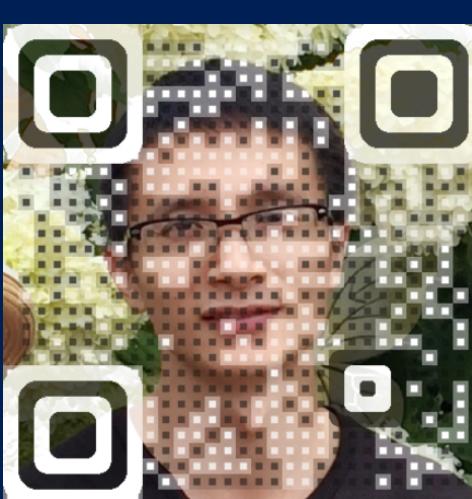


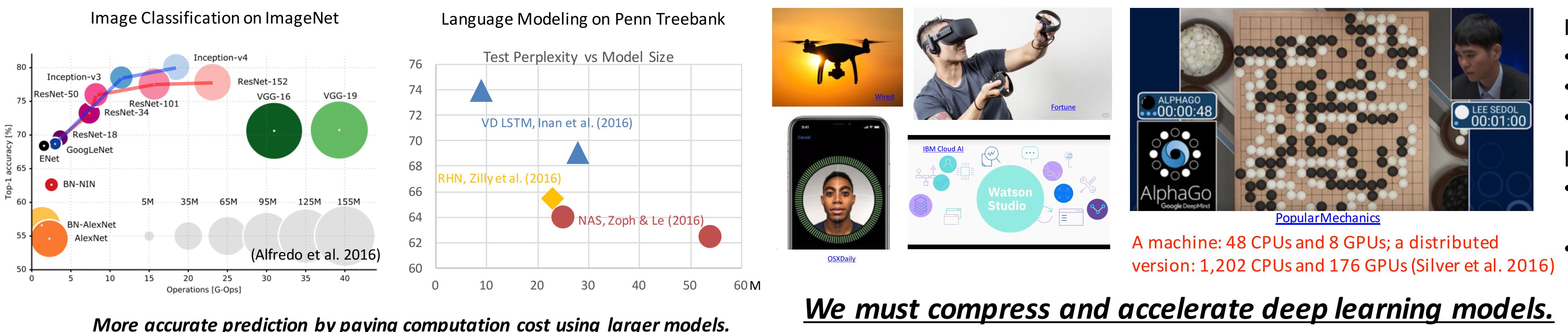
Efficient and Scalable Deep Learning

Wei Wen, Yiran Chen, Hai (Helen) Li, Duke University

{wei.wen, yiran.chen, hai.li}@duke.edu



Efficient Deep Learning for Faster Inference on the Edge and in the Cloud

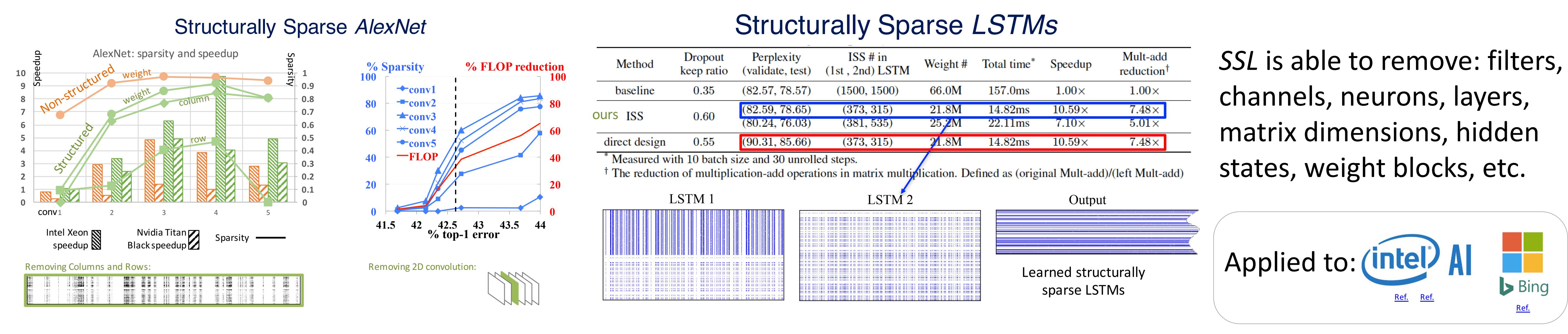
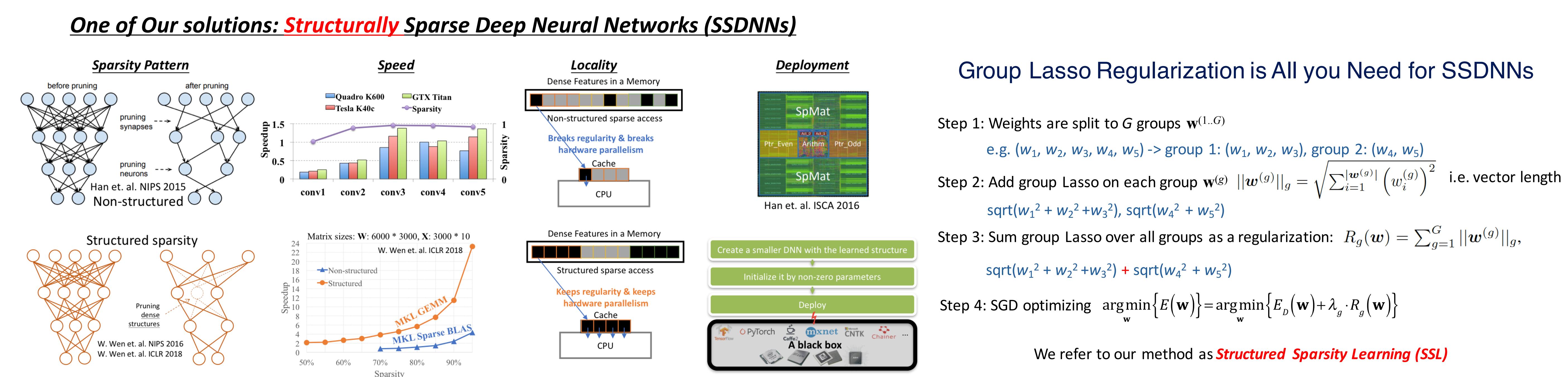


Inference on the edge:

- Limited computing capability
- Limited memory
- Limited battery energy

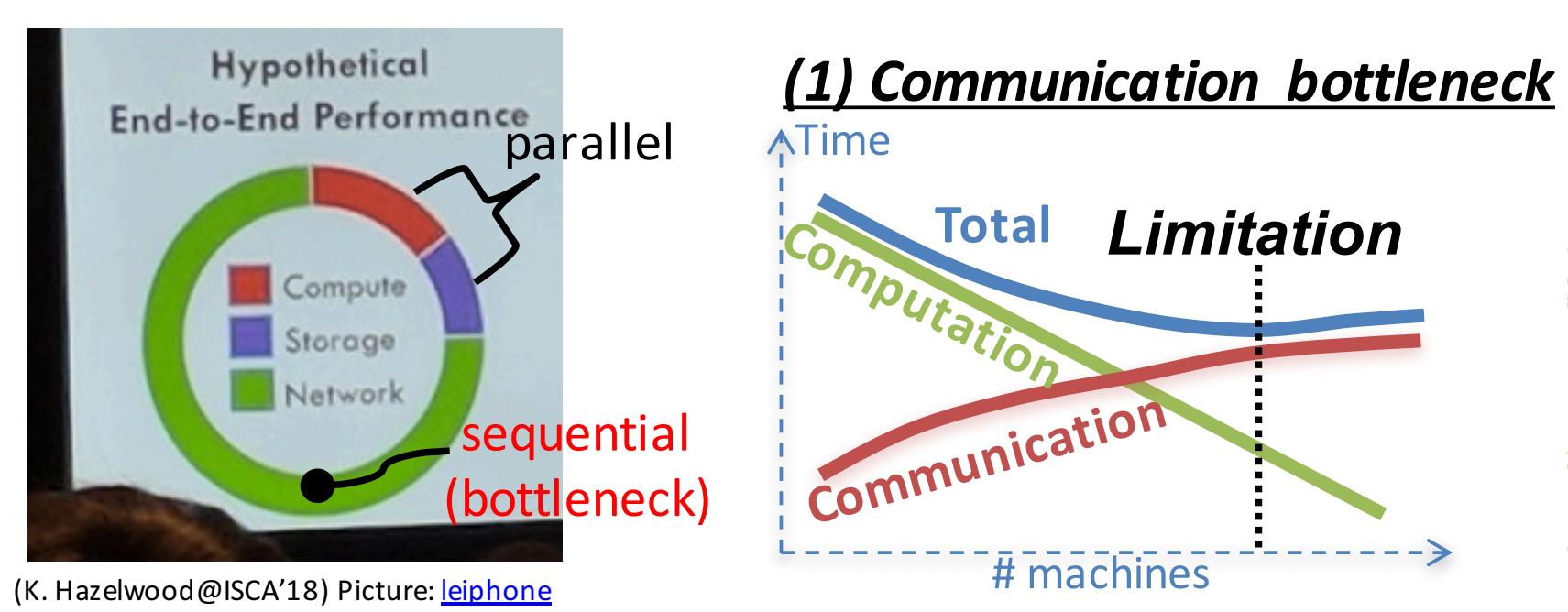
Inference in the cloud:

- A challenge: Real-time response to trillions of AI service requests
- Facebook (K. Hazelwood@ISCA'18)
 - 200+ Trillion predictions per day
 - 5+ Billion language translations per day



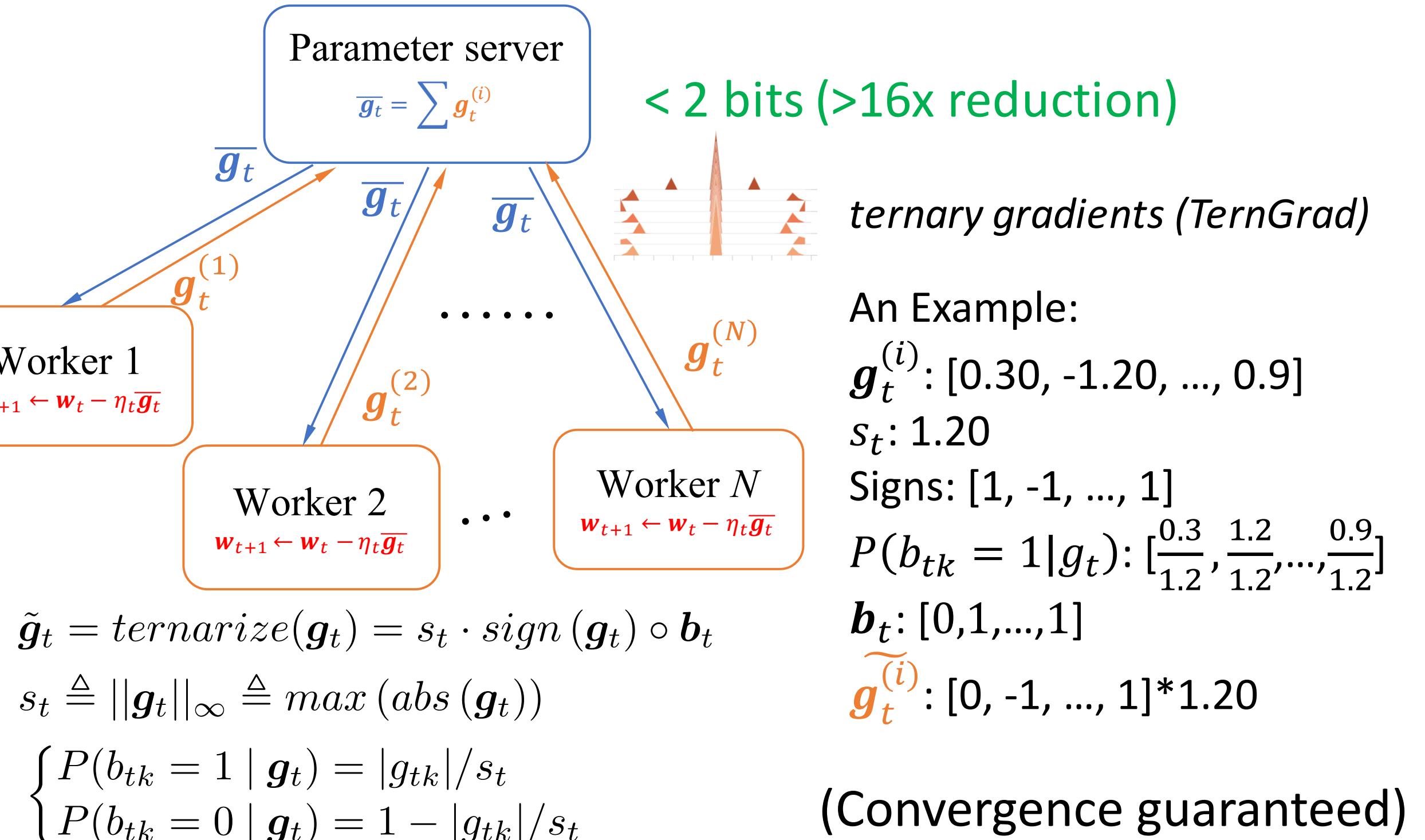
Our works on Efficient Deep Learning: NIPS 2016, ICLR 2018, ICCV 2017, ICLR 2017, CVPR 2017, ASP-DAC 2017 (Best Paper Award), DAC 2015 & 2016 (Two Best Paper Nominations).

Scalable Deep Learning for Faster Training in Distributed Systems



Solution to (1)

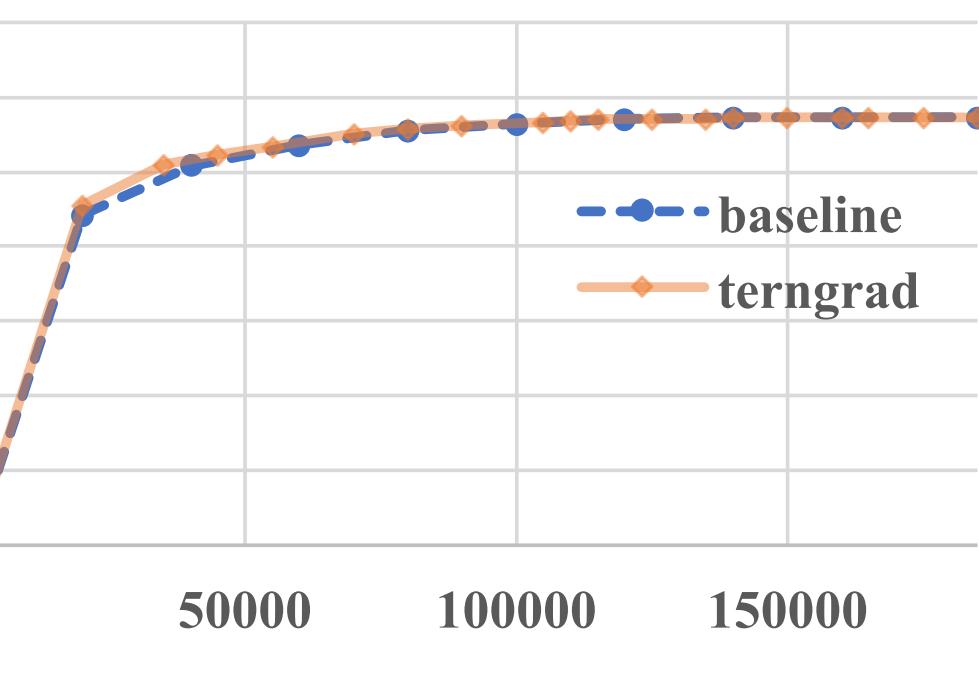
TernGrad: quantizing gradients to reduce communication (NIPS 2017, Oral)



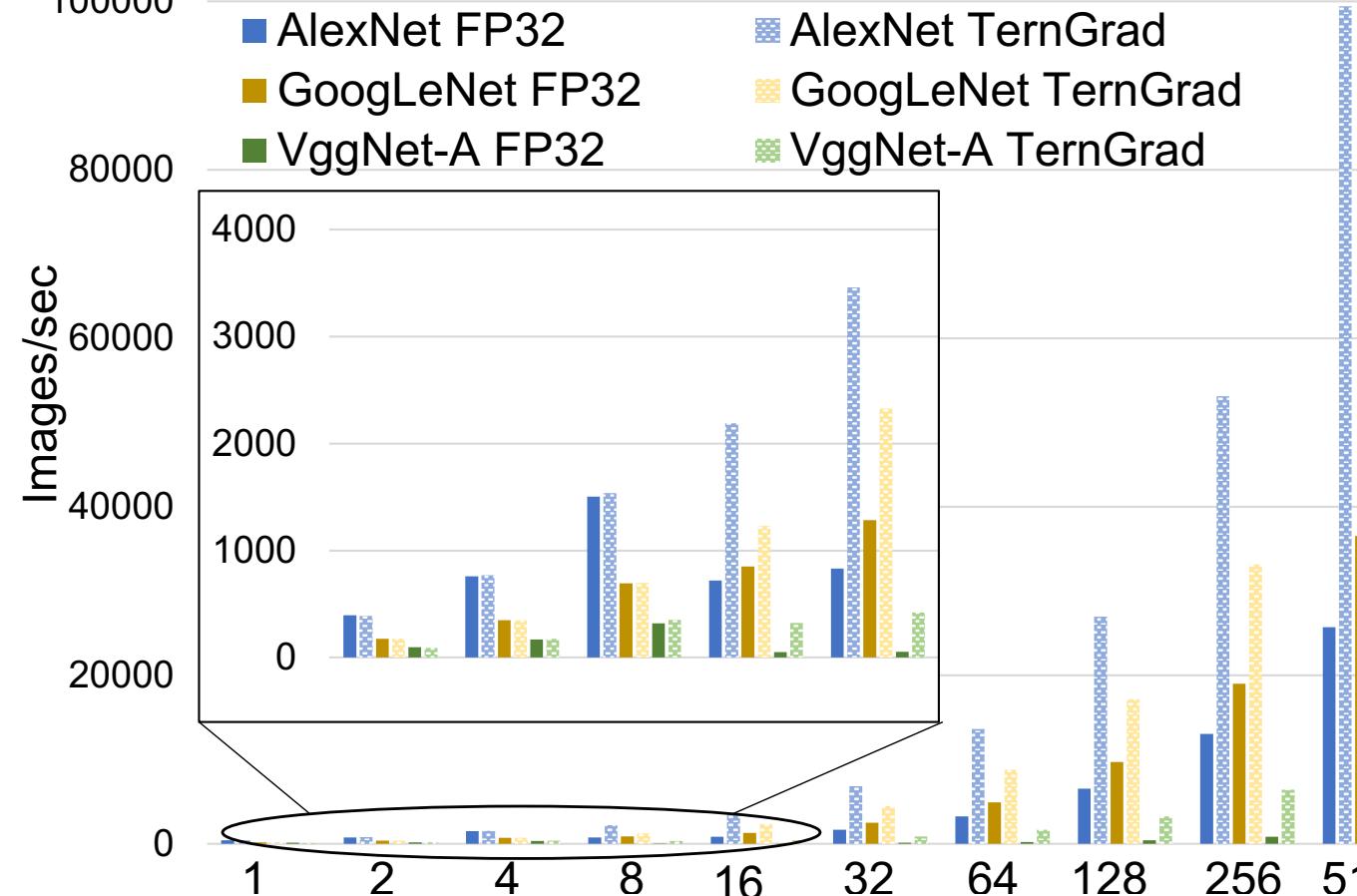
TernGrad Summary:

- No loss in AlexNet;
- <2% loss in GoogLeNet

(a) top-1 accuracy vs iteration



Training throughput on GPU cluster with Ethernet and PCI switch



TernGrad is now in PyTorch/Caffe2 and Adopted in Facebook AI Infra.



Solution to (2)

W. Wen, Y. Wang, F. Yan, C. Xu, Y. Chen, H. Li, “SmoothOut: Smoothing Out Sharp Minima to Improve Generalization in Deep Learning”, AAAI 2018 Submission.

