

Bandwidth-Efficient Deep Learning — Challenges and Trade-offs

Song Han

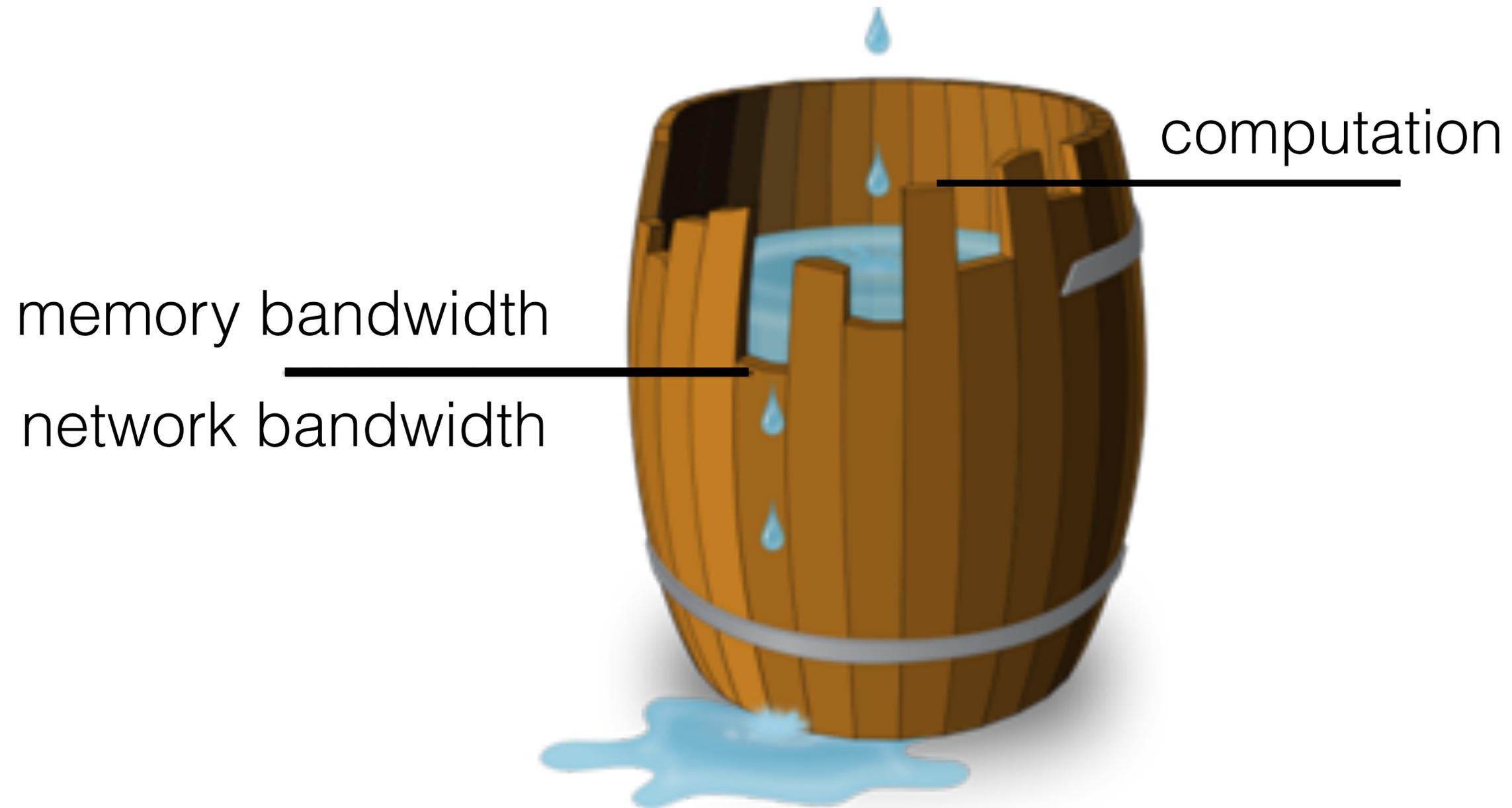
Assistant Professor, MIT EECS Department*

77 Massachusetts Avenue, Room 38-344
Cambridge, MA, 02139
songhan@mit.edu

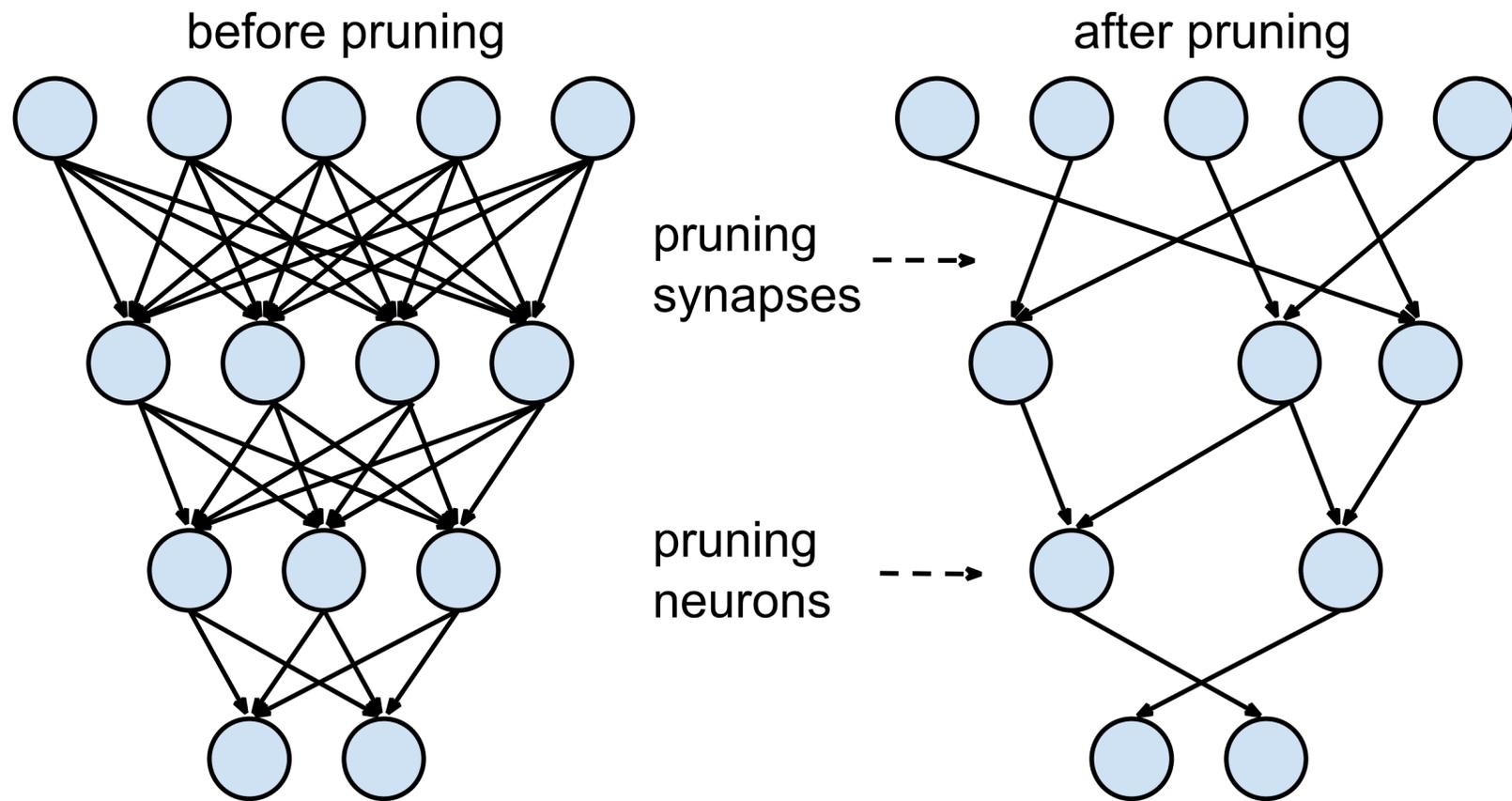
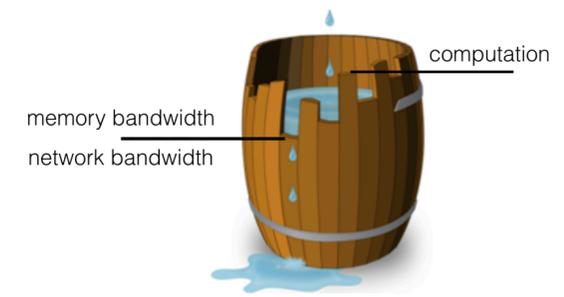
Deep Learning Hardware: Challenges and Trade-offs

- Training vs. Inference
- Latency vs. Throughput
- Edge vs. Cloud
- Accuracy vs. Model Size
- Flexibility/Programmability vs. Efficiency
- Computation bounded vs. Memory bounded
- Time to market vs. Efficiency
- Mature ecosystem vs. Emerging technology

Data Movement is the Bottleneck

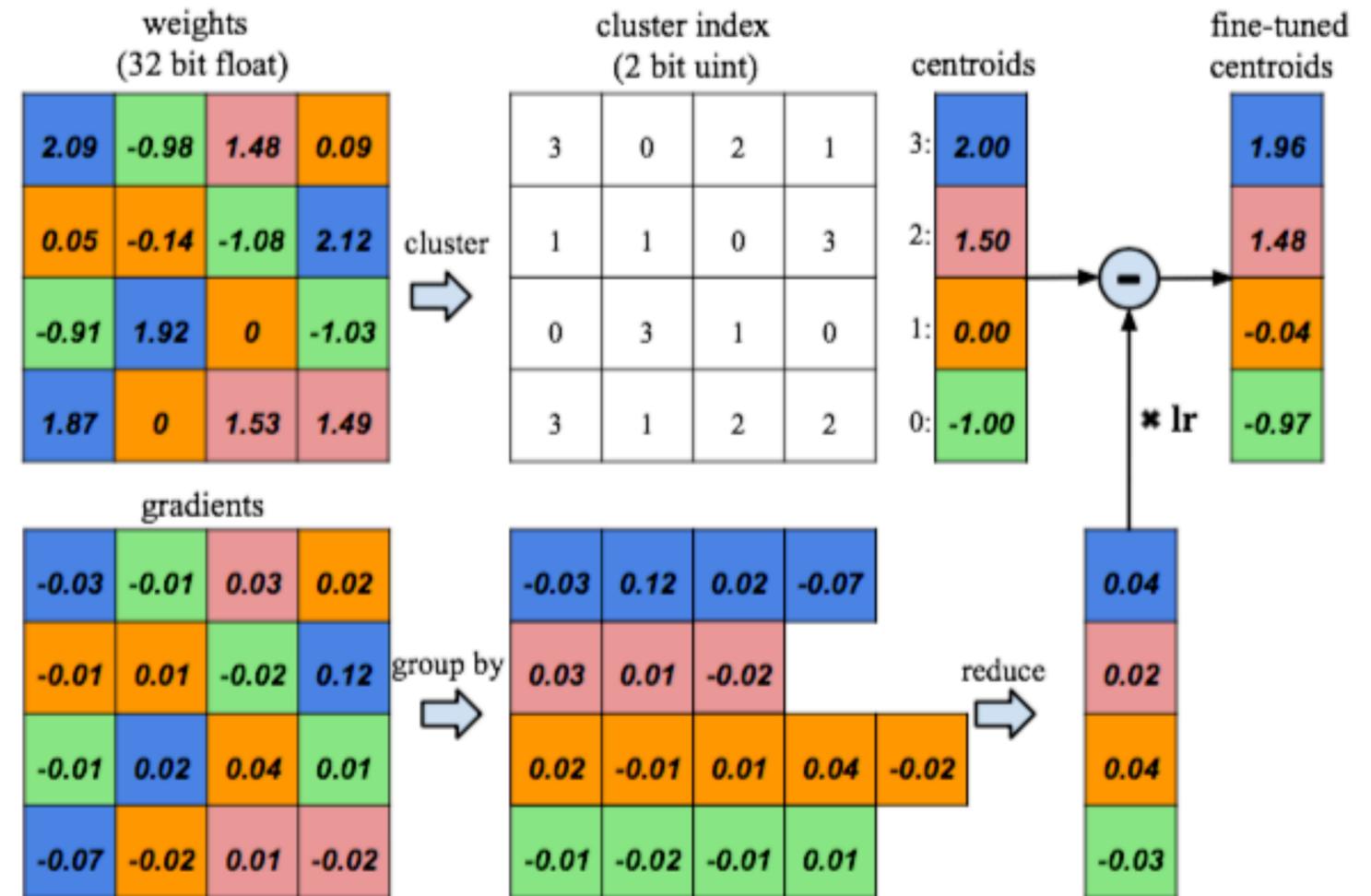


Save Memory Bandwidth: Deep Compression



Pruning

Han et al [NIPS'15]



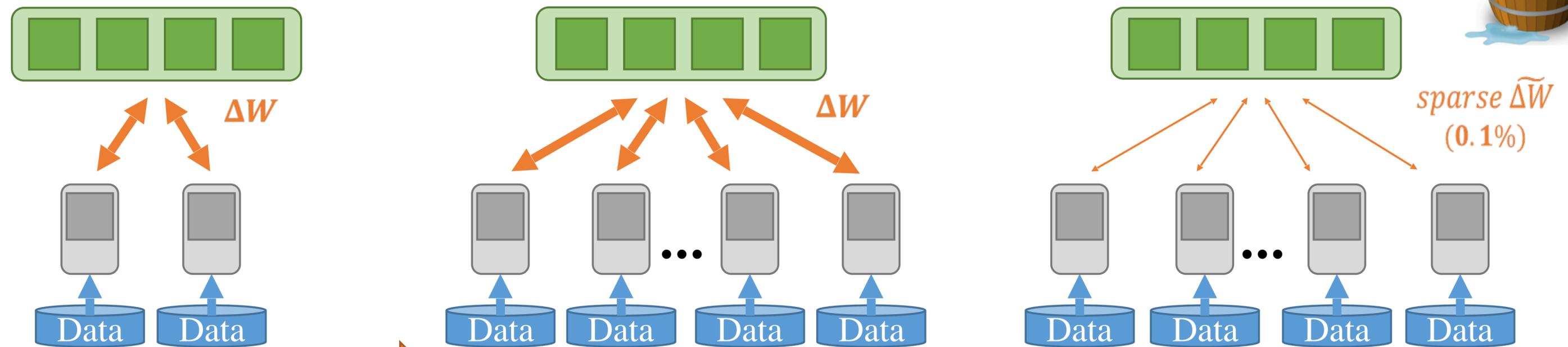
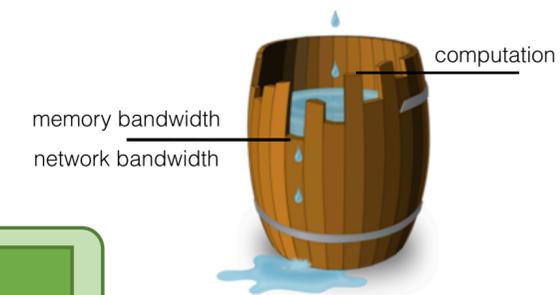
Trained Quantization

Han et al [ICLR'16]



Save Network Bandwidth: Deep Gradient Compression

saving the communication bandwidth of distributed training by 500x



More Training Nodes

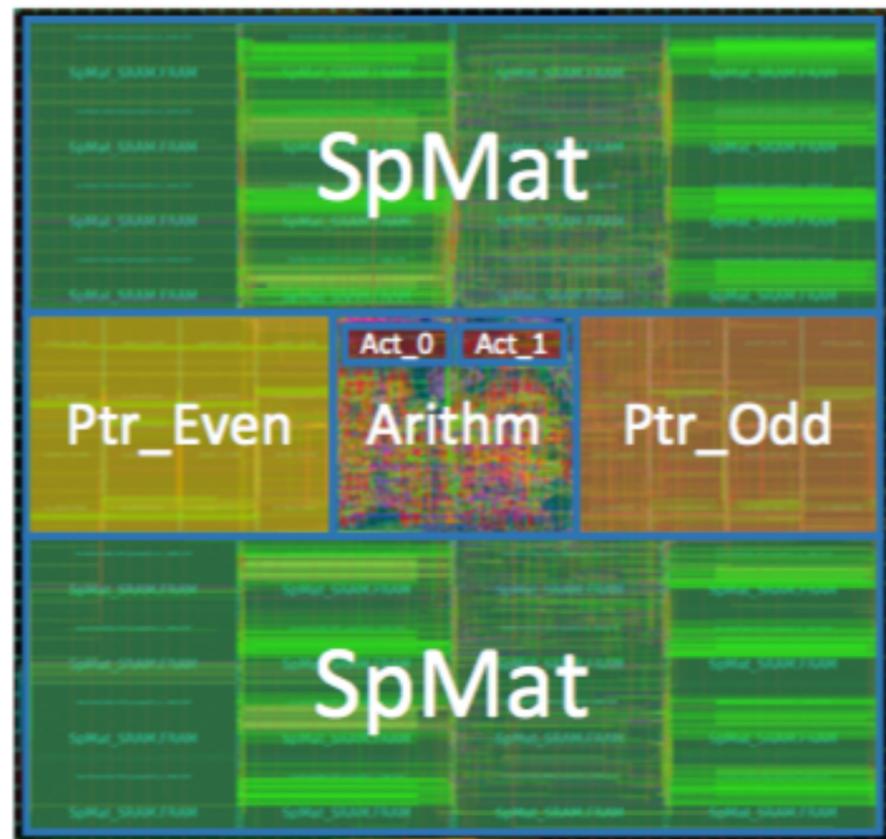
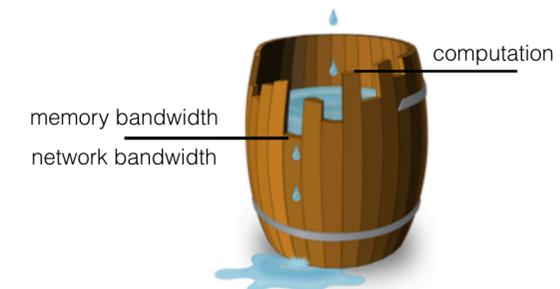
Deep Gradient Compression

Lin et al [ICLR'18]

Time:

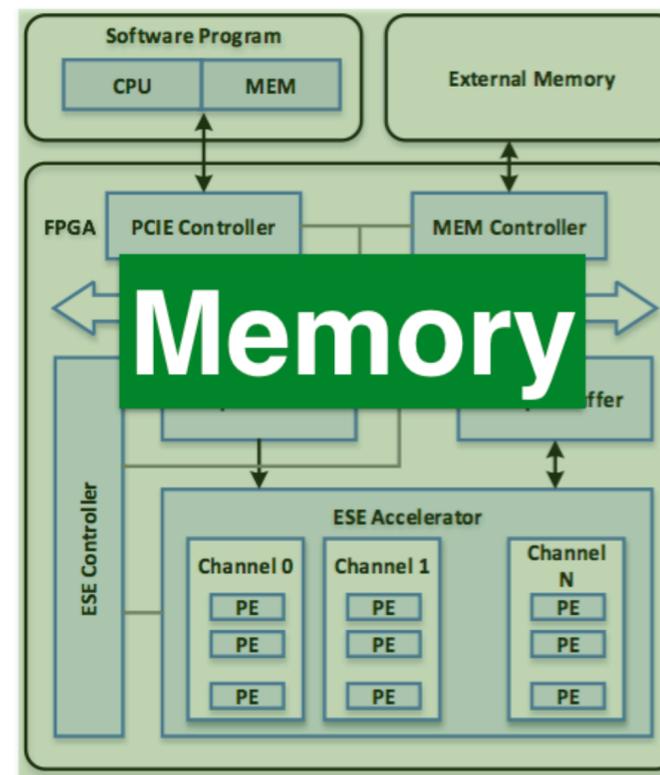


Efficient Inference Engine on Compressed DNN



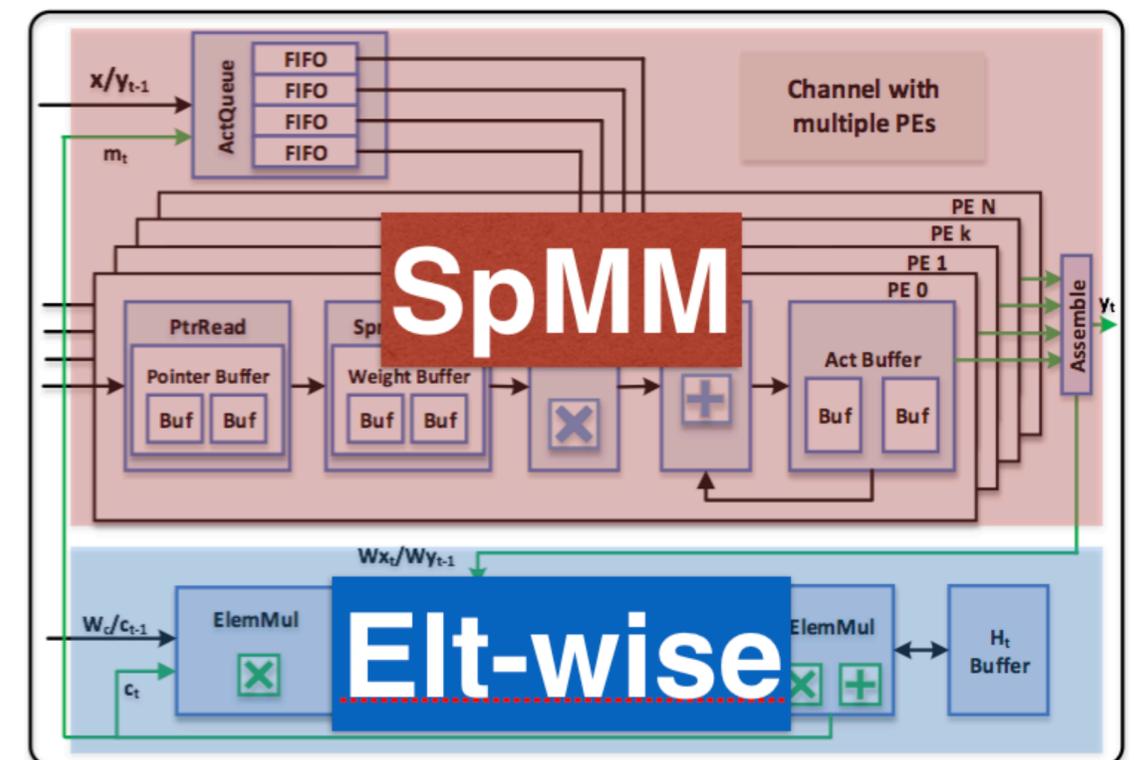
EIE Accelerator

Han et al [ISCA'16]



ESE Accelerator

Han et al [FPGA'17]



MIT Intelligent Hardware Lab

I'm founding the MIT Intelligent Hardware Lab starting in July 2018.

My research areas include:

- **H**: High performance, High energy efficiency Hardware
- **A**: Architectures and Accelerators for Artificial Intelligence
- **N**: Novel algorithms for Neural Networks and Deep Learning

Looking forward to industry/academic collaborations and talented PhD students!

