

High-Performance Hardware for Machine Learning

U.C. Berkeley
October 19, 2016

William Dally
NVIDIA Corporation
Stanford University

Machine learning is transforming computing

Speech

Natural Language Understanding

Question Answering

Game Playing (Go)

Vision

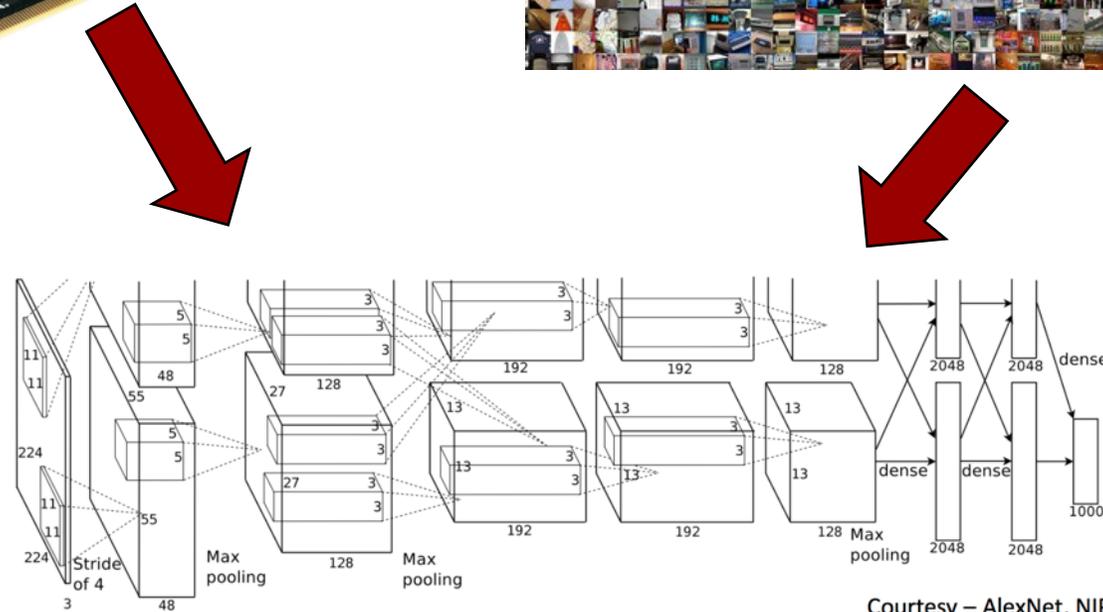
Autonomous Vehicles

Control

Ad Placement

Whole research fields rendered irrelevant

Hardware and Data enable DNNs



Courtesy – AlexNet, NIPS 2012

The Need for Speed

Important Property of Neural Networks

Results get better with

**more data +
bigger models +
more computation**

(Better algorithms, new insights and improved techniques always help, too!)



IMAGE RECOGNITION

16X
Model



2012
AlexNet



2015
ResNet



SPEECH RECOGNITION

10X
Training Ops



2014
Deep Speech 1

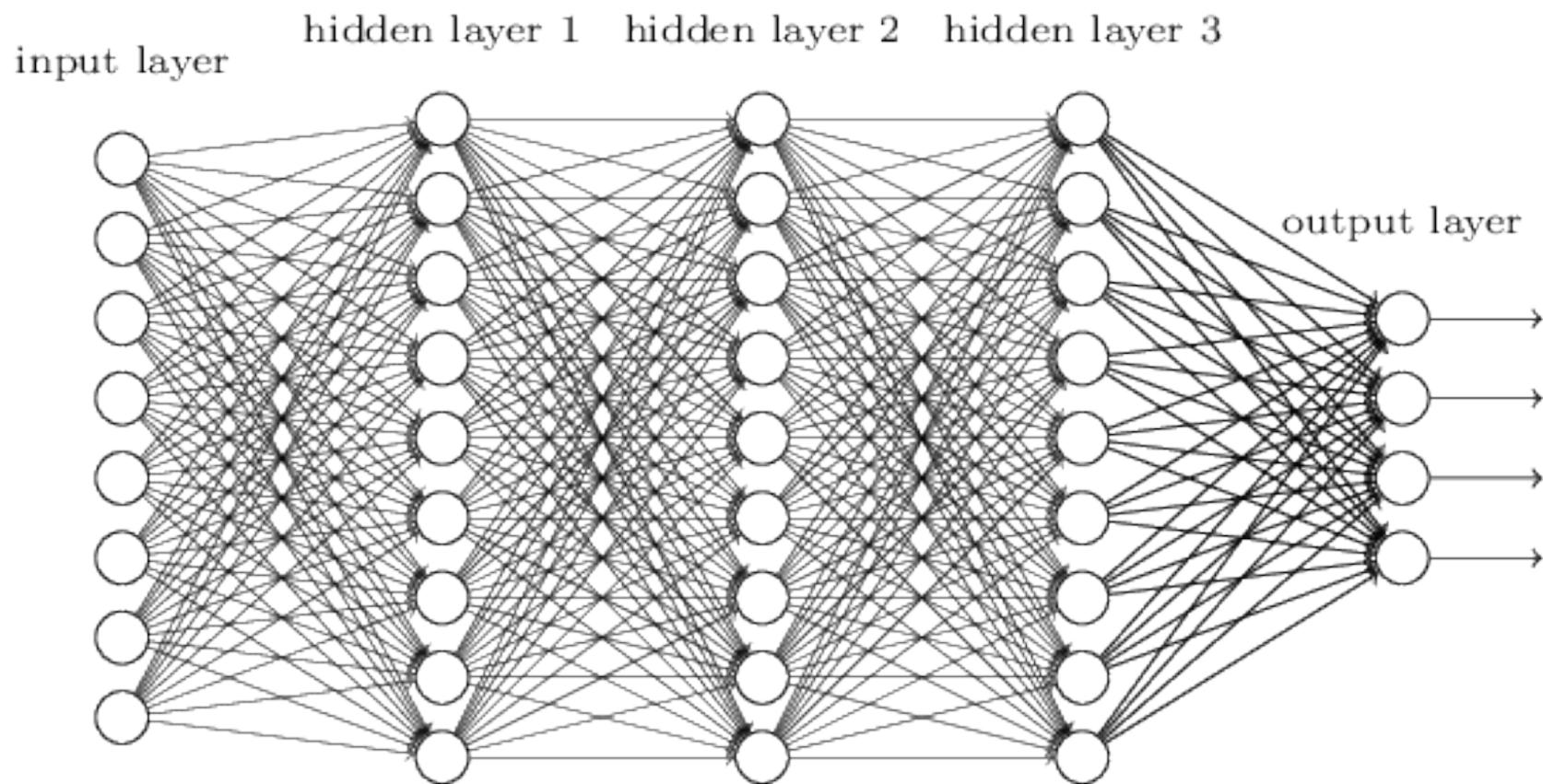


2015
Deep Speech 2

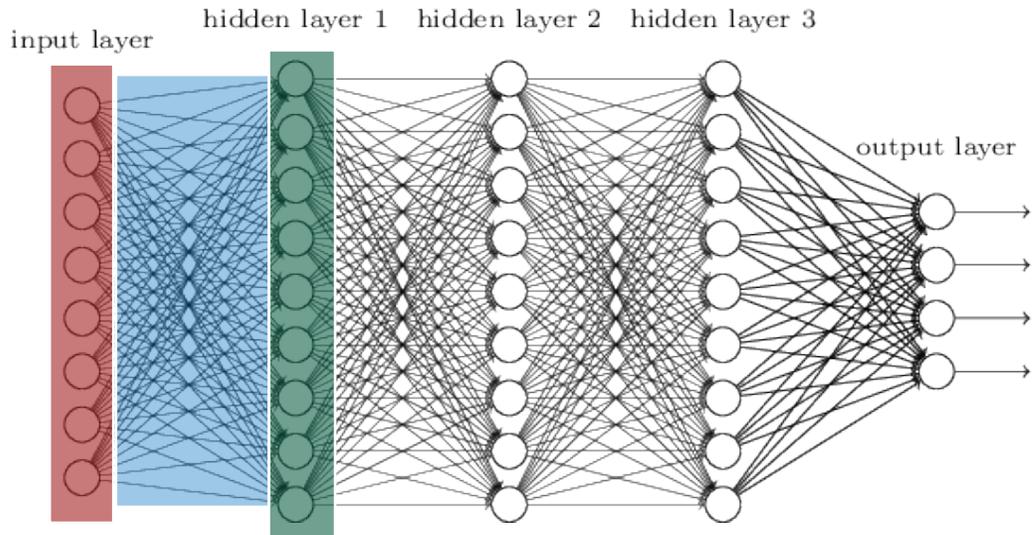


DNN primer

WHAT NETWORK? DNNS, CNNS, AND RNNS



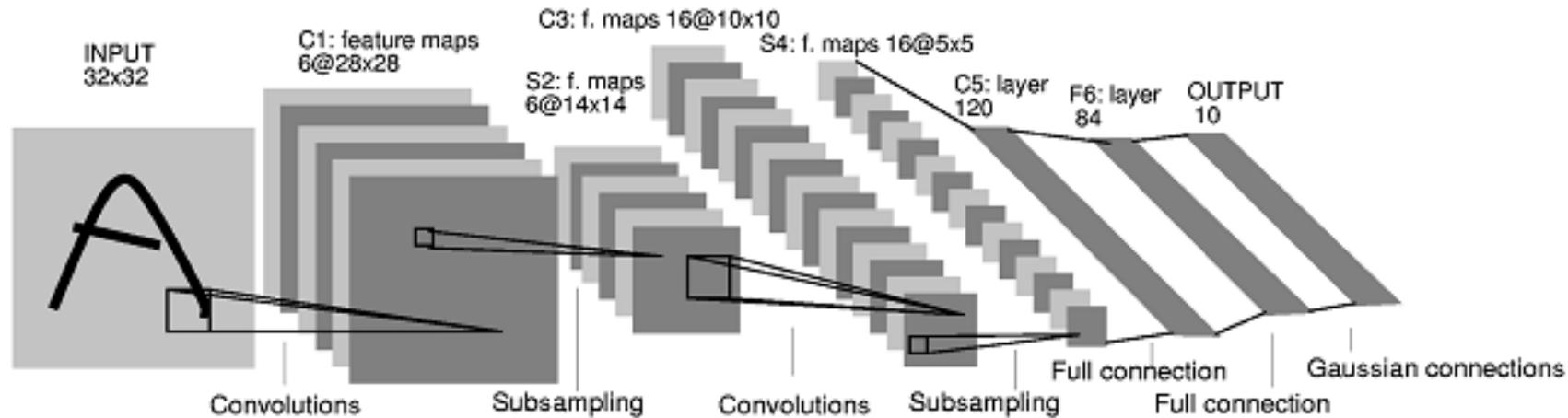
DNN, KEY OPERATION IS DENSE $M \times V$



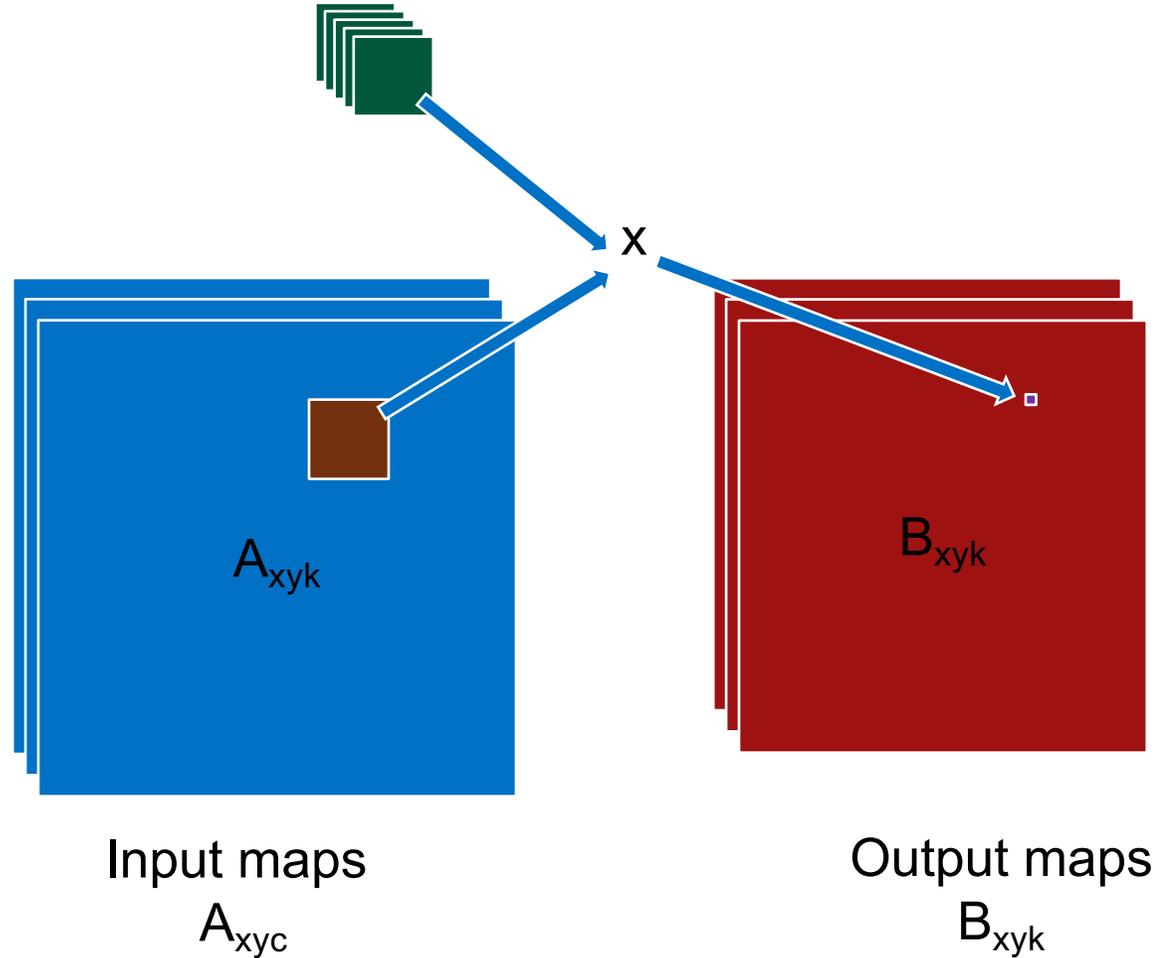
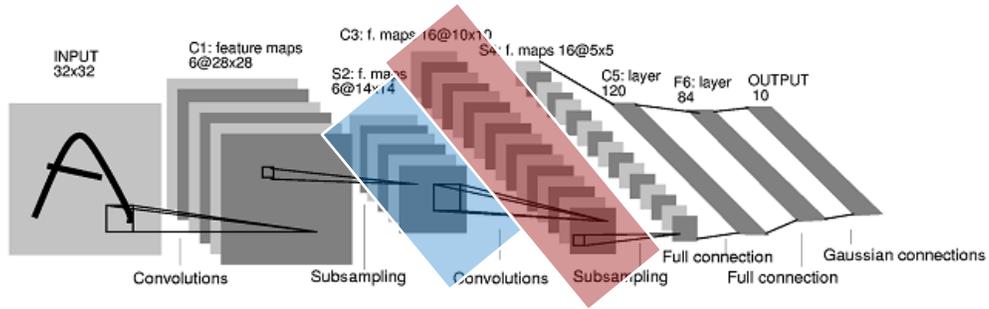
$$b_i = W_{ij} \times a_j$$

The equation shows a dark green vertical bar labeled b_i on the left, followed by an equals sign, a large blue square labeled W_{ij} in the center, followed by a multiplication sign, and a dark red vertical bar labeled a_j on the right.

CNNS - For image inputs, convolutional stages act as trained feature detectors



CNNs require convolution in addition to $M \times V$

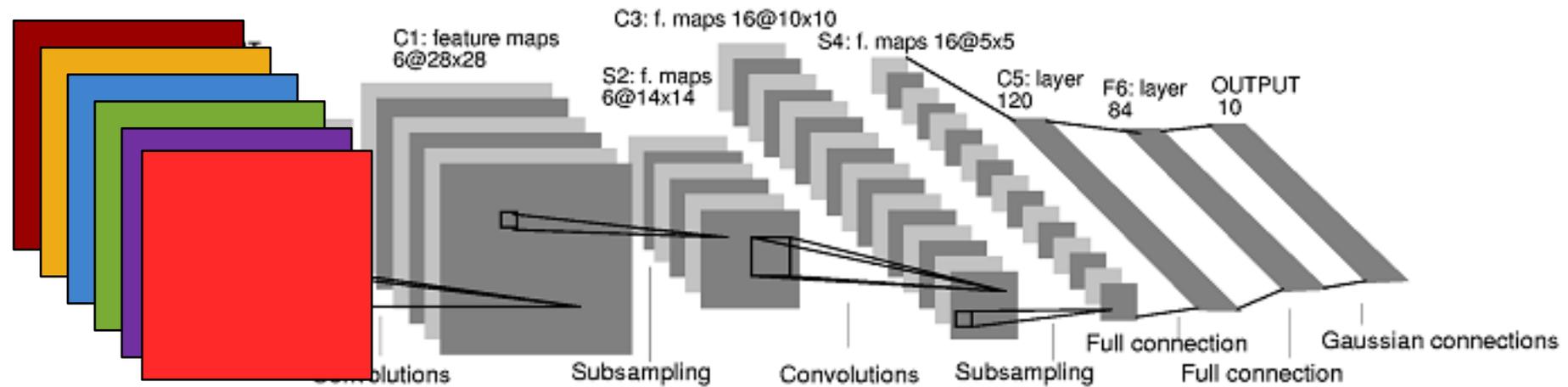


4 Distinct Sub-problems

	Training	Inference	
Convolutional	Train Conv	Inference Conv	B x S Weight Reuse Act Dominated
Fully-Conn.	Train FC	Inference FC	B Weight Reuse Weight Dominated
	32b FP - large batches Large Memory Footprint Minimize Training Time	8b Int - small (unit) batches Meet real-time constraint	

DNNs are Trivially Parallelized

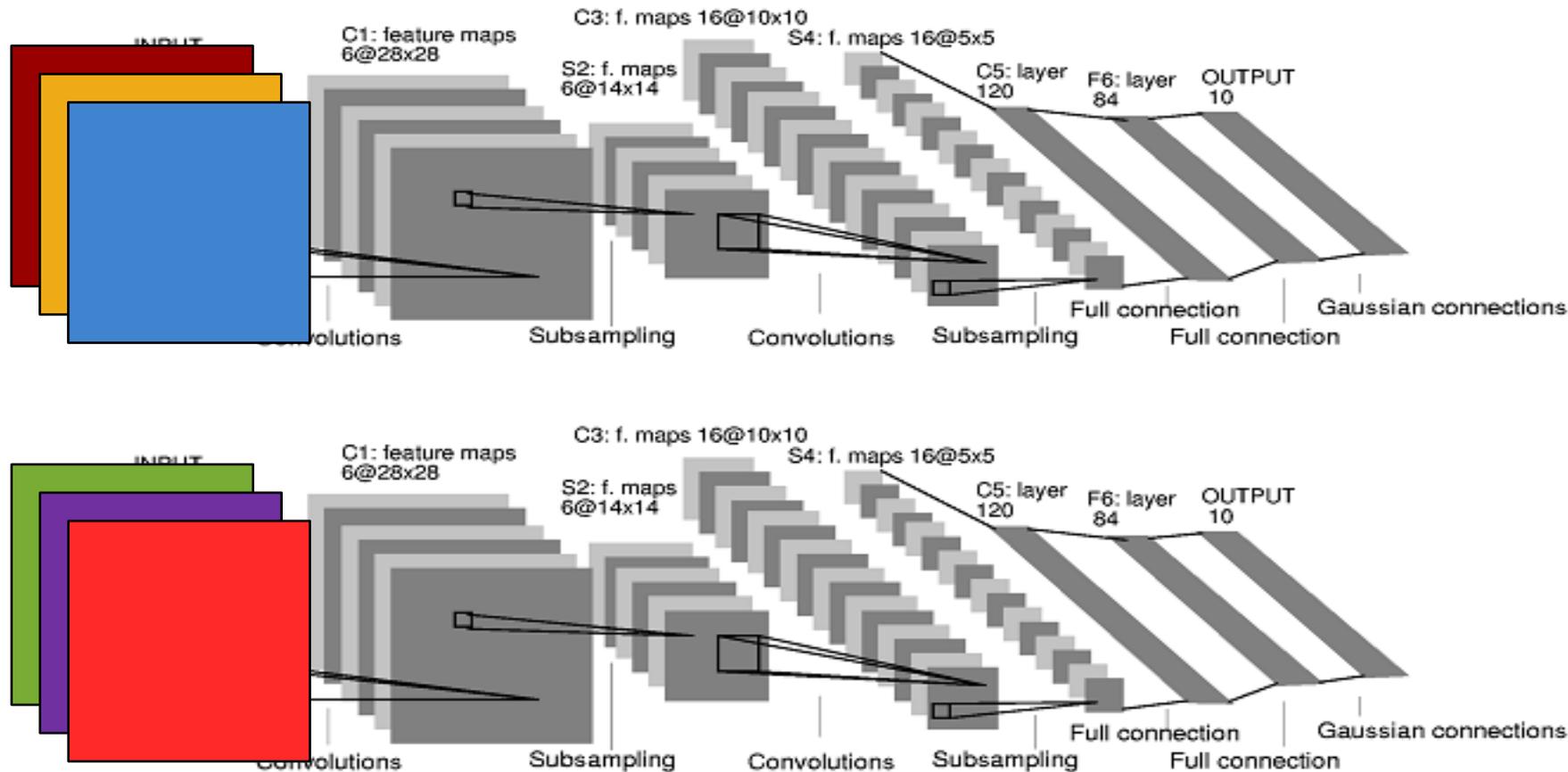
Lots of parallelism in a DNN



- Inputs
- Points of a feature map
- Filters
- Elements within a filter

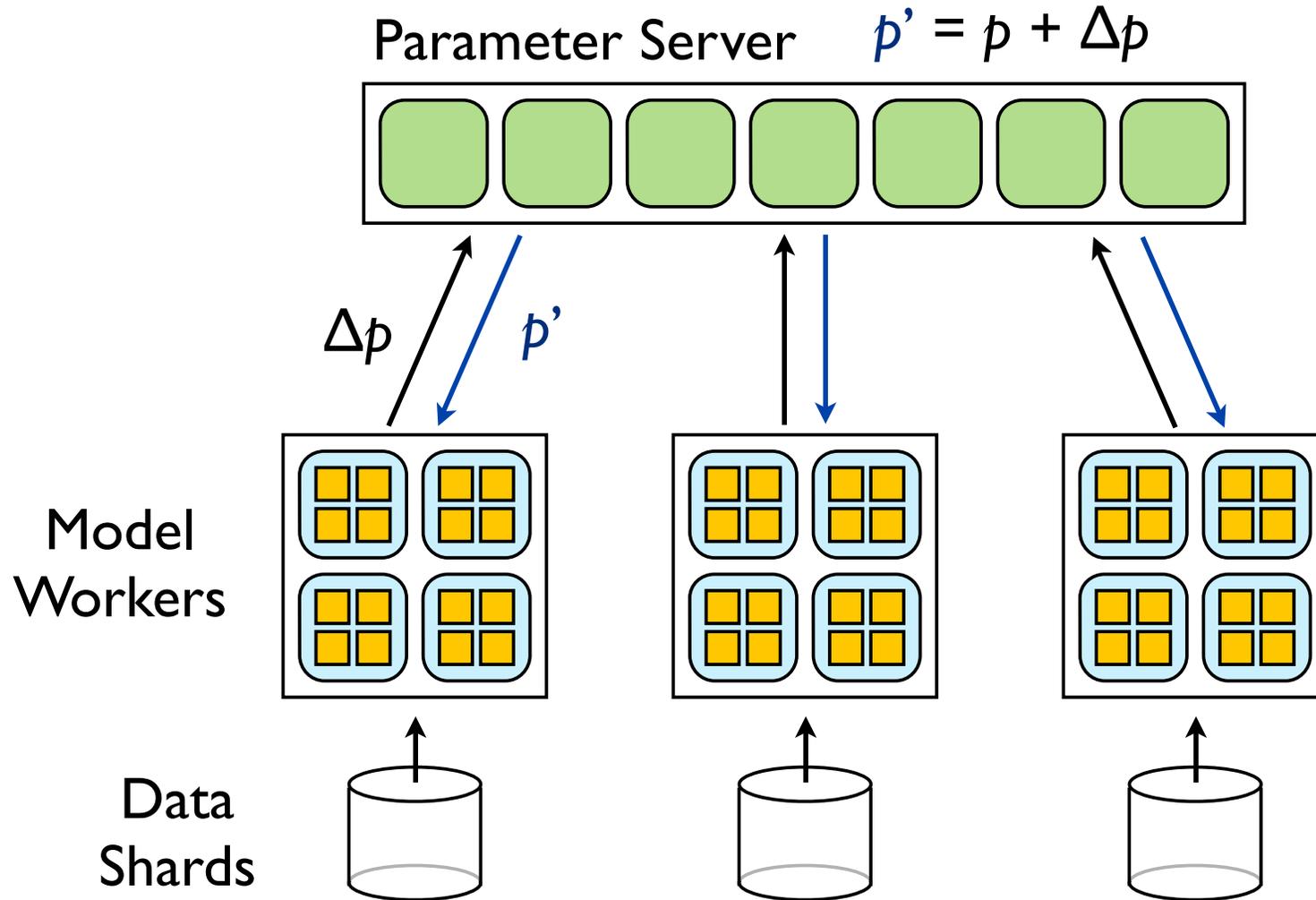
- Multiplies within layer are independent
- Sums are reductions
- Only layers are dependent
- No data dependent operations
=> can be statically scheduled

Data Parallel – Run multiple inputs in parallel

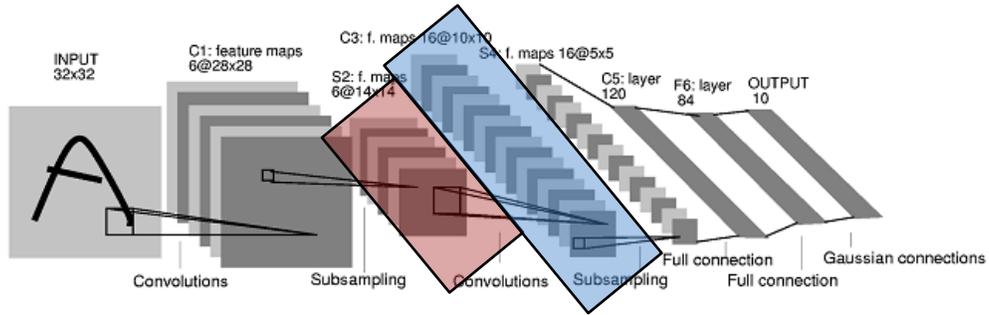


- Doesn't affect latency for one input
- Requires P-fold larger batch size
- For training requires coordinated weight update

Parameter Update



Model-Parallel Convolution – by output region (x,y)



Kernels
 Multiple 3D
 K_{uvkj}



x

6D Loop

For all region XY

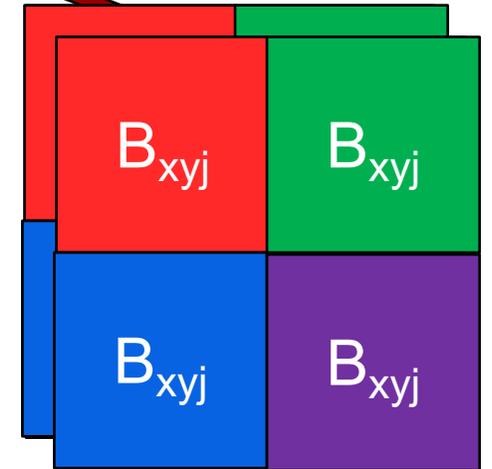
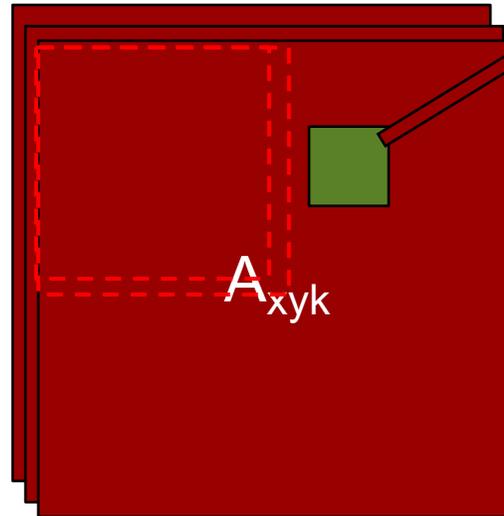
For each output map j

For each input map k

For each pixel x,y in XY

For each kernel element u,v

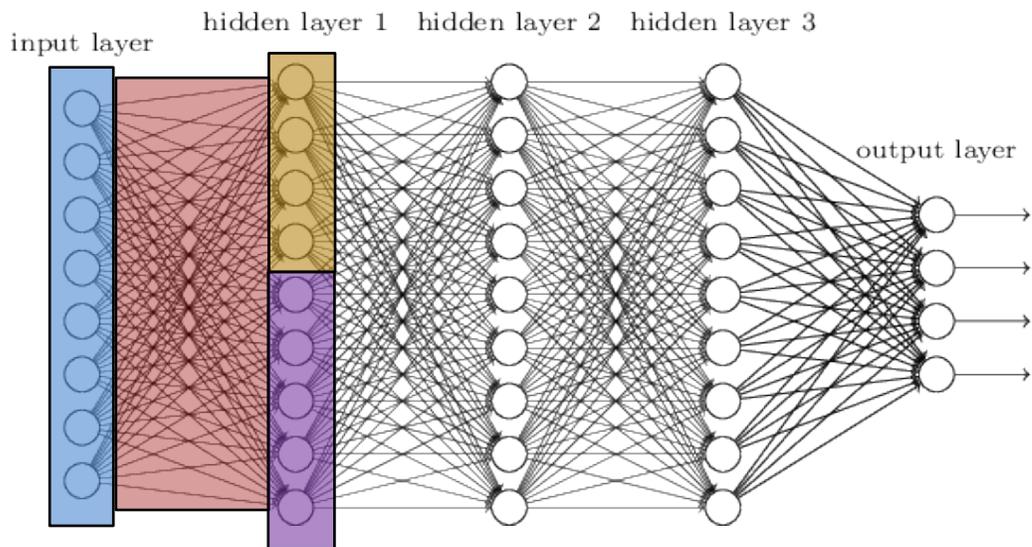
$$B_{xyj} += A_{(x-u)(y-v)k} \times K_{uvkj}$$



Input maps
 A_{xyk}

Output maps
 B_{xyj}

Model Parallel Fully-Connected Layer ($M \times V$)

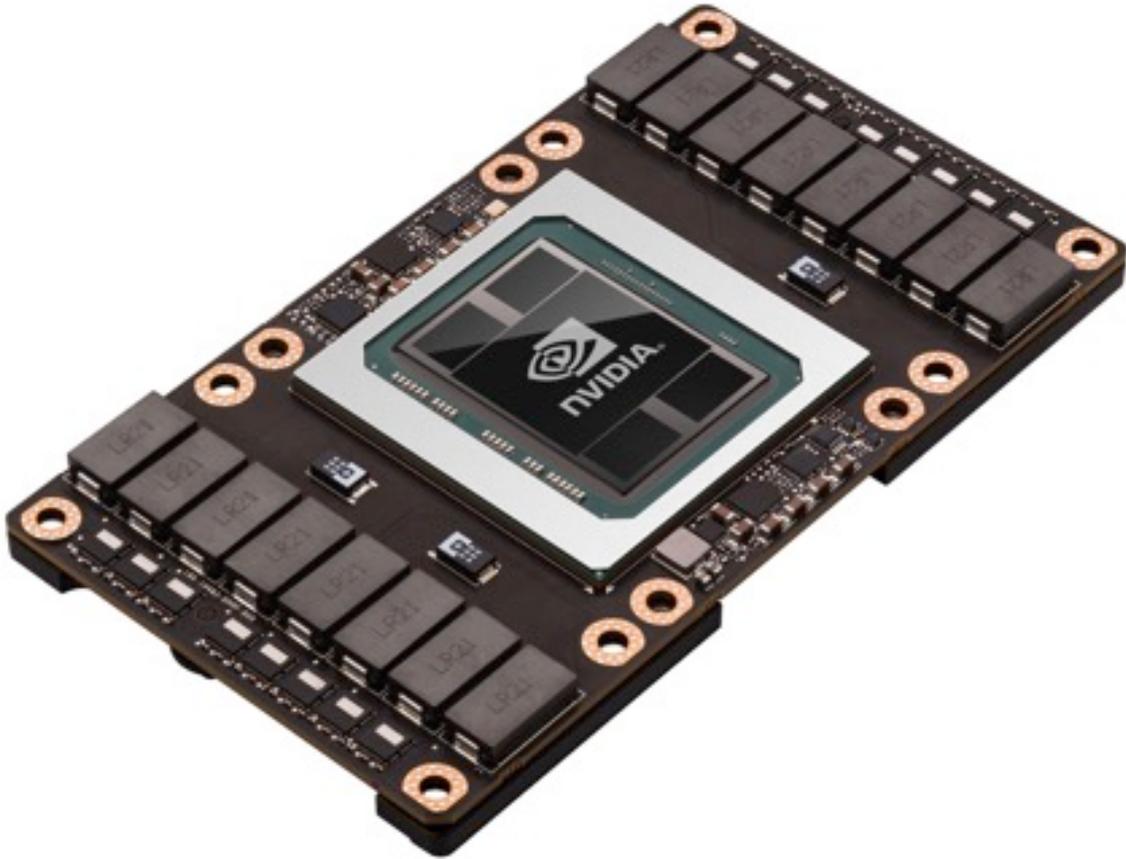


$$\begin{array}{c} \text{Output activations} \\ \begin{array}{c} b_i \\ b_i \end{array} \end{array} = \begin{array}{c} \text{weight matrix} \\ \begin{array}{c} W_{ij} \\ W_{ij} \end{array} \end{array} \times \begin{array}{c} \text{Input activations} \\ a_j \end{array}$$

The equation shows the relationship between the output activations, the weight matrix, and the input activations. The output activations are represented by a vertical stack of two colored boxes: a yellow box labeled b_i on top and a purple box labeled b_i on the bottom. The weight matrix is a large square divided into two horizontal sections: a dark red top section labeled W_{ij} and a red bottom section labeled W_{ij} . The input activations are represented by a blue vertical box labeled a_j . The entire equation is preceded and followed by equals signs and a multiplication sign.

GPUs

Pascal GP100



- 10 TeraFLOPS FP32
- 20 TeraFLOPS FP16
- 16GB HBM – 750GB/s
- 300W TDP
- 67GFLOPS/W (FP16)
- 16nm process
- 160GB/s NV Link

NVIDIA DGX-1

WORLD'S FIRST DEEP LEARNING SUPERCOMPUTER



170 TFLOPS

8x Tesla P100 16GB

NVLink Hybrid Cube Mesh

Optimized Deep Learning
Software

Dual Xeon

7 TB SSD Deep Learning Cache

Dual 10GbE, Quad IB 100Gb

3RU - 3200W

Facebook's deep learning machine

- Purpose-Built for Deep Learning Training



2x Faster Training for Faster Deployment

2x Larger Networks for Higher Accuracy

Powered by Eight Tesla M40 GPUs

Open Rack Compliant

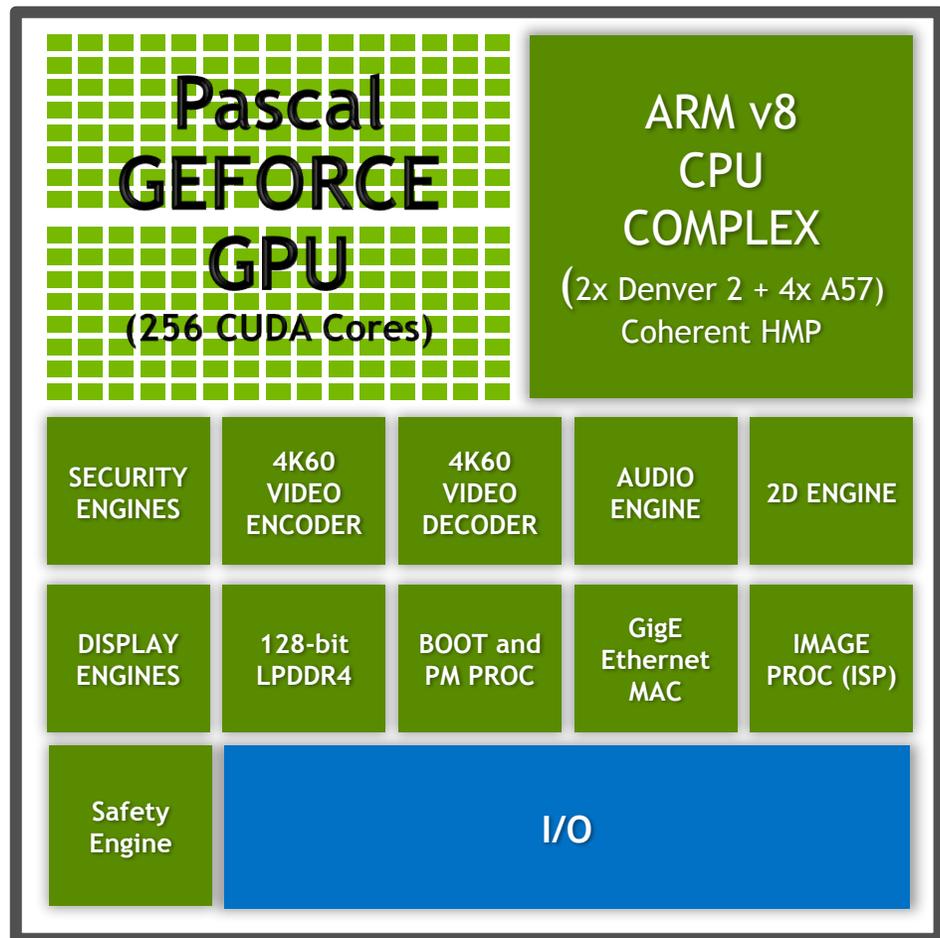
*“Most of the major advances in machine learning and AI in the past few years have been contingent on **tapping into powerful GPUs** and huge data sets to build and train advanced models”*



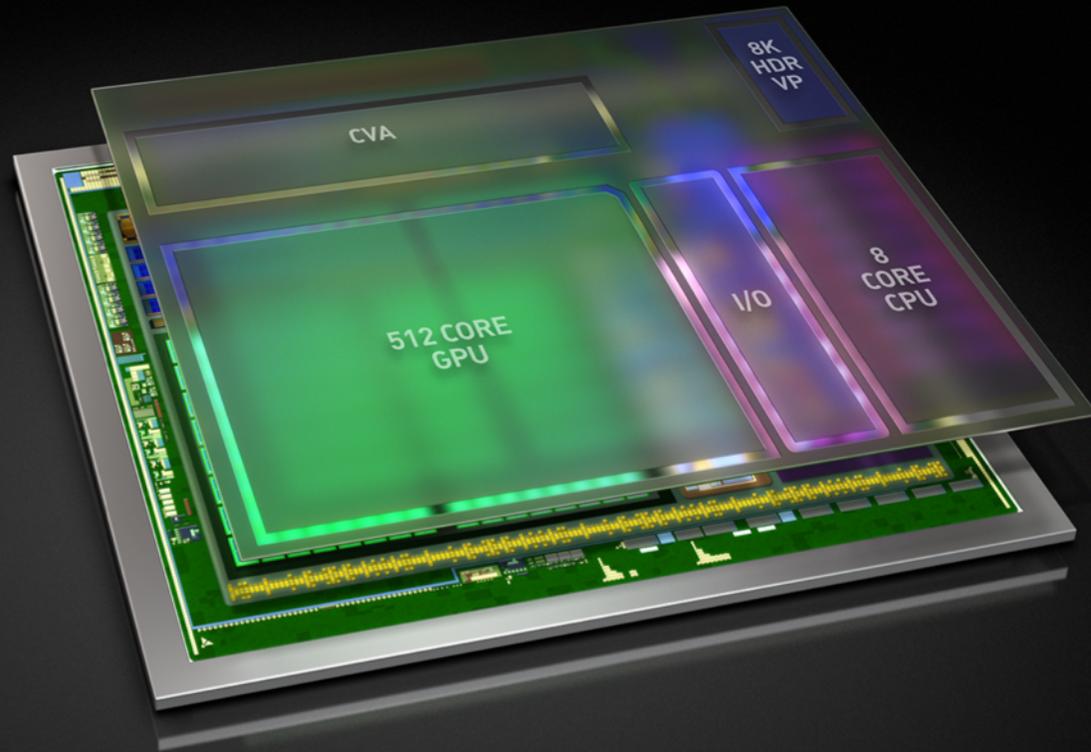
Serkan Piantino

Engineering Director of Facebook AI Research

NVIDIA Parker



- 1.5 Teraflop FP16
- 4GB of LPDDR4 @ 25.6 GB/s
- 15 W TDP (1W idle, <10W typical)
- 100GFLOPS/W (FP16)
- 16nm process



XAVIER

AI SUPERCOMPUTER SOC

7 Billion Transistors 16nm FF

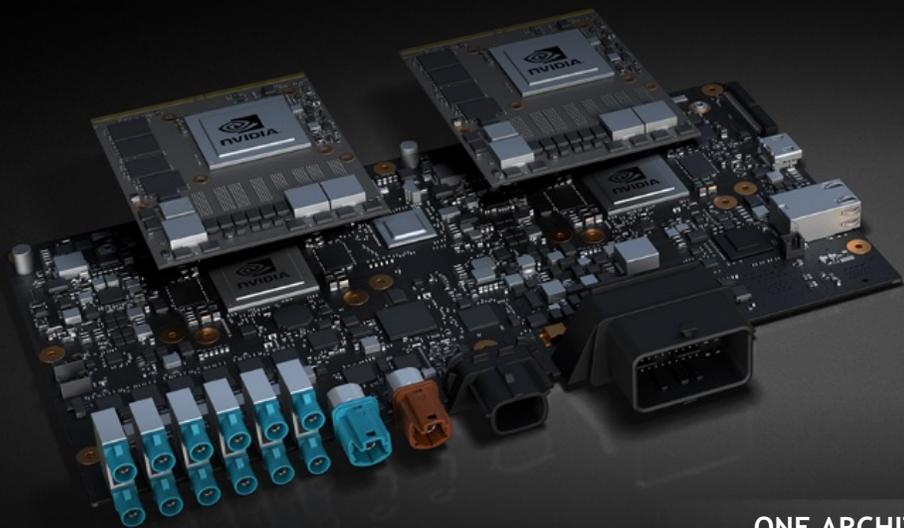
8 Core Custom ARM64 CPU

512 Core Volta GPU

New Computer Vision Accelerator

Dual 8K HDR Video Processors

Designed for ASIL C Functional Safety



DRIVE PX 2

2 PARKER + 2 PASCAL GPU | 20 TOPS DL | 120 SPECINT | 80W

ONE ARCHITECTURE



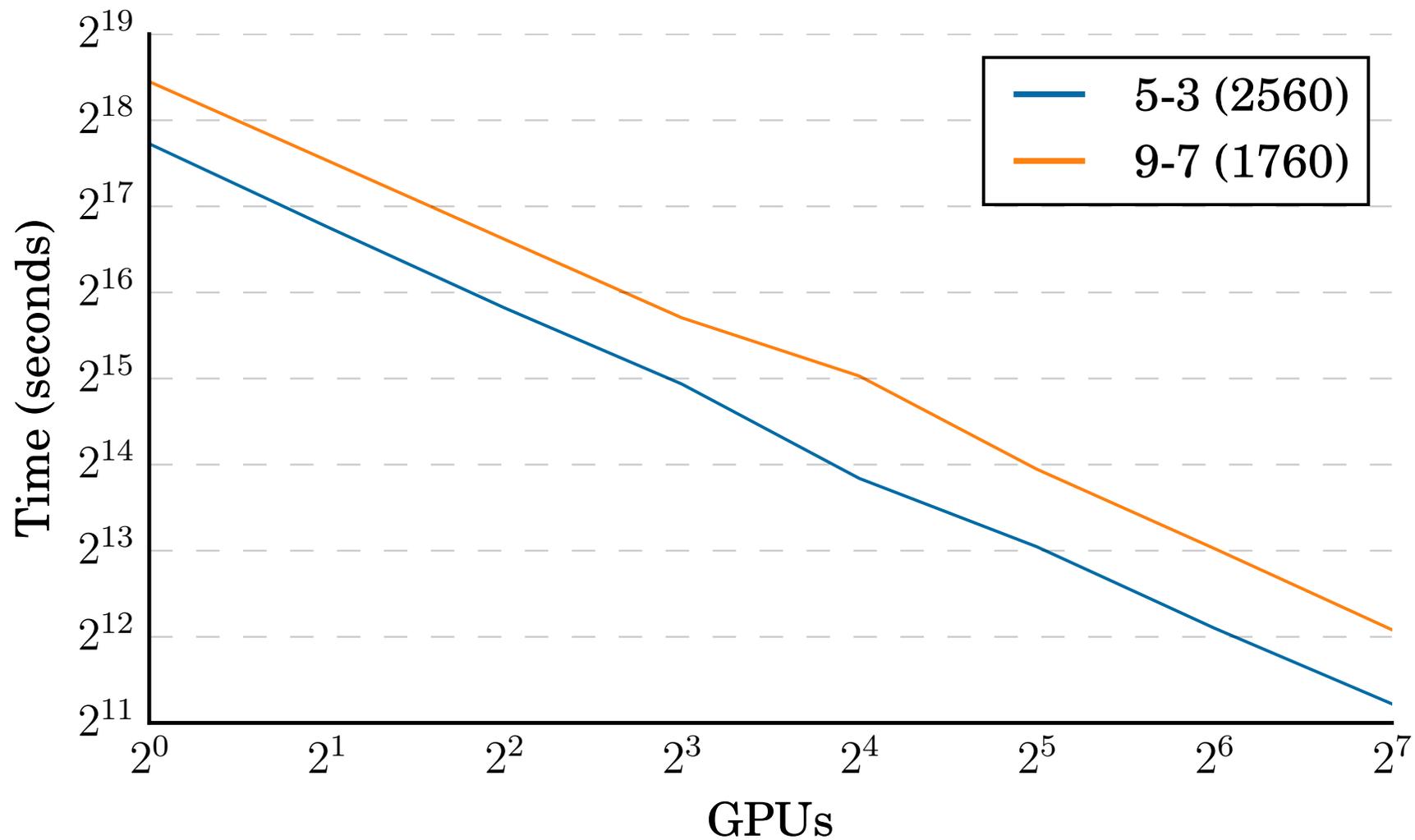
XAVIER

20 TOPS DL | 160 SPECINT | 20W

XAVIER

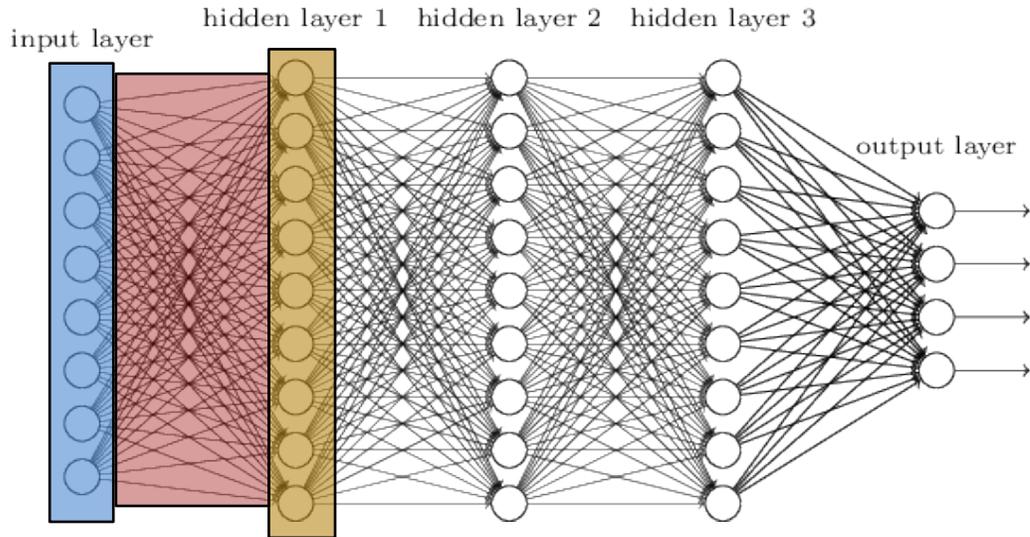
AI SUPERCOMPUTER SOC

Parallel GPUs on Deep Speech 2

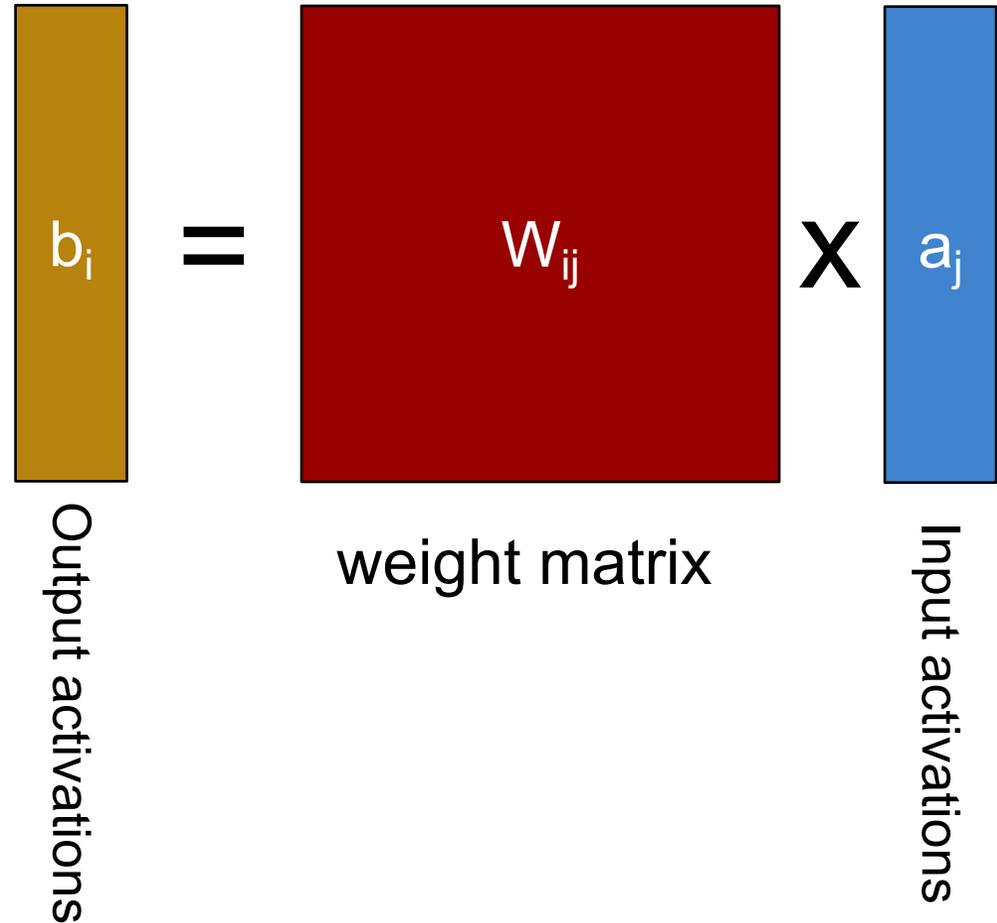


Reduced Precision

How Much Precision is Needed for Dense M x V?



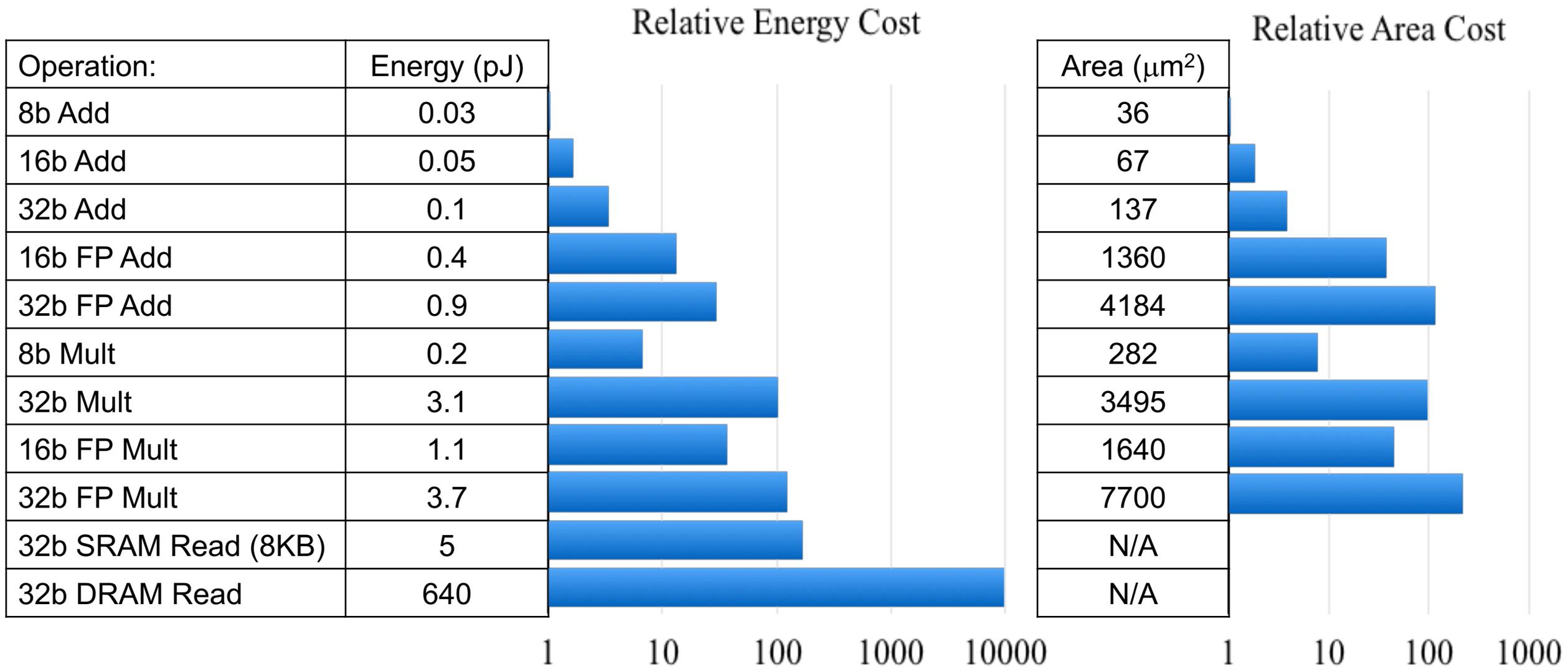
$$b_i = f\left(\sum_j w_{ij} a_j\right)$$



Number Representation

		Range	Accuracy
FP32		$10^{-38} - 10^{38}$.000006%
FP16		$6 \times 10^{-5} - 6 \times 10^4$.05%
Int32		$0 - 2 \times 10^9$	$\frac{1}{2}$
Int16		$0 - 6 \times 10^4$	$\frac{1}{2}$
Int8		$0 - 127$	$\frac{1}{2}$

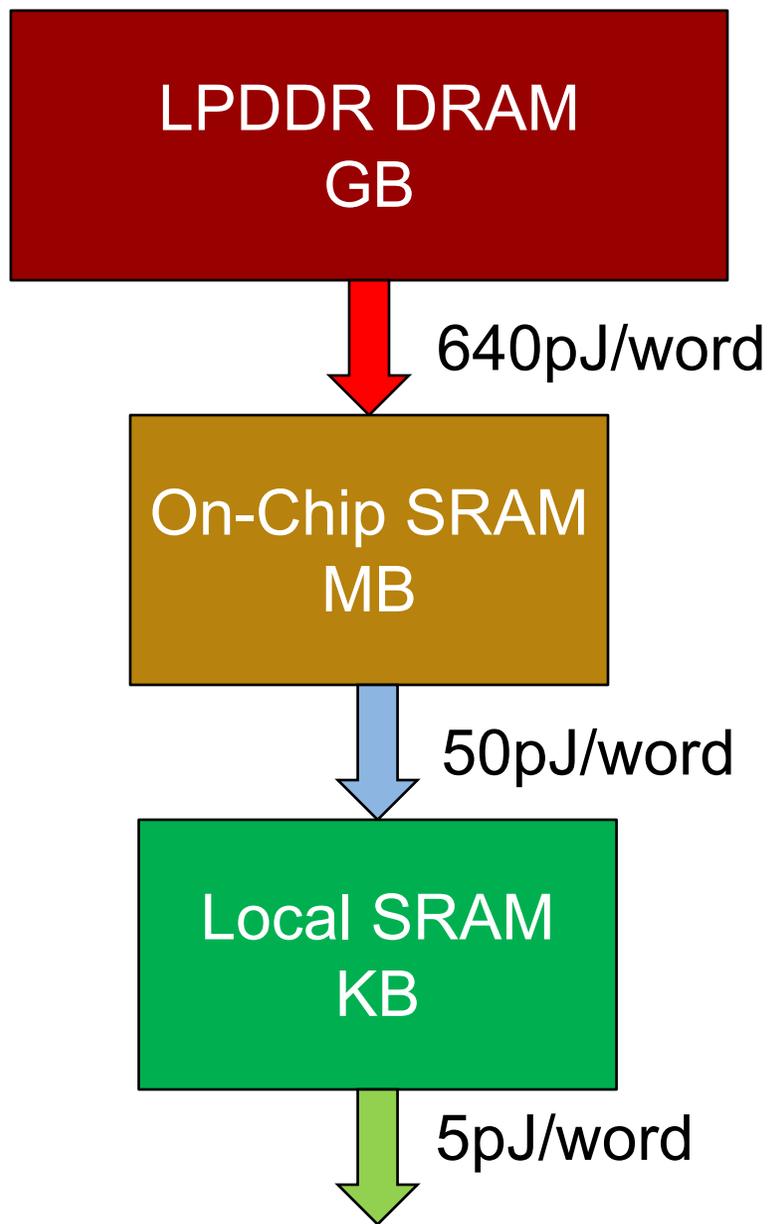
Cost of Operations



Energy numbers are from Mark Horowitz “Computing’s Energy Problem (and what we can do about it)”, ISSCC 2014

Area numbers are from synthesized result using Design Compiler under TSMC 45nm tech node. FP units used DesignWare Library.

The Importance of Staying Local



Mixed Precision

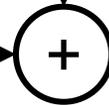
Store weights as 4b using
Trained quantization,
decode to 16b



Store activations as 16b



16b x 16b multiply
round result to 16b

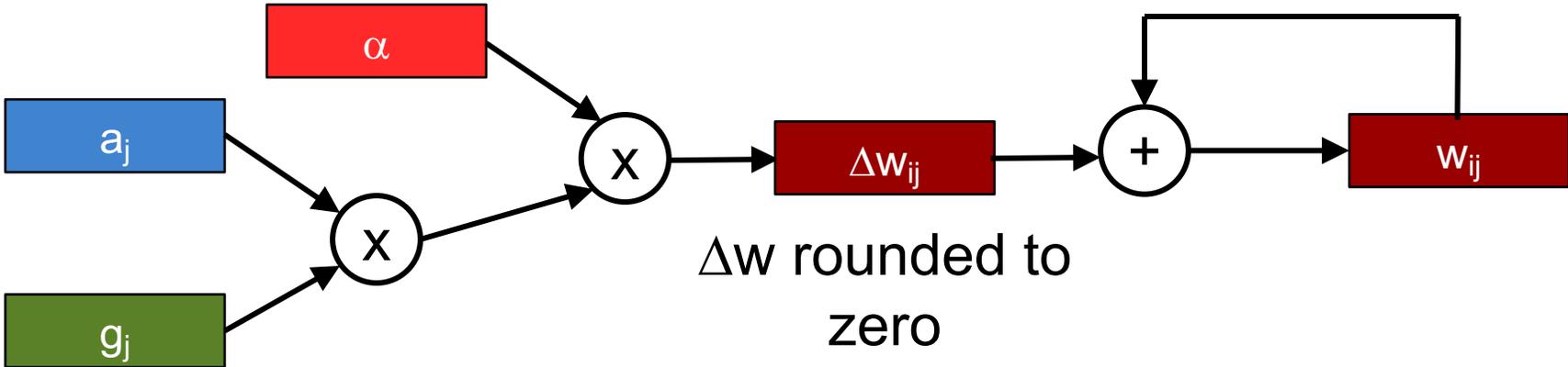


accumulate 24b or 32b
to avoid saturation

Batch normalization important to 'center' dynamic range

Weight Update

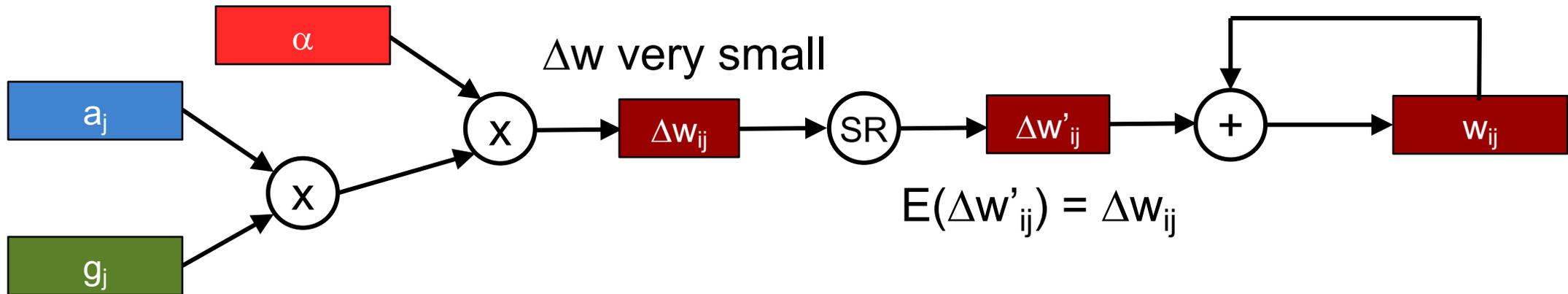
Learning rate may be very small (10^{-5} or less)



No learning!

Stochastic Rounding

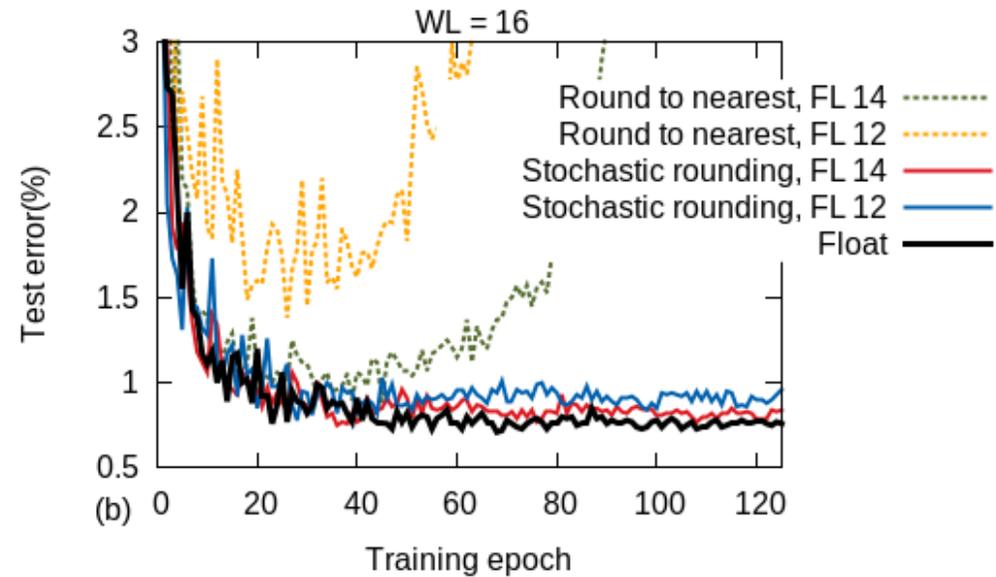
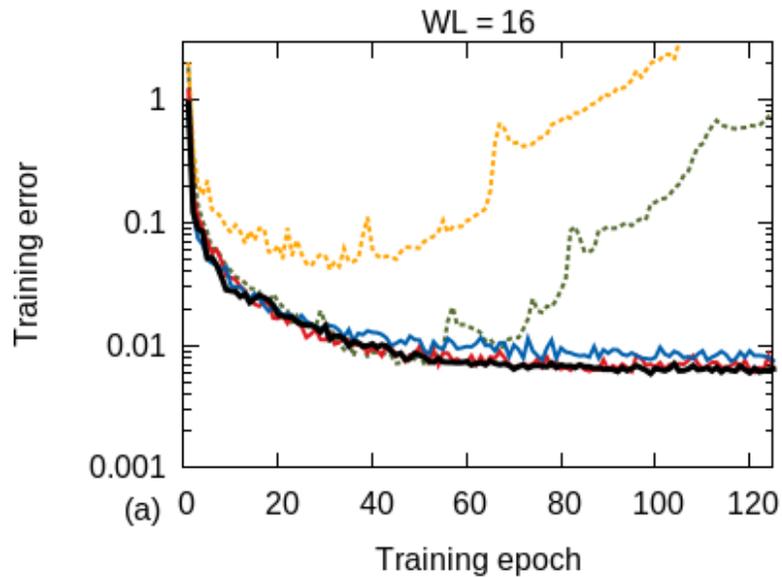
Learning rate may
be very small
(10^{-5} or less)



Reduced Precision For Training

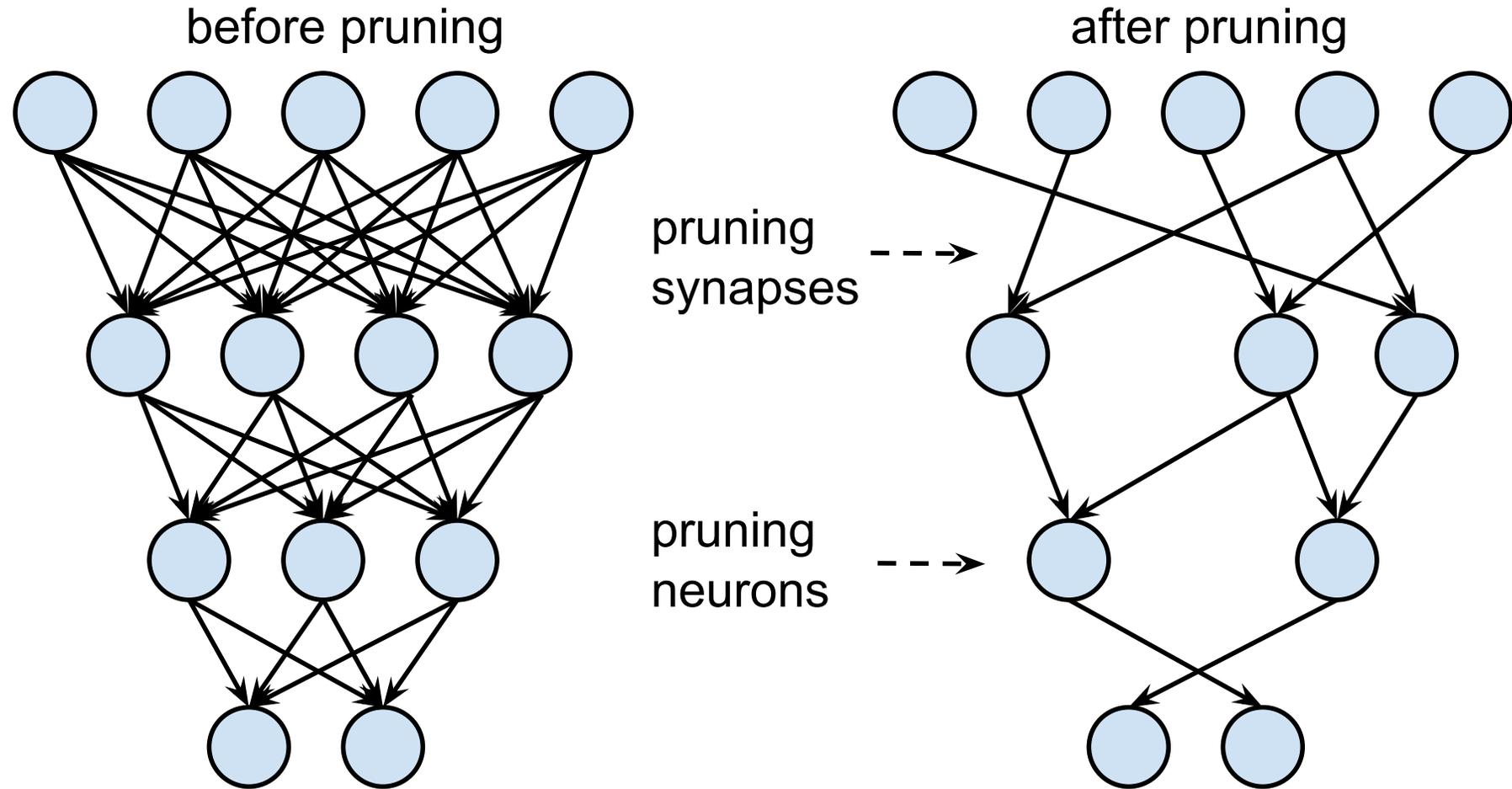
$$b_i = f \left(\sum_j w_{ij} a_j \right)$$

$$w_{ij} = w_{ij} + \alpha a_i g_j$$

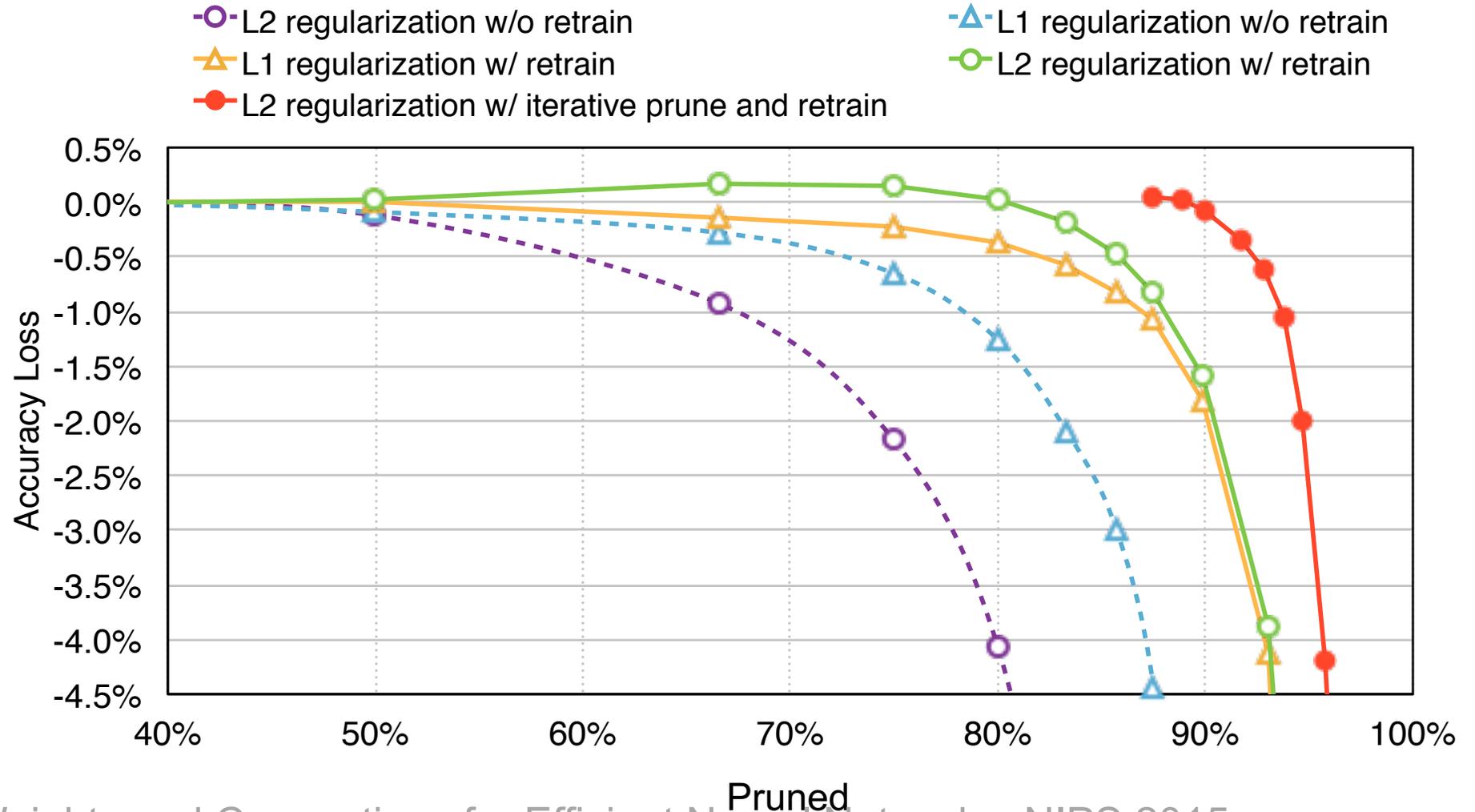
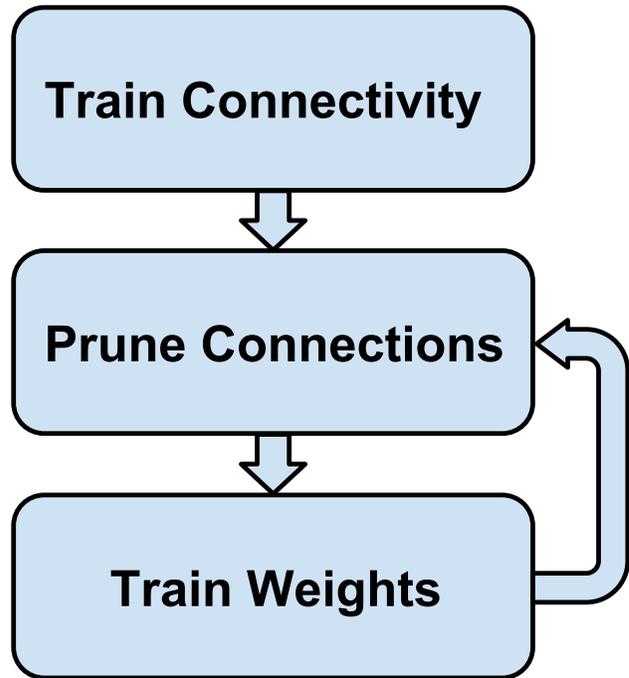


Pruning

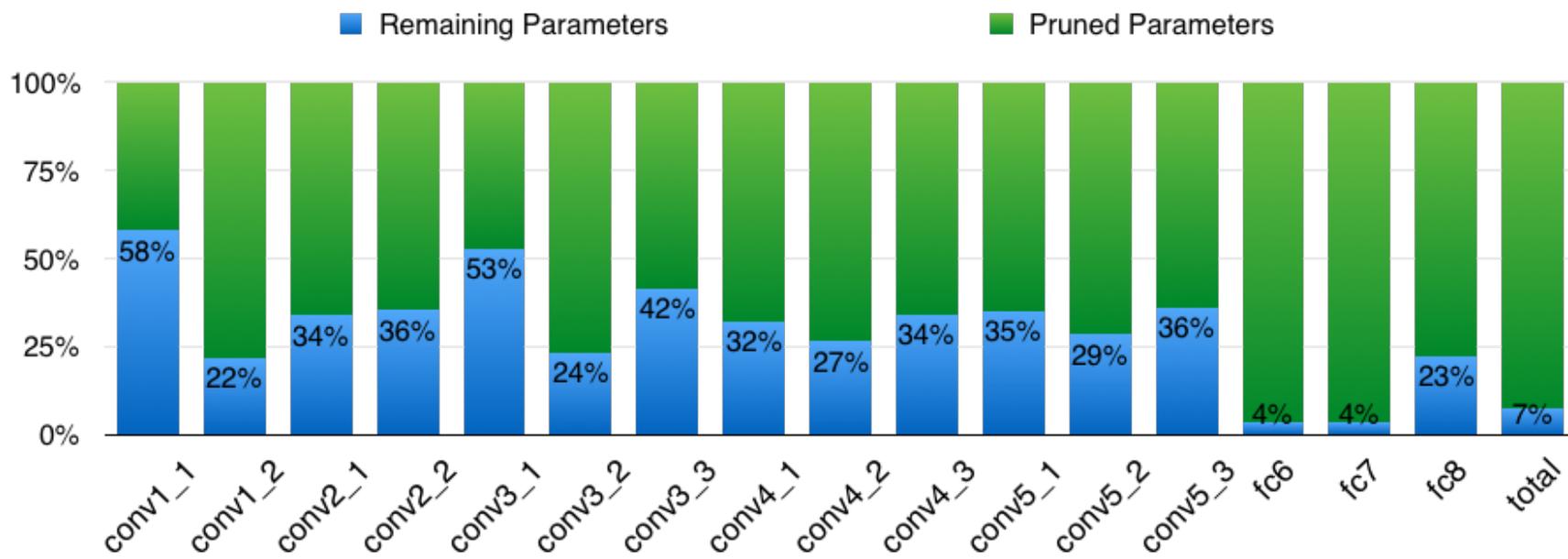
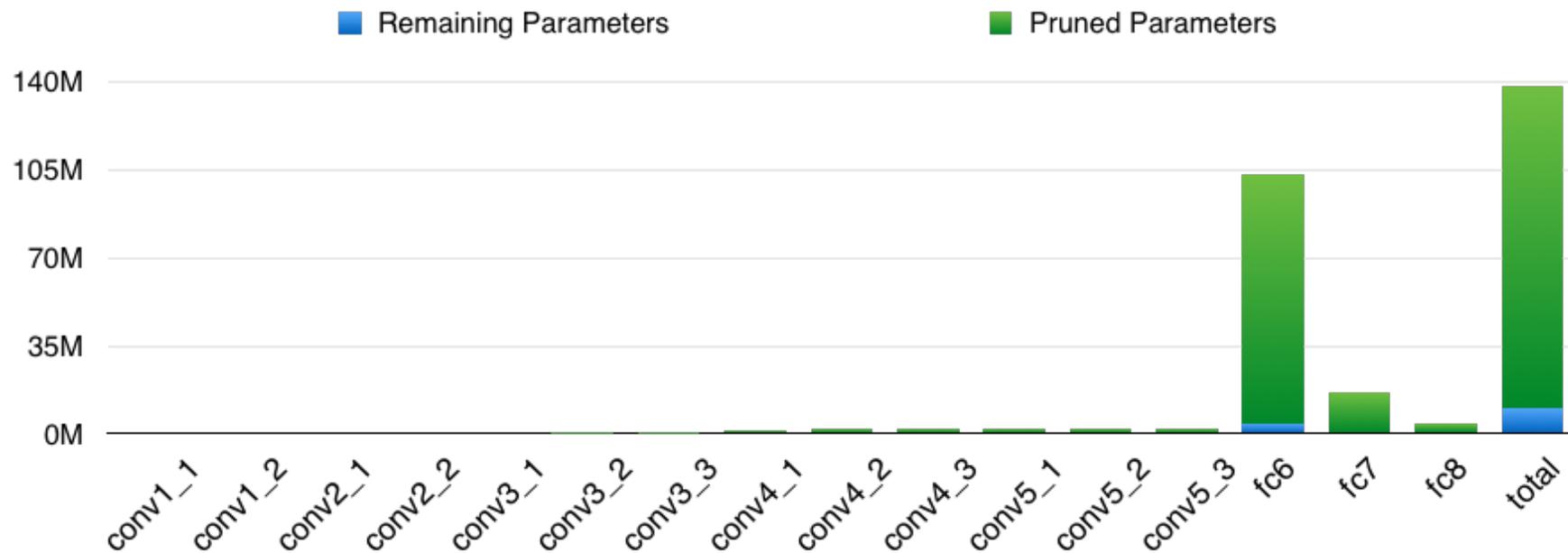
Pruning



Retrain to Recover Accuracy



Pruning of VGG-16



Pruning Neural Talk and LSTM



- **Original:** a basketball player in a white uniform is playing with a **ball**
- **Pruned 90%:** a basketball player in a white uniform is playing with a **basketball**



- **Original :** a brown dog is running through a grassy **field**
- **Pruned 90%:** a brown dog is running through a grassy **area**



- **Original :** a man is riding a surfboard on a wave
- **Pruned 90%:** a man in a wetsuit is riding a wave **on a beach**



- **Original :** a soccer player in red is running in the field
- **Pruned 95%:** a man in a **red shirt and black and white black shirt** is running through a field

Speedup of Pruning on CPU/GPU

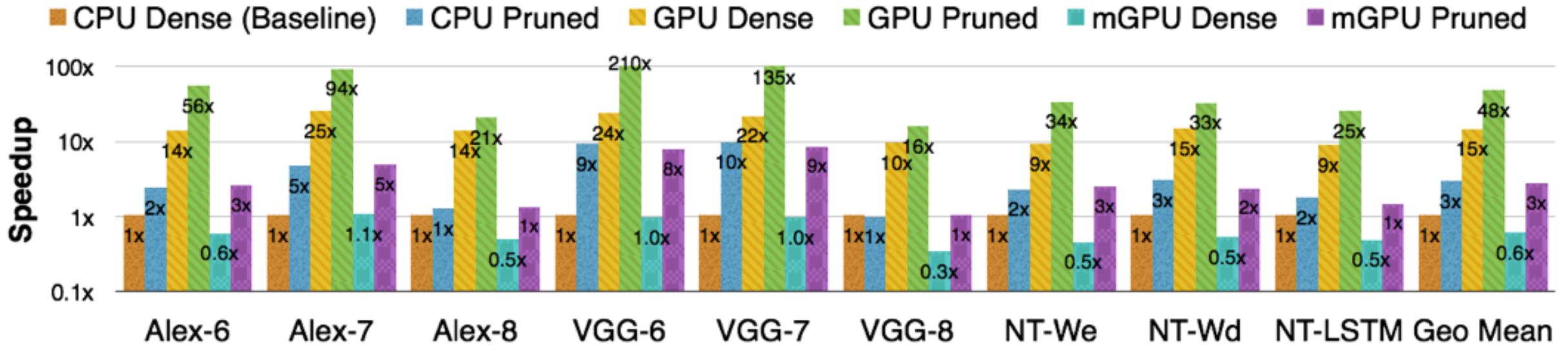
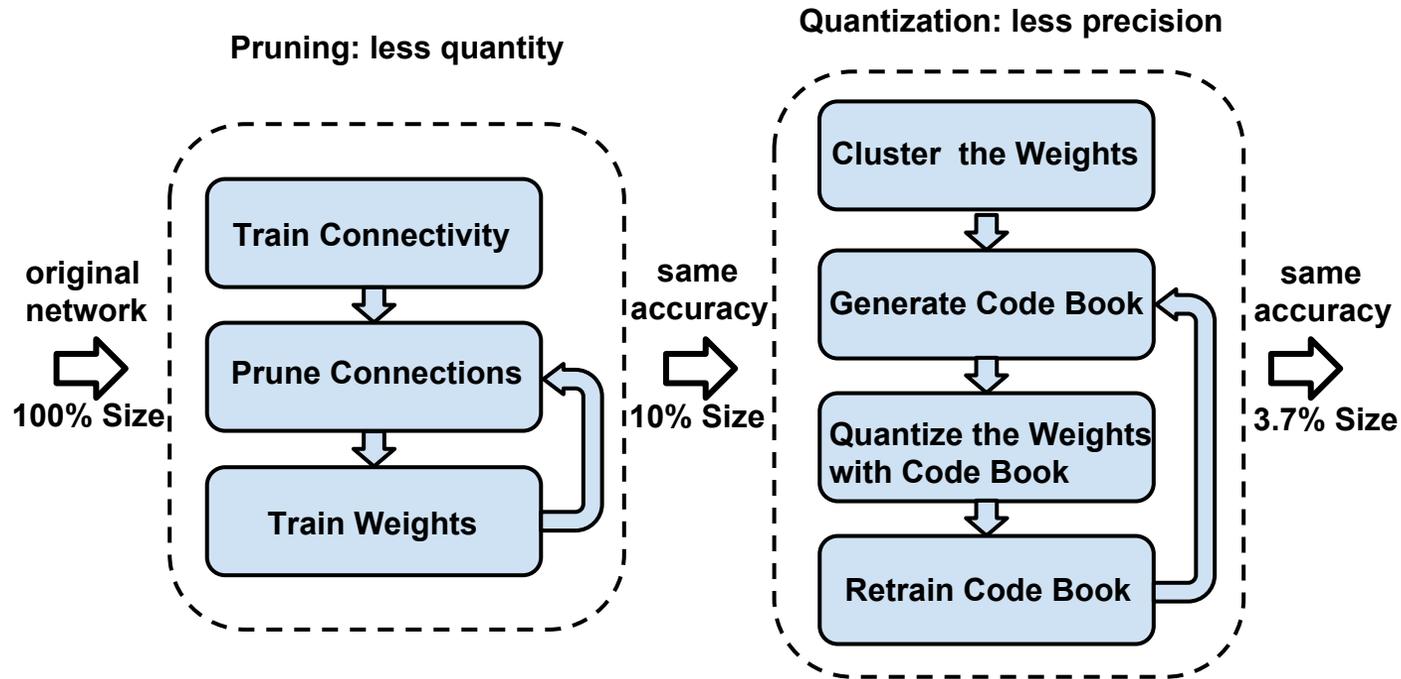


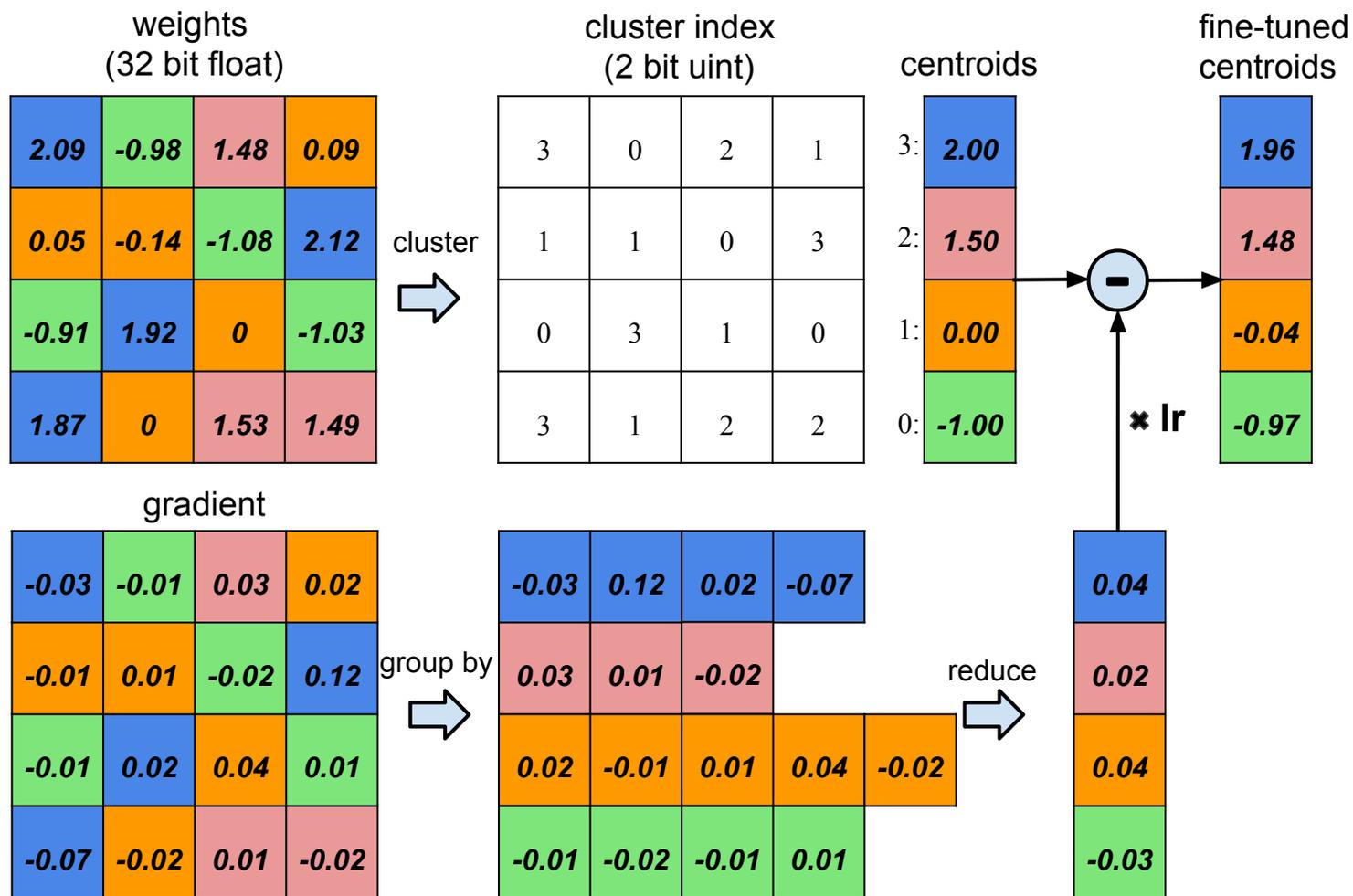
Figure 9: Compared with the original network, pruned network layer achieved $3\times$ speedup on CPU, $3.5\times$ on GPU and $4.2\times$ on mobile GPU on average. Batch size = 1 targeting real time processing. Performance number normalized to CPU.

Intel Core i7 5930K: MKL CBLAS GEMV, MKL SPBLAS CSRMMV
NVIDIA GeForce GTX Titan X: cuBLAS GEMV, cuSPARSE CSRMMV
NVIDIA Tegra K1: cuBLAS GEMV, cuSPARSE CSRMMV

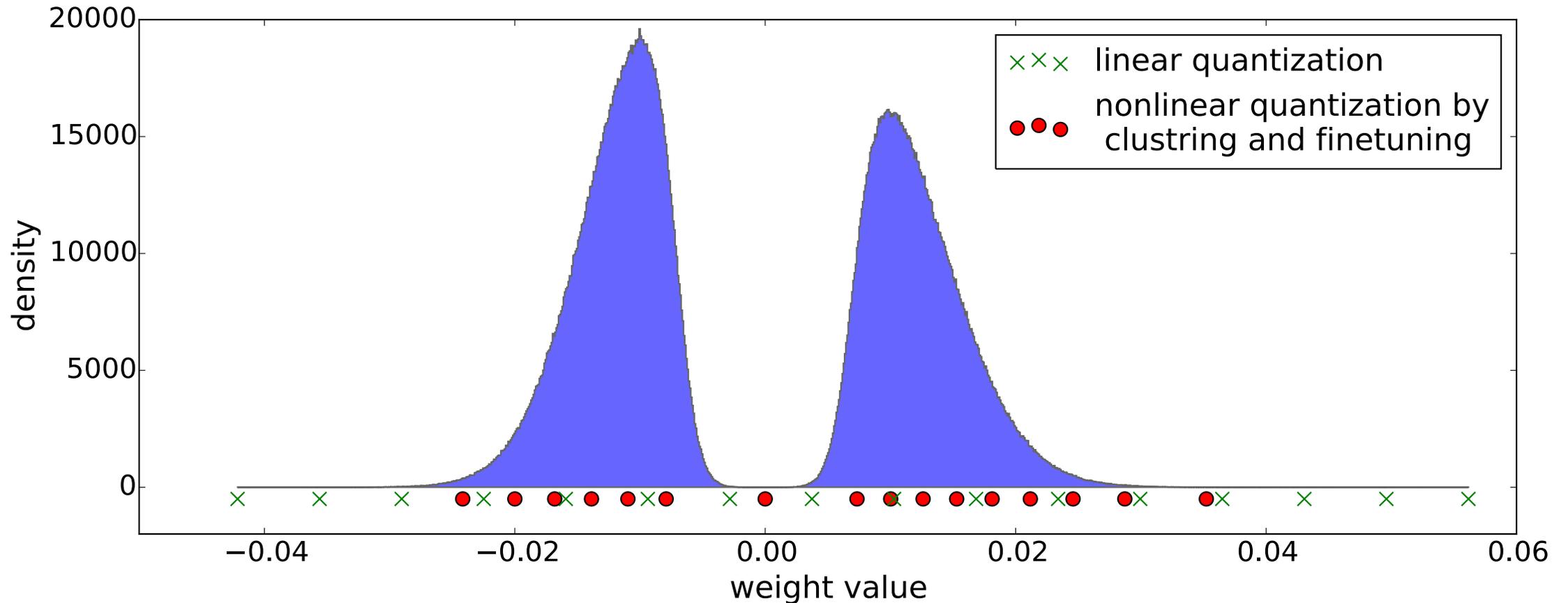
Trained Quantization (Weight Sharing)



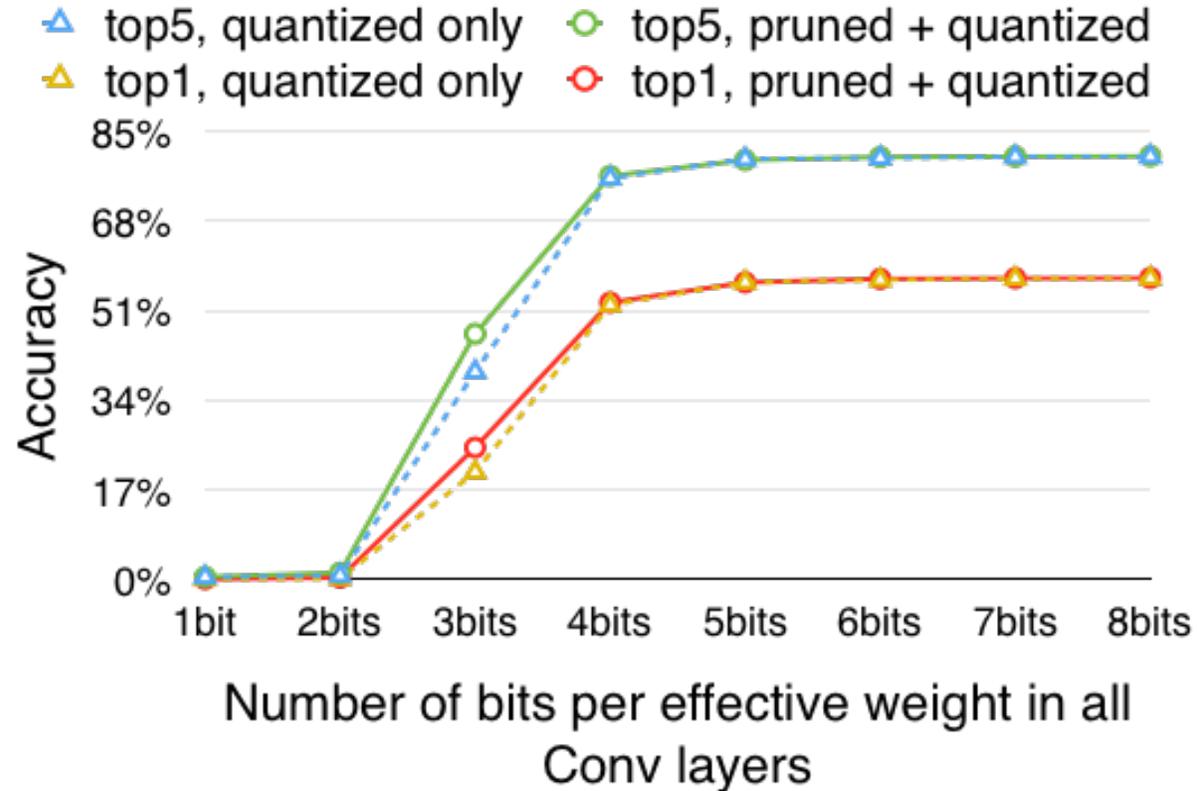
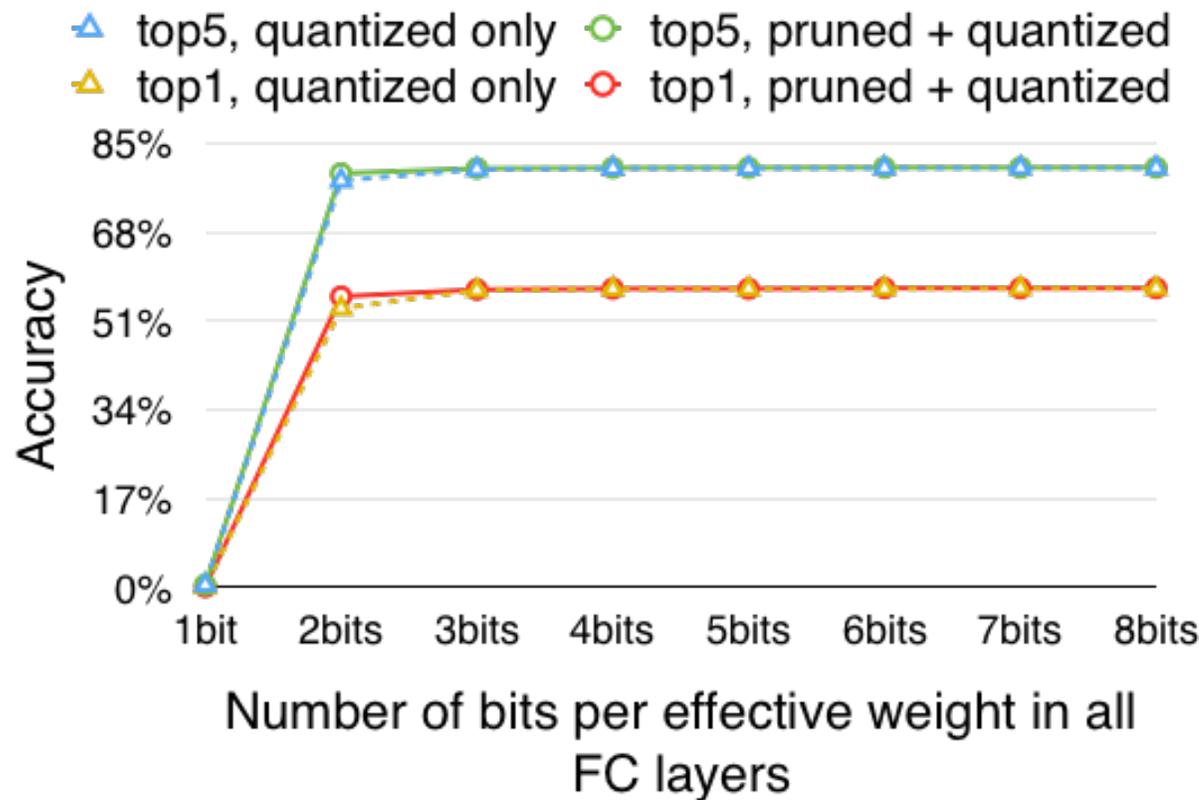
Weight Sharing via K-Means



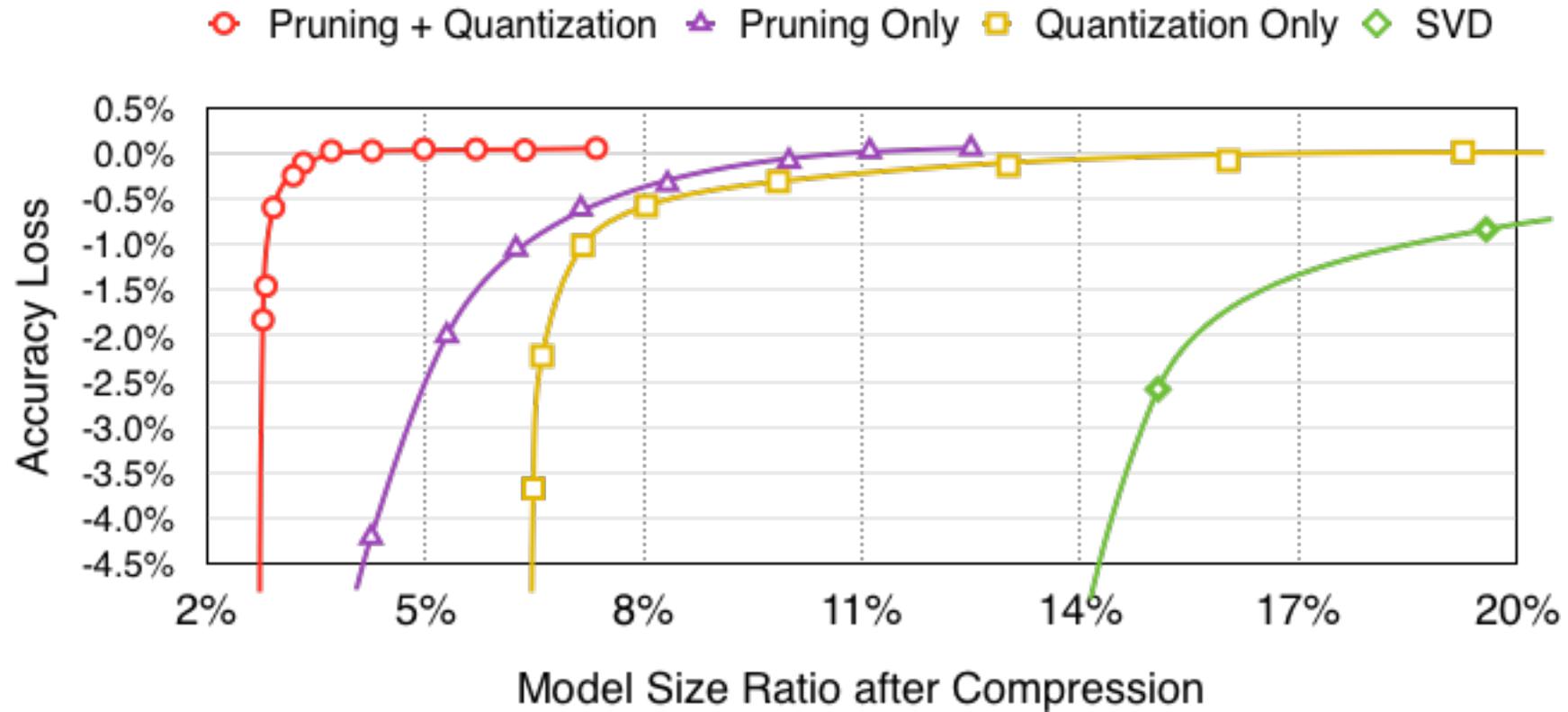
Trained Quantization



Bits per Weight



Pruning + Trained Quantization



30x – 50x Compression Means

- Complex DNNs can be put in mobile applications (<100MB total)
 - 1GB network (250M weights) becomes 20-30MB
- Memory bandwidth reduced by 30-50x
 - Particularly for FC layers in real-time applications with no reuse
- Memory working set fits in on-chip SRAM
 - 5pJ/word access vs 640pJ/word

Efficient Inference Engine

Sparse Matrix Representation

$$\vec{a} \begin{pmatrix} 0 & a_1 & 0 & a_3 \end{pmatrix} \times \begin{pmatrix} PE0 & w_{0,0} & w_{0,1} & 0 & w_{0,3} \\ PE1 & 0 & 0 & w_{1,2} & 0 \\ PE2 & 0 & w_{2,1} & 0 & w_{2,3} \\ PE3 & 0 & 0 & 0 & 0 \\ & 0 & 0 & w_{4,2} & w_{4,3} \\ & w_{5,0} & 0 & 0 & 0 \\ & 0 & 0 & 0 & w_{6,3} \\ & 0 & w_{7,1} & 0 & 0 \end{pmatrix} = \begin{pmatrix} b_0 \\ b_1 \\ -b_2 \\ b_3 \\ -b_4 \\ b_5 \\ b_6 \\ -b_7 \end{pmatrix} \xrightarrow{ReLU} \vec{b} \begin{pmatrix} b_0 \\ b_1 \\ 0 \\ b_3 \\ 0 \\ b_5 \\ b_6 \\ 0 \end{pmatrix}$$

Sparse Matrix Representation

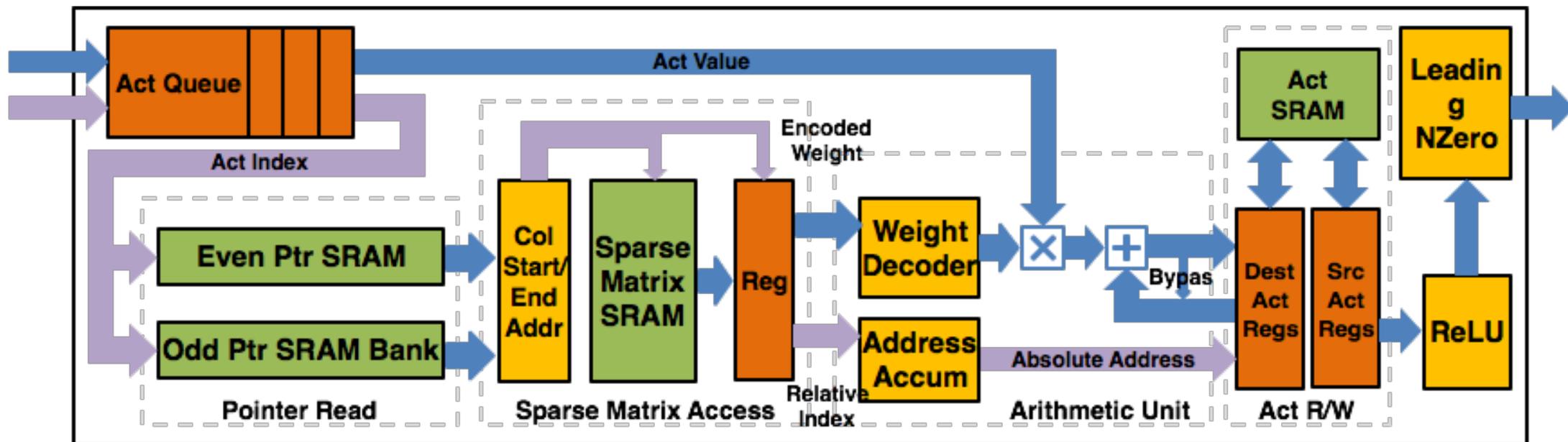
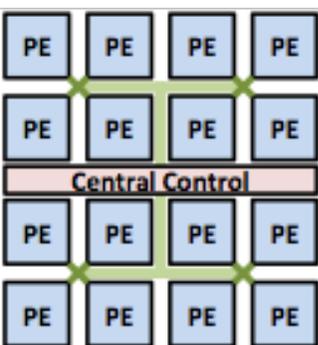
$$\vec{a} \begin{pmatrix} 0 & a_1 & 0 & a_3 \end{pmatrix}$$

×

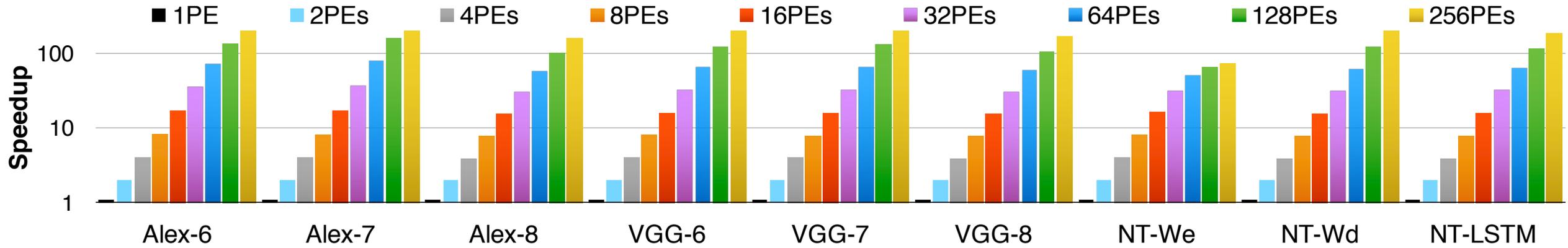
$$\begin{pmatrix} PE0 \\ PE1 \\ PE2 \\ PE3 \\ \\ \\ \\ \end{pmatrix} \begin{pmatrix} w_{0,0} & w_{0,1} & 0 & w_{0,3} \\ 0 & 0 & w_{1,2} & 0 \\ 0 & w_{2,1} & 0 & w_{2,3} \\ 0 & 0 & 0 & 0 \\ 0 & 0 & w_{4,2} & w_{4,3} \\ w_{5,0} & 0 & 0 & 0 \\ 0 & 0 & 0 & w_{6,3} \\ 0 & w_{7,1} & 0 & 0 \end{pmatrix}$$

Virtual Weight	$W_{0,0}$	$W_{0,1}$	$W_{4,2}$	$W_{0,3}$	$W_{4,3}$
Relative Index	0	1	2	0	0
Column Pointer	0	1	2	3	

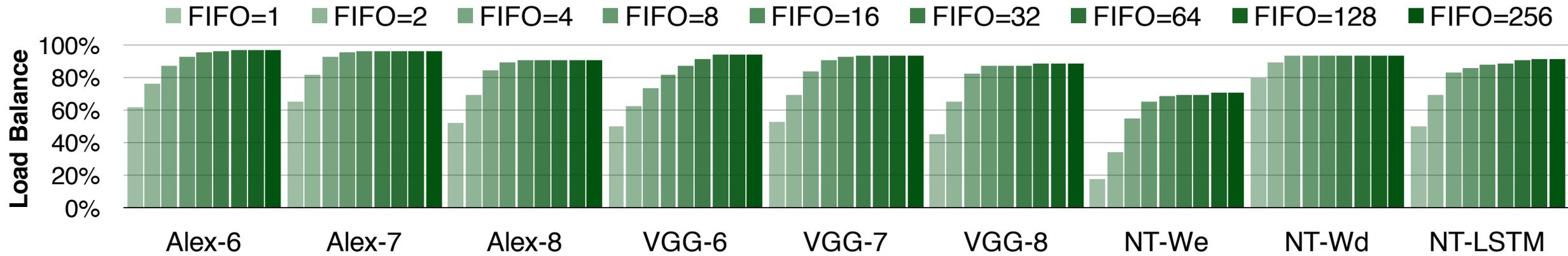
EIE Architecture



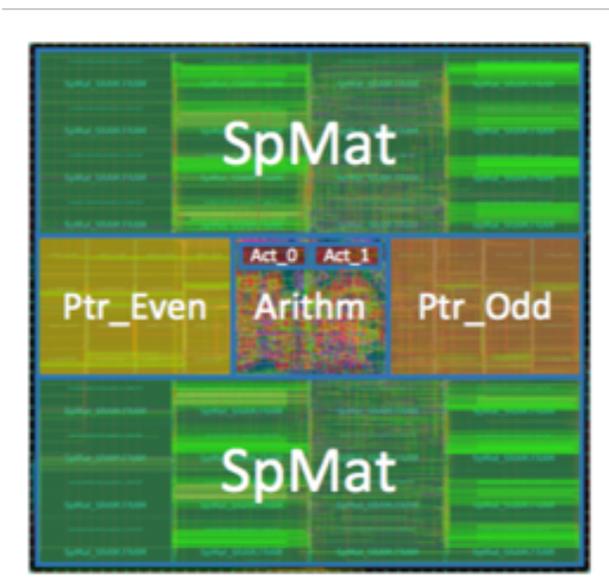
Scalability



Load Balance



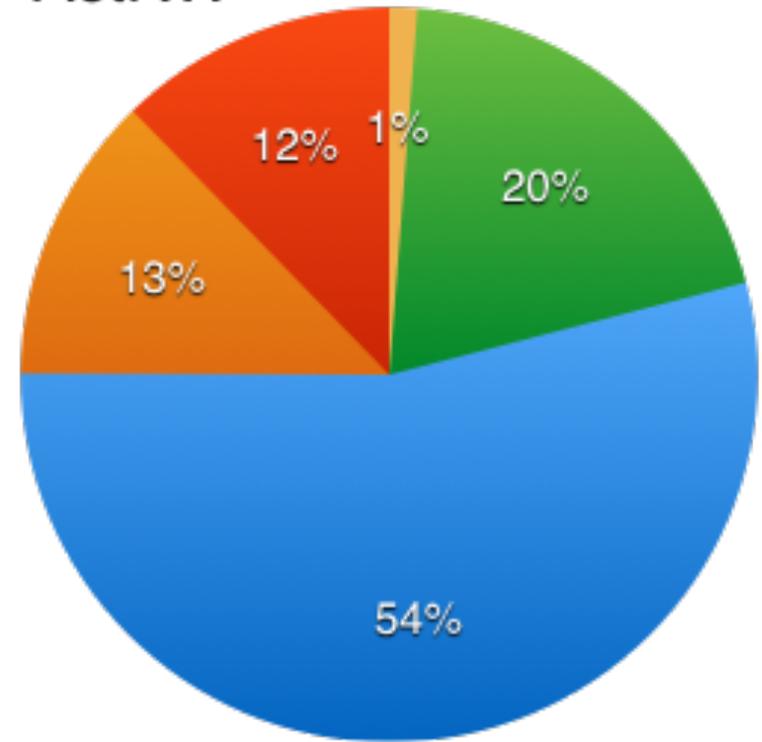
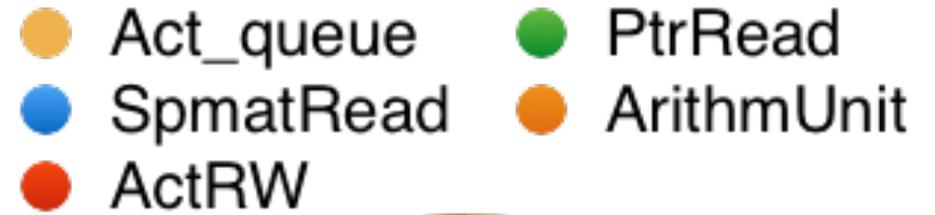
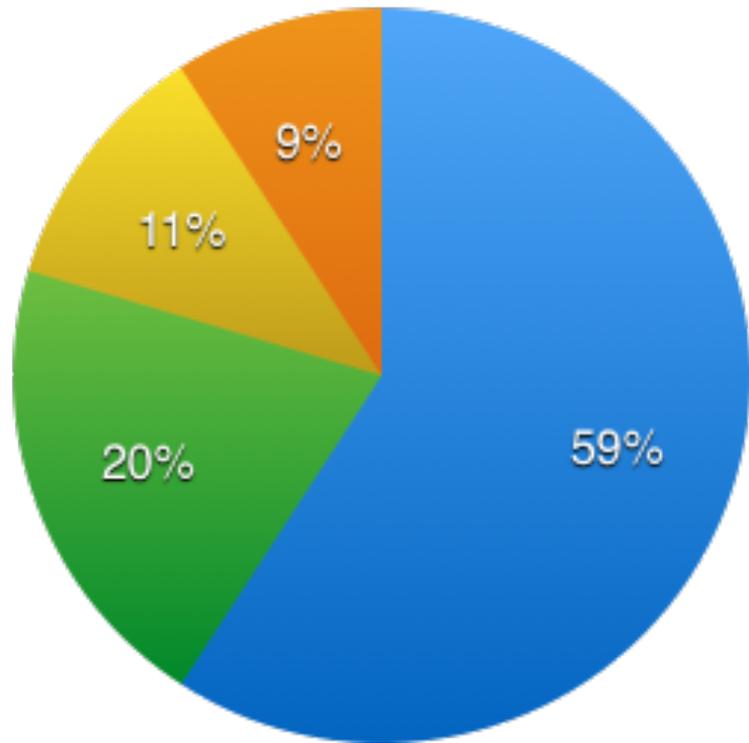
Implementation



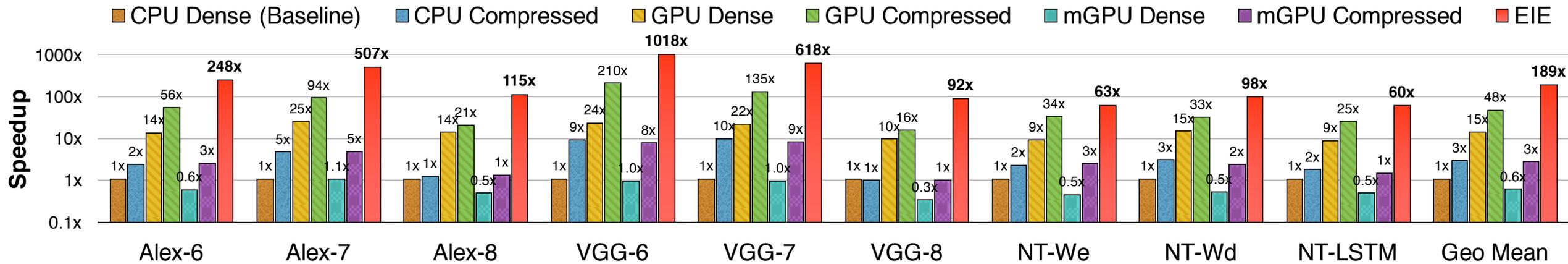
Technology	45 nm
# PEs	64
on-chip SRAM	8 MB
Max Model Size	84 Million
Static Sparsity	10x
Dynamic Sparsity	3x
Quantization	4-bit
ALU Width	16-bit
Area	40.8 mm ²
MxV Throughput	81,967 layers/s
Power	586 mW

1. Post layout result
2. Throughput on AlexNet FC-7

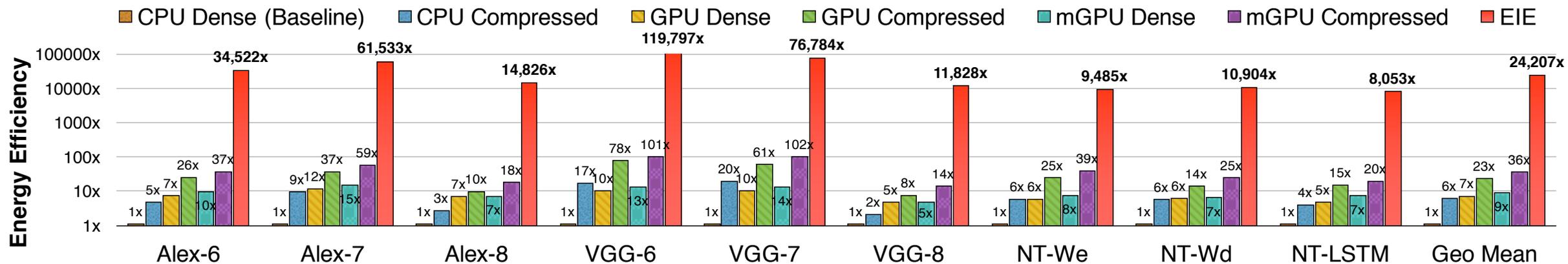
Energy Distribution



FC Layer: Speedup on EIE

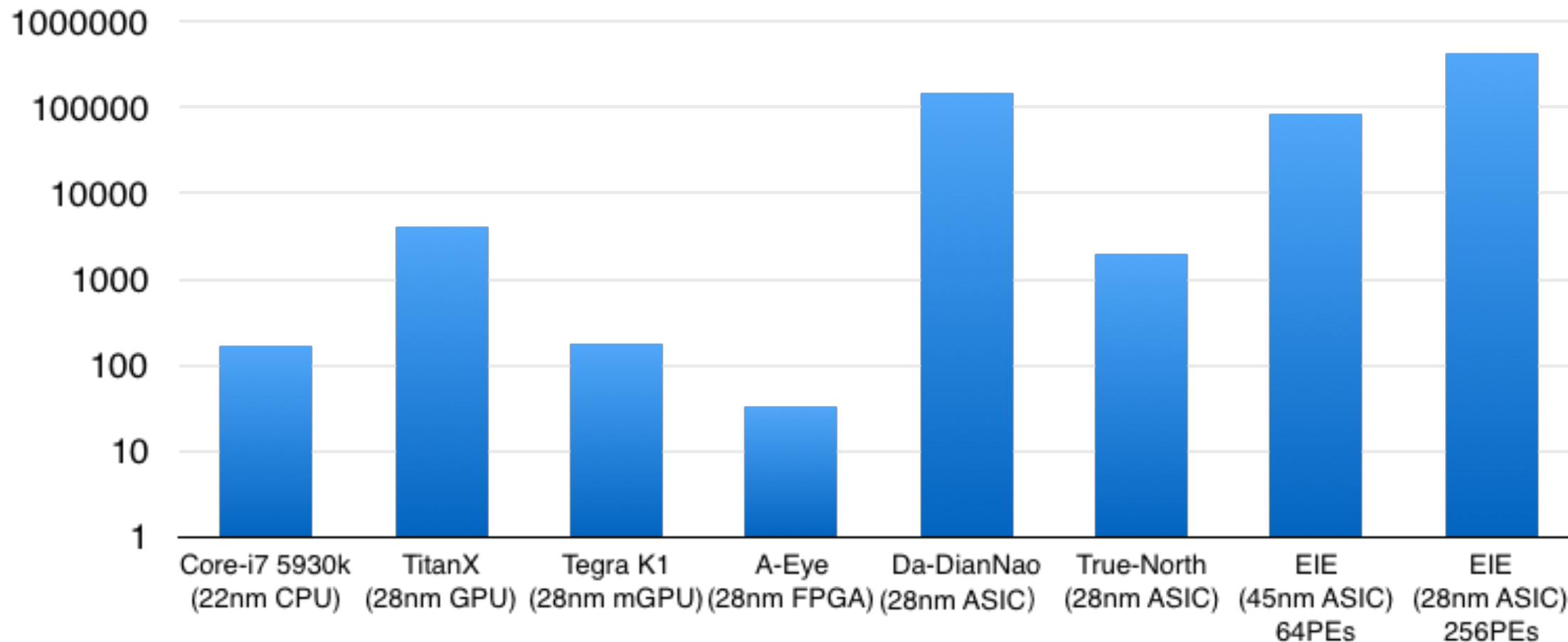


FC Layer: Energy Efficiency on EIE

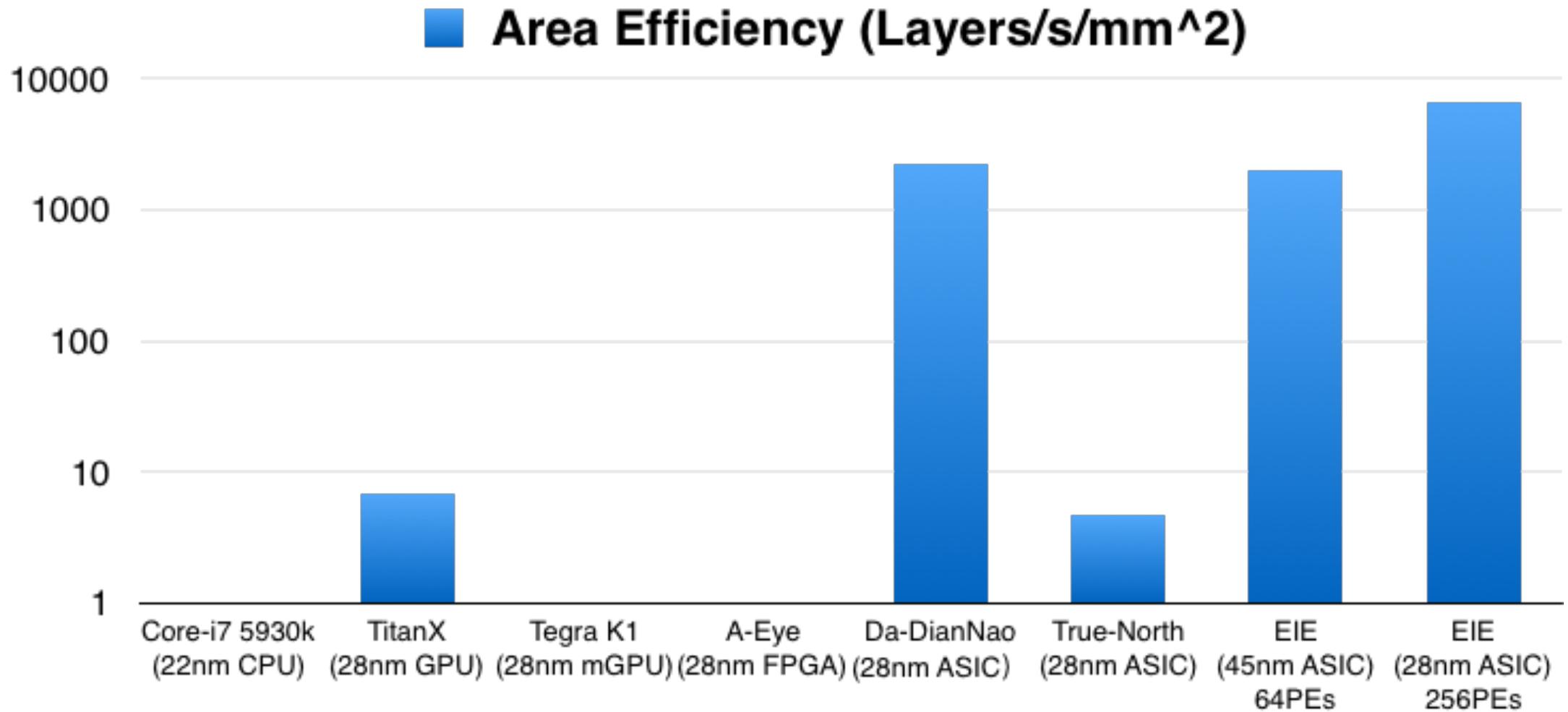


Comparison: Throughput

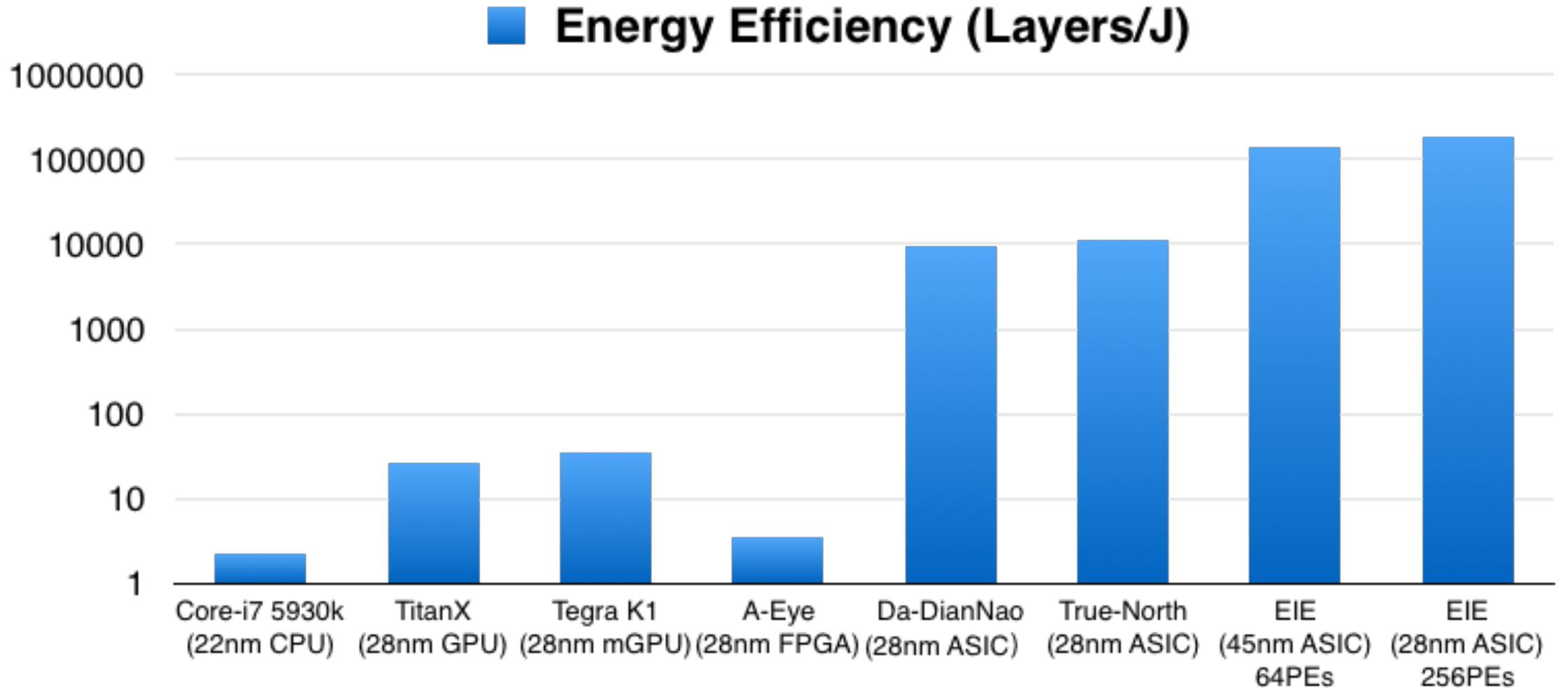
■ MxV Throughput (Layers/s)



Comparison: Area Efficiency



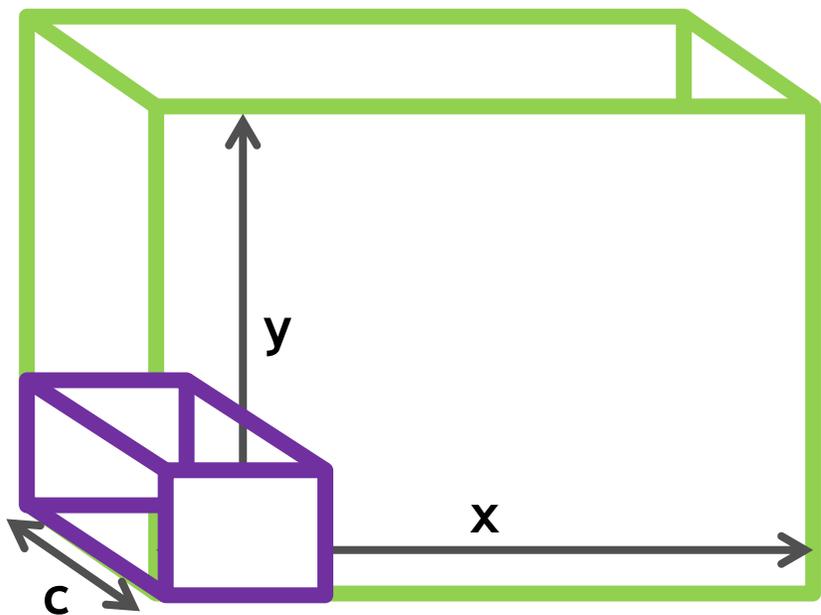
Comparison: Energy Efficiency



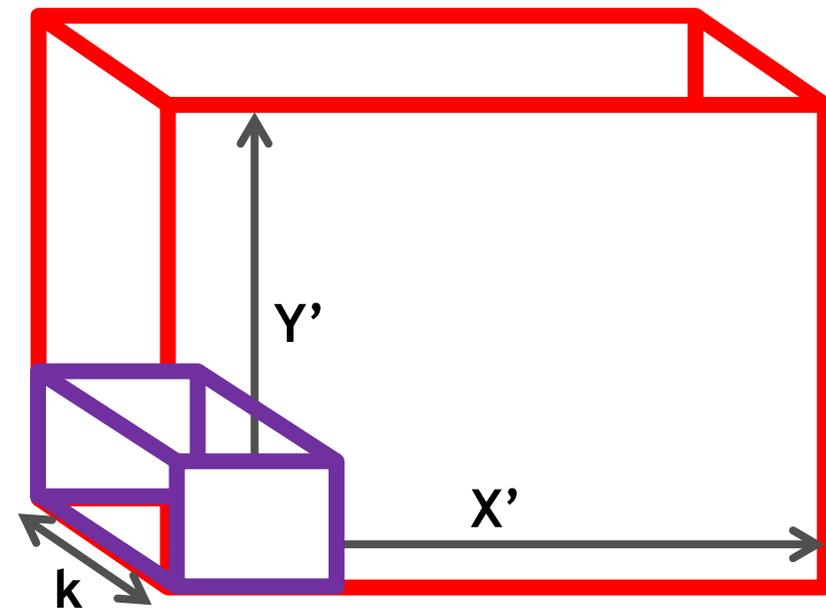
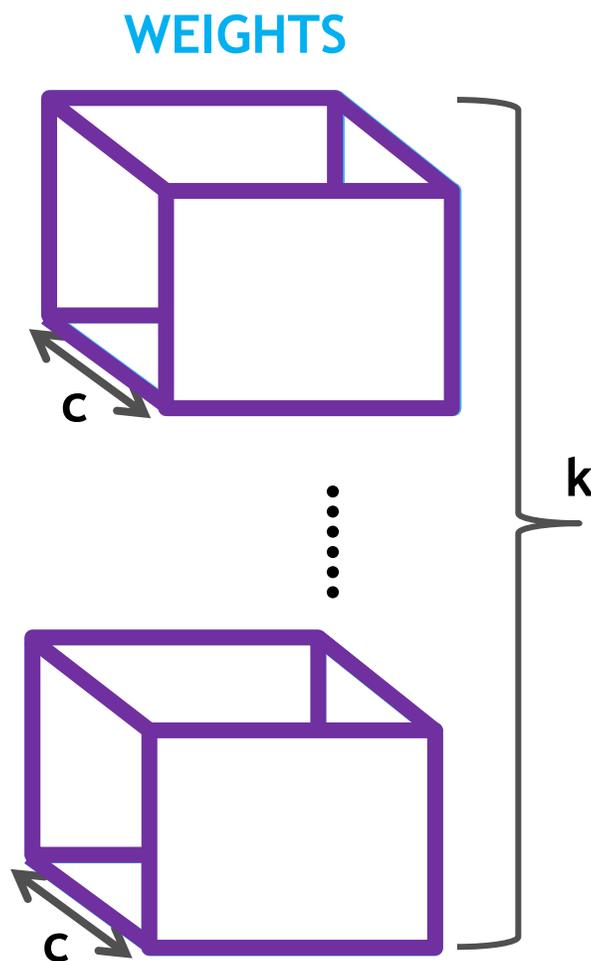
Sparse Convolutional Accelerator

Blocking CNN Inference

PURPLE: allocated to 1 PE



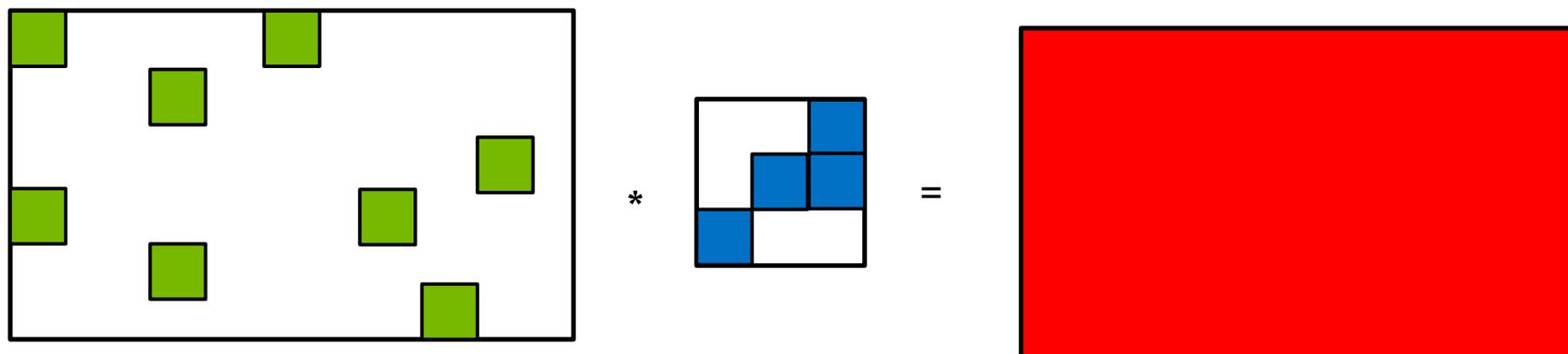
INPUT ACTIVATIONS



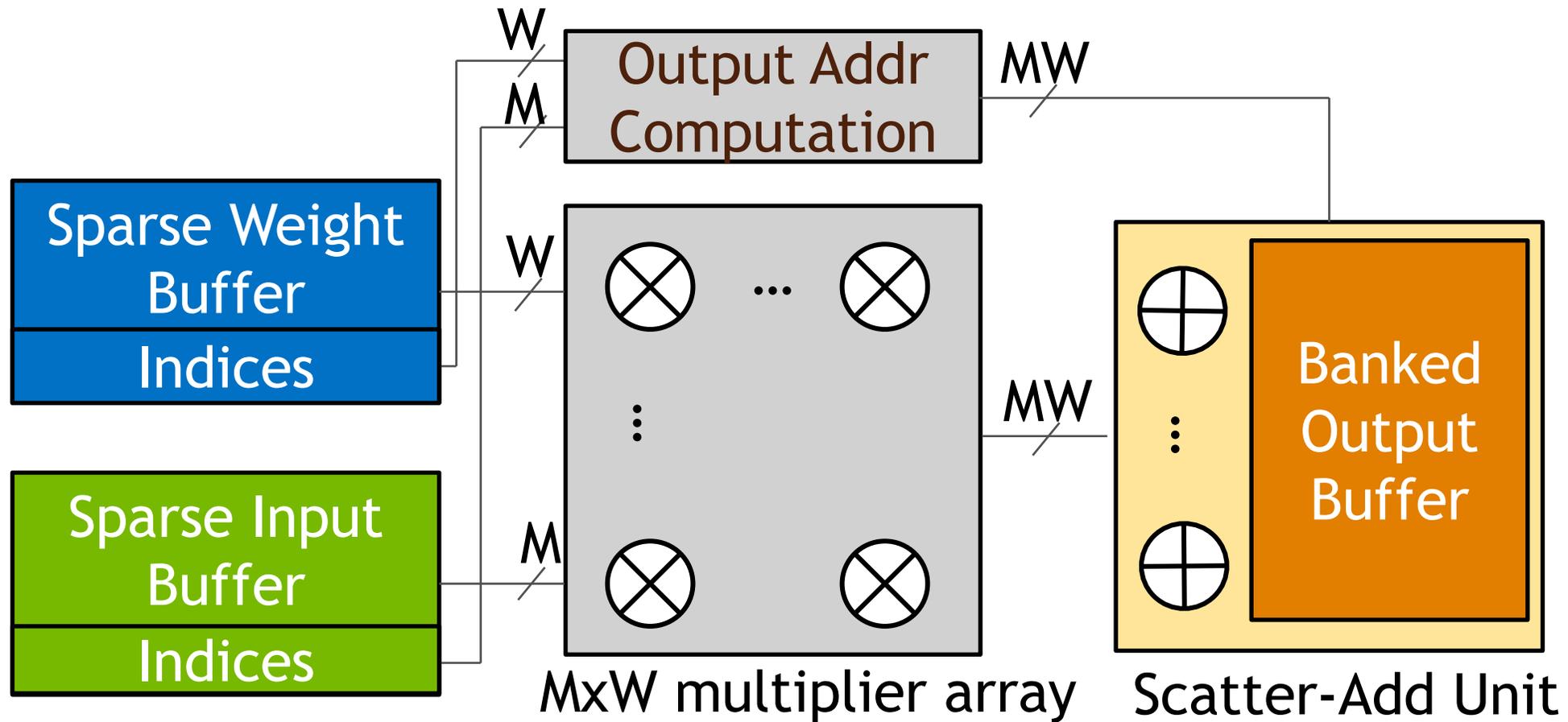
OUTPUT ACTIVATIONS

Sparse Convolution

- Only compute where both operands are nonzero
- 10-30x Reduction in work

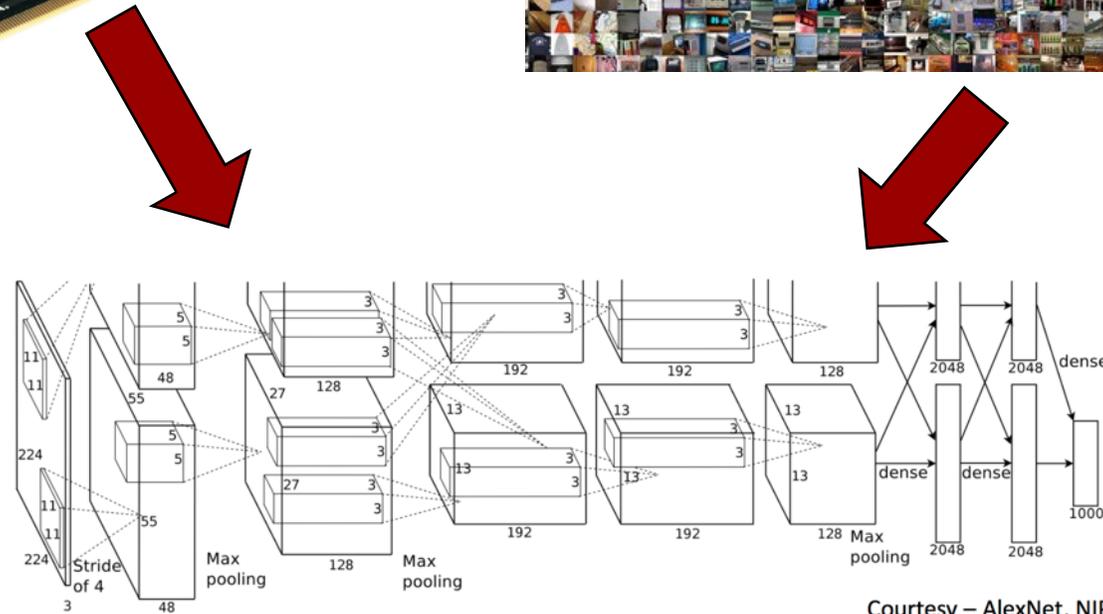


Sparse Convolution Engine



Conclusion

Hardware and Data enable DNNs



Courtesy – AlexNet, NIPS 2012

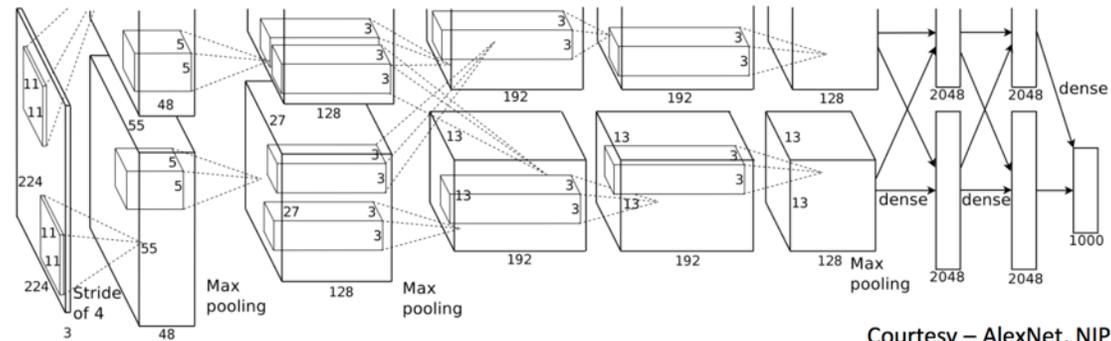
Summary

- **Hardware** has enabled the current resurgence of DNNs
 - And limits the size of today's networks
- **Inference**
 - Dynamically **sparse** activations x statically **sparse** weights
 - **8b weights** sufficient (can be compressed to **2-4b**)
 - Energy dominated by **data movement** and **buffering**
 - **Fixed-function hardware** will dominate **inference**
- **Training**
 - Only **dynamic sparsity** (3x activations, 2x dropout)
 - Medium precision (**FP16** – for weights)
 - **Large memory footprint** (batch x retained activations) – can be 10s – 100s of GB
 - Parallelism to 10PF today 100PF in near future (Communication BW)
 - **GPUs** will dominate **training**

4 Distinct Sub-problems

	Training	Inference	
Convolutional	32b FP Batch Activation Storage GPUs ideal Comm for Parallelism	Low-Precision Compressed Latency-Sensitive Fixed-Function HW Arithmetic Dominated	B x S Weight Reuse Act Dominated
Fully-Conn.	32b FP Batch Weight Storage GPUs ideal Comm. for Parallelism	Low-Precision Compressed Latency-Sensitive No weight reuse Fixed-Function HW Storage dominated	B Weight Reuse Weight Dominated
	32b FP - large batches Minimize Training Time Enables larger networks	8b Int - small (unit) batches Meet real-time constraint	

Thank You



Courtesy – AlexNet, NIPS 2012