

Network and Hardware Co-Design

MICRO Tutorial (2016)


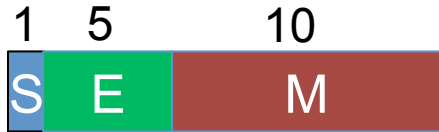

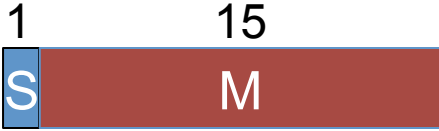

Website: <http://eyeriss.mit.edu/tutorial.html>

Joel Emer, Vivienne Sze, Yu-Hsin Chen

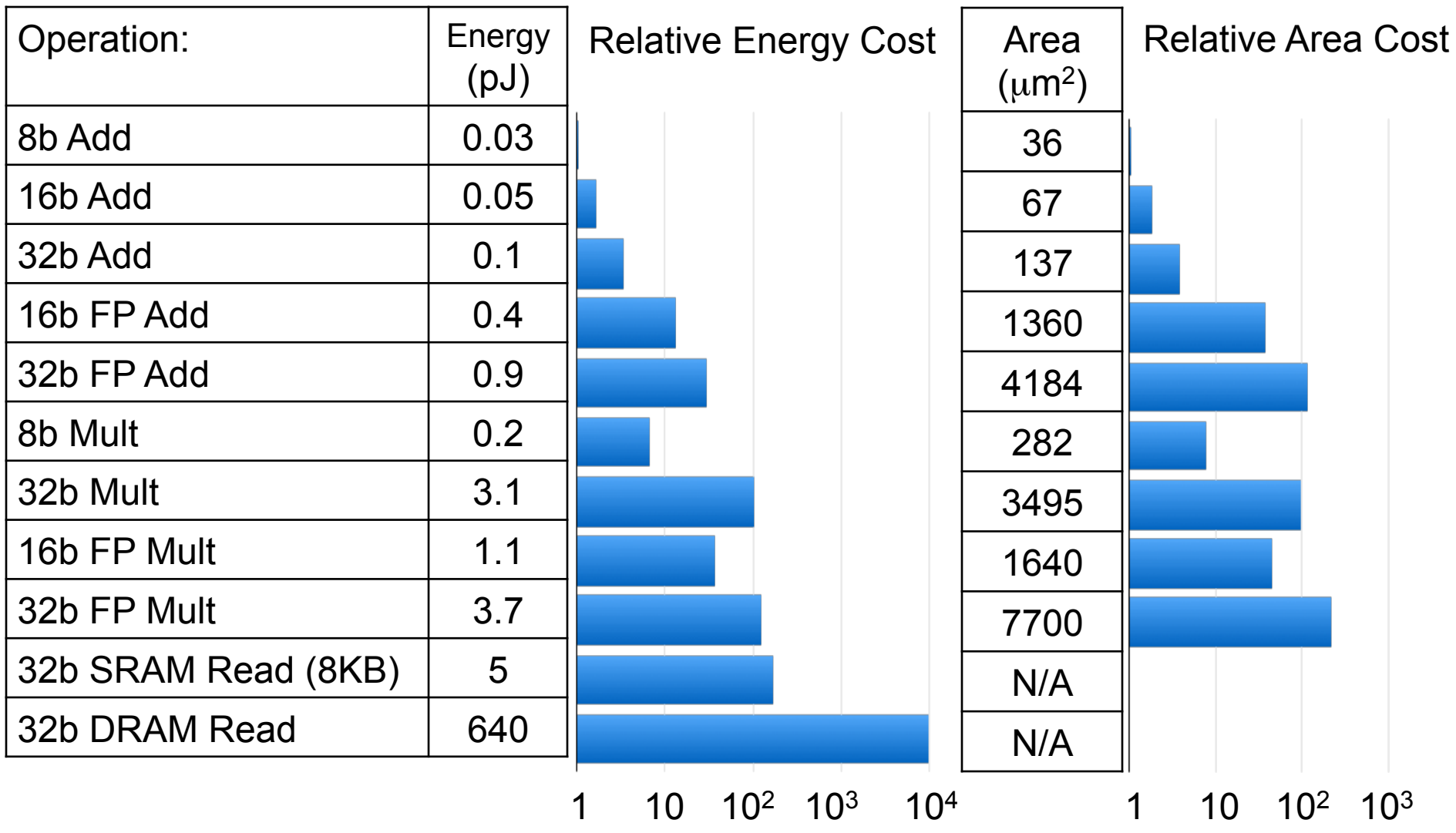
Network Optimization

- **Reduce precision of operations and operands**
 - Fixed and Floating point
 - Bit-width
- **Reduce number of operations and storage of weights**
 - Compression
 - Pruning
 - Network Architectures

Number Representation

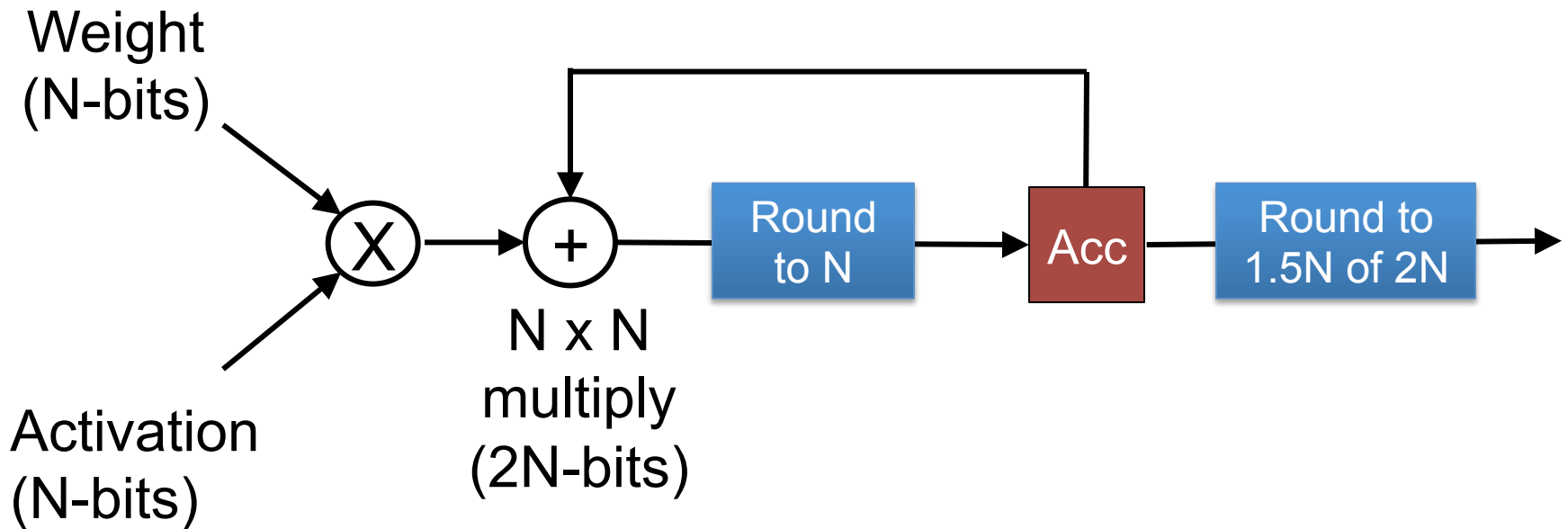
| | | Range | Accuracy |
|-------|---|------------------------------------|---------------|
| FP32 |  | $10^{-38} - 10^{38}$ | .000006% |
| FP16 |  | $6 \times 10^{-5} - 6 \times 10^4$ | .05% |
| Int32 |  | $0 - 2 \times 10^9$ | $\frac{1}{2}$ |
| Int16 |  | $0 - 6 \times 10^4$ | $\frac{1}{2}$ |
| Int8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

Cost of Operations



[Horowitz, "Computing's Energy Problem (and what we can do about it)", ISSCC 2014]

N-bit Precision

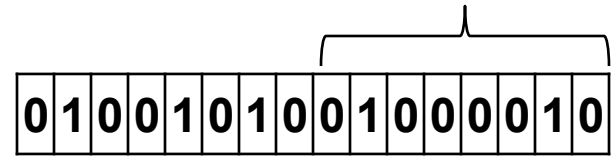


Methods to Reduce Bits

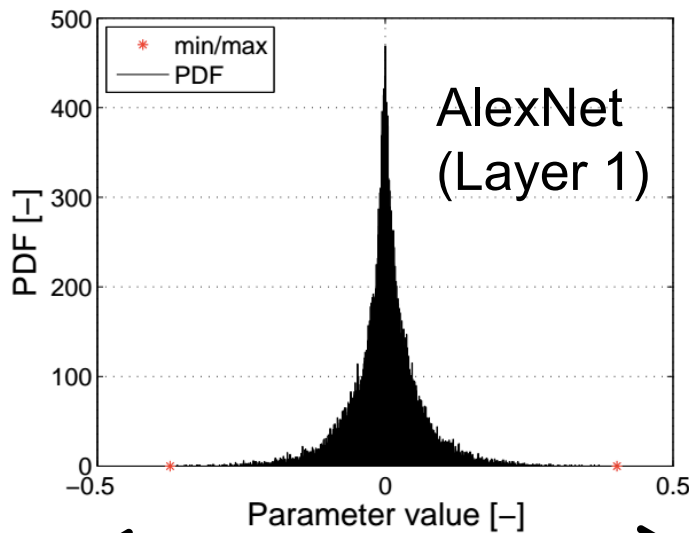
- Quantization/Rounding
- Dynamic Fixed Point
 - Rescale and Reduce bits
- Fine-tuning: Retrain Weights

Example: 16-bit \rightarrow 8-bits

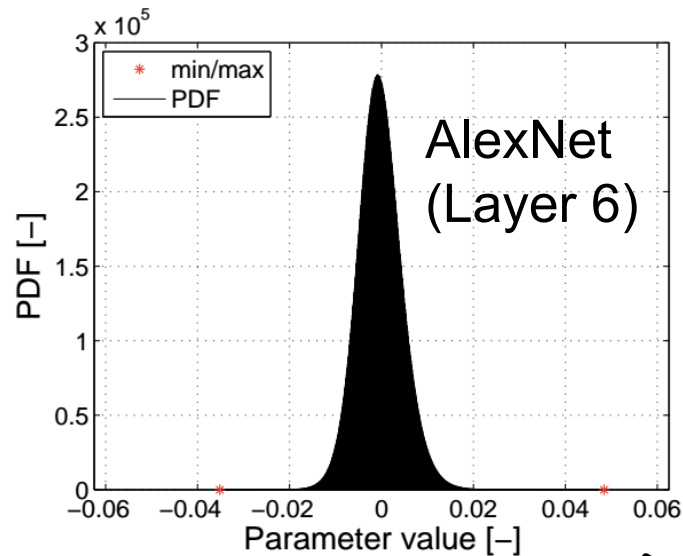
$$2^{11} + 2^9 + 2^6 + 2^1 = 2626 \text{ (overflow)}$$



$$2^{10} + 2^7 + 2^5 + 2^2 + 2^0 = 1189$$



Dynamic range = 1

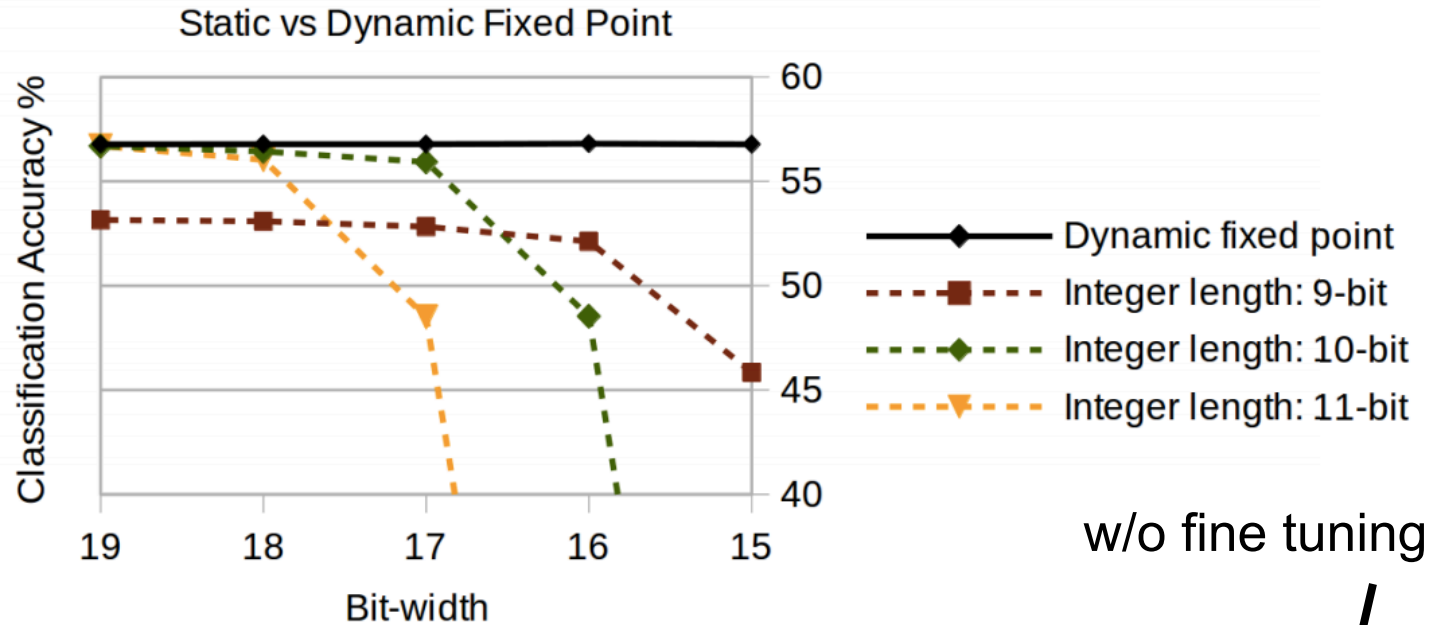


Dynamic range = 0.125

Image Source:
Moons et al,
WACV 2016

Impact on Accuracy

Top-1 accuracy
of CaffeNet
on ImageNet



| | Layer outputs | CONV parameters | FC parameters | 32-bit floating point baseline | Fixed point accuracy |
|------------------|---------------|-----------------|---------------|--------------------------------|----------------------|
| LeNet (Exp 1) | 4-bit | 4-bit | 4-bit | 99.1% | 99.0% (98.7%) |
| LeNet (Exp 2) | 4-bit | 2-bit | 2-bit | 99.1% | 98.8% (98.0%) |
| Full CIFAR-10 | 8-bit | 8-bit | 8-bit | 81.7% | 81.4% (80.6%) |
| SqueezeNet top-1 | 8-bit | 8-bit | 8-bit | 57.7% | 57.1% (55.2%) |
| CaffeNet top-1 | 8-bit | 8-bit | 8-bit | 56.9% | 56.0% (55.8%) |
| GoogLeNet top-1 | 8-bit | 8-bit | 8-bit | 68.9% | 66.6% (66.1%) |

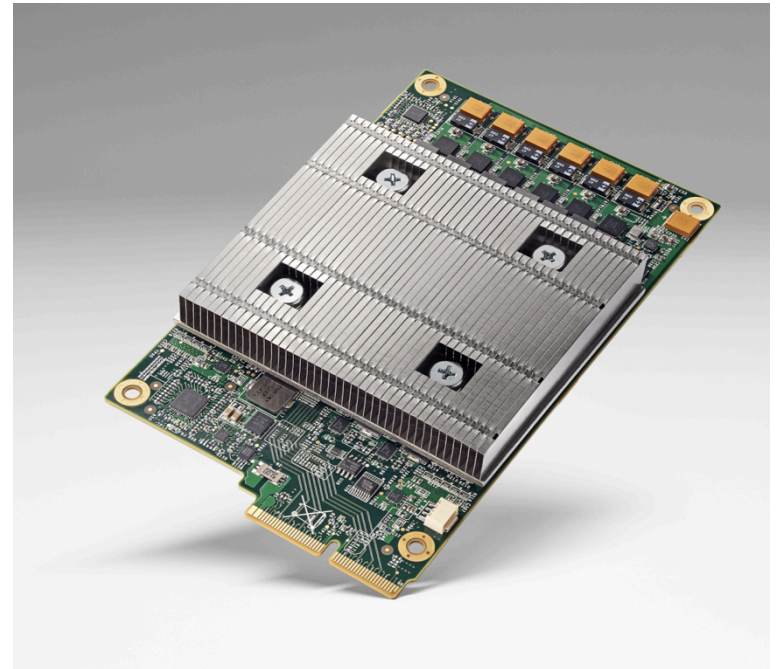
Google's Tensor Processing Unit (TPU)

“ With its TPU Google has seemingly focused on delivering the data really quickly by cutting down on precision. Specifically, it doesn't rely on floating point precision like a GPU

....

Instead the chip uses integer math...TPU used 8-bit integer.”

- Next Platform (May 19, 2016)



Nvidia PASCAL

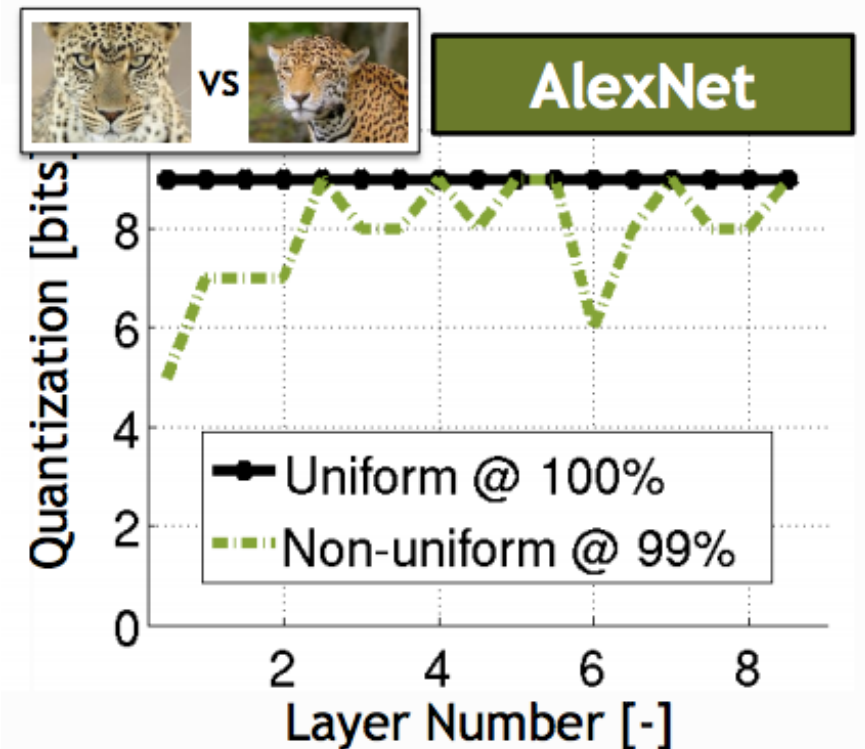
“New half-precision, **16-bit floating point instructions** deliver over **21 TeraFLOPS** for unprecedented training performance. **With 47 TOPS (tera-operations per second) of performance, new 8-bit integer instructions** in Pascal allow AI algorithms to deliver real-time responsiveness for deep learning inference.”

– Nvidia.com (April 2016)



Precision Varies from Layer to Layer

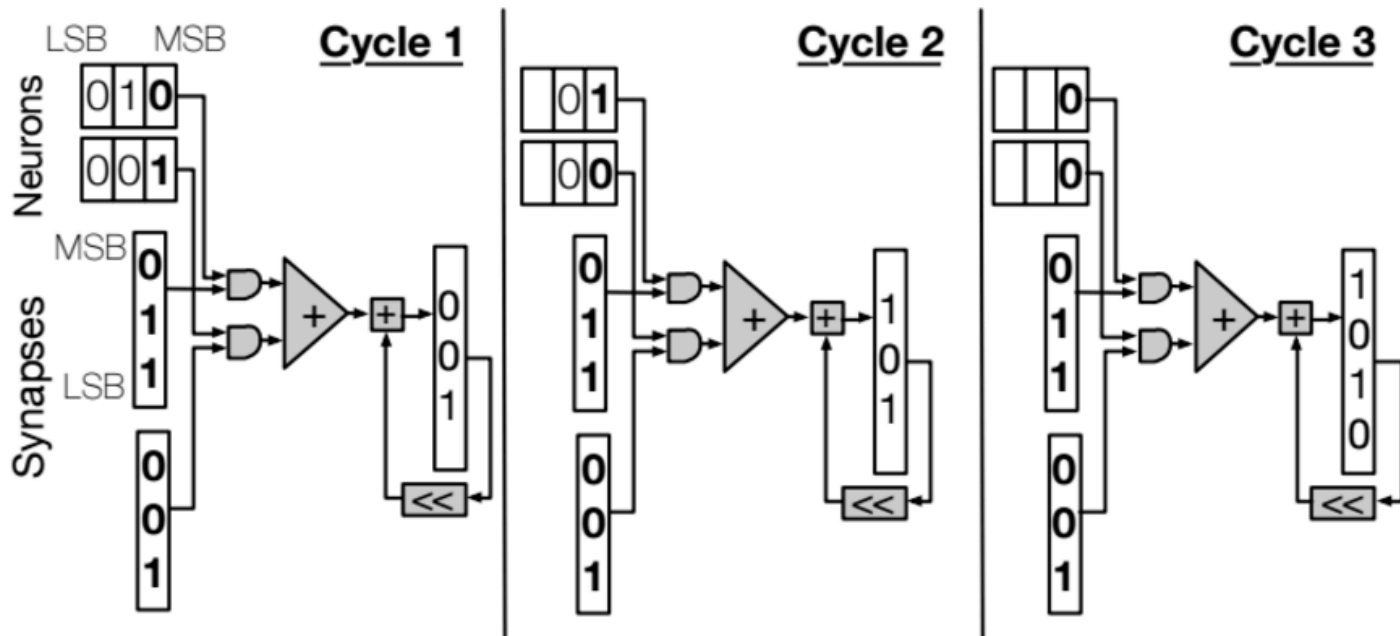
| Tolerance | Bits per layer (I+F) |
|----------------------|----------------------|
| AlexNet (F=0) | |
| 1% | 10-8-8-8-8-8-6-4 |
| 2% | 10-8-8-8-8-8-5-4 |
| 5% | 10-8-8-8-7-7-5-3 |
| 10% | 9-8-8-8-7-7-5-3 |



Bitwidth Scaling (Speed)

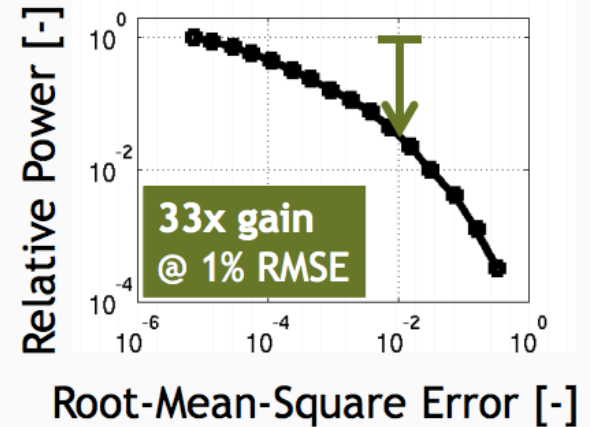
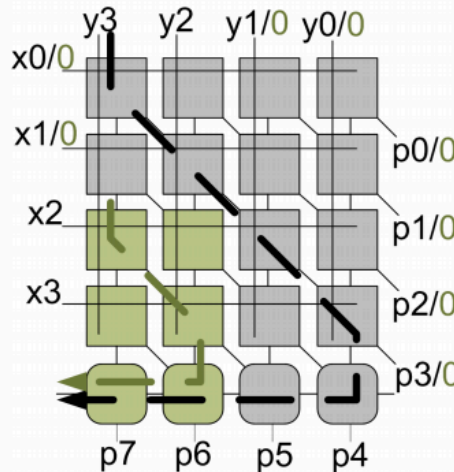
Bit-Serial Processing: Reduce Bit-width → Skip Cycles
Speed up of 2.24x vs. 16-bit fixed

$$\sum_{i=0}^{N_i-1} s_i \times n_i = \sum_{i=0}^{N_i-1} s_i \times \sum_{b=0}^{P-1} n_i^b \times 2^b = \sum_{b=0}^{P-1} 2^b \times \sum_{i=0}^{N_i-1} n_i^b \times S_i$$



Bitwidth Scaling (Power)

Reduce Bit-width →
Shorter Critical Path
→ Reduce Voltage



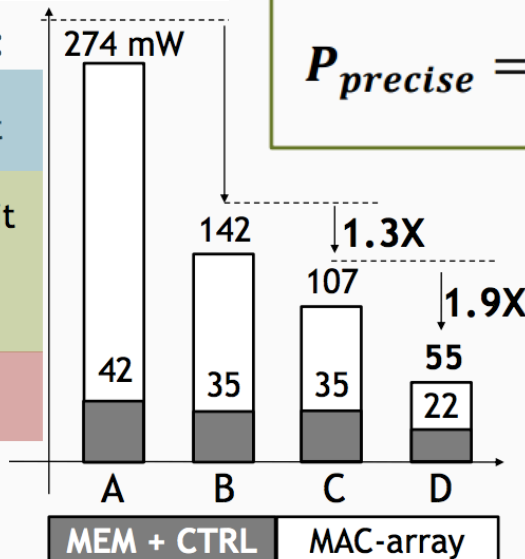
AlexNet Layer 2 example:

A. 2D-baseline @ 16 bit

B. Precision-Scaling @ 7-7 bit

C. Voltage-Scaling @ 0.9 V

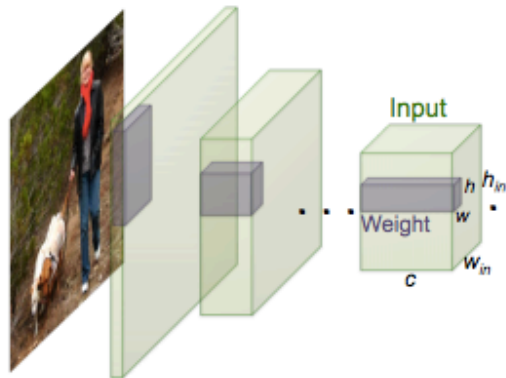
D. Sparse operation guarding



$$P_{precise} = \alpha C f V^2 \Rightarrow P_{imprecise} = \frac{\alpha}{k_1} C f \left(\frac{V}{k_2}\right)^2$$

Power reduction of
2.56x vs. 16-bit fixed
On AlexNet Layer 2

Binary Nets



| | Network Variations | Operations used in Convolution | Memory Saving (Inference) | Computation Saving (Inference) | Accuracy on ImageNet (AlexNet) |
|--------------------------------------|---|--------------------------------|---------------------------|--------------------------------|--------------------------------|
| Standard Convolution | Real-Value Inputs $\begin{bmatrix} 0.11 & -0.21 & \dots & -0.34 \\ -0.25 & 0.61 & \dots & 0.52 \end{bmatrix}$ Real-Value Weights $\begin{bmatrix} 0.12 & -1.2 & \dots & 0.41 \\ -0.2 & 0.5 & \dots & 0.68 \end{bmatrix}$ | $+, -, \times$ | 1x | 1x | %56.7 |
| Binary Weight | Real-Value Inputs $\begin{bmatrix} 0.11 & -0.21 & \dots & -0.34 \\ -0.25 & 0.61 & \dots & 0.52 \end{bmatrix}$ Binary Weights $\begin{bmatrix} 1 & -1 & \dots & 1 \\ -1 & 1 & \dots & 1 \end{bmatrix}$ | $+, -$ | $\sim 32x$ | $\sim 2x$ | %56.8 |
| BinaryWeight Binary Input (XNOR-Net) | Binary Inputs $\begin{bmatrix} 1 & -1 & \dots & -1 \\ -1 & 1 & \dots & 1 \end{bmatrix}$ Binary Weights $\begin{bmatrix} 1 & -1 & \dots & 1 \\ -1 & 1 & \dots & 1 \end{bmatrix}$ | XNOR, bitcount | $\sim 32x$ | $\sim 58x$ | %44.2 |

Classification Accuracy(%)

| Binary-Weight | | Binary-Input-Binary-Weight | | | | Full-Precision | | | |
|---------------|-------------|----------------------------|-------|-------------|-------------|----------------|-------|------------|-------|
| BWN | | BC[11] | | XNOR-Net | | BNN[11] | | AlexNet[1] | |
| Top-1 | Top-5 | Top-1 | Top-5 | Top-1 | Top-5 | Top-1 | Top-5 | Top-1 | Top-5 |
| 56.8 | 79.4 | 35.4 | 61.0 | 44.2 | 69.2 | 27.9 | 50.42 | 56.6 | 80.2 |

BinaryConnect (BC) = [Courbariaux et al., ArXiv 2015]

Binary Neural Networks (BNN) = [Courbariaux et al., ArXiv 2016]

[Rastegari et al., BWN & XNOR-Net, ECCV 2016]

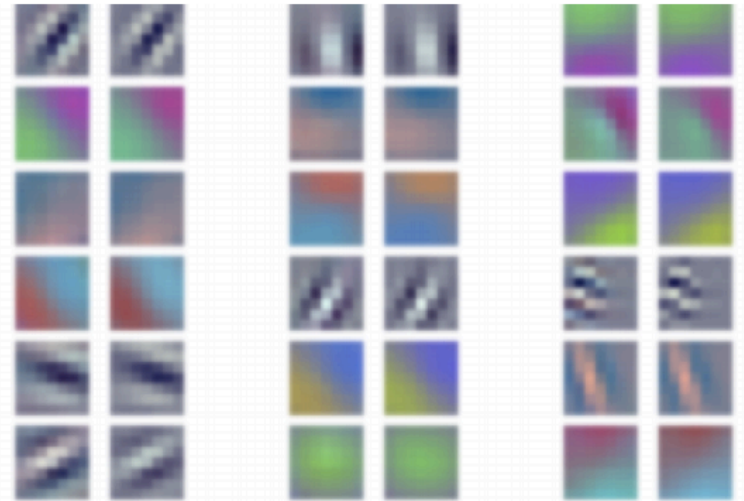
Reduce Number of Ops and Weights

- **Network Compression**
 - Low Rank Approximation
 - Weight Sharing and Vector Quantization
- **Pruning**
 - Weights
 - Activations
- **Network Architectures**

Low Rank Approximation

- **Low Rank approximation**

- Tensor decomposition based on singular value decomposition (SVD)
- Filter Clustering with modified K-means
- Fine Tuning



- **Speed up by 1.6 – 2.7x** on CPU/GPU for CONV1, CONV2 layers
- **Reduce size by 5 - 13x** for FC layer
- **< 1% drop in accuracy**

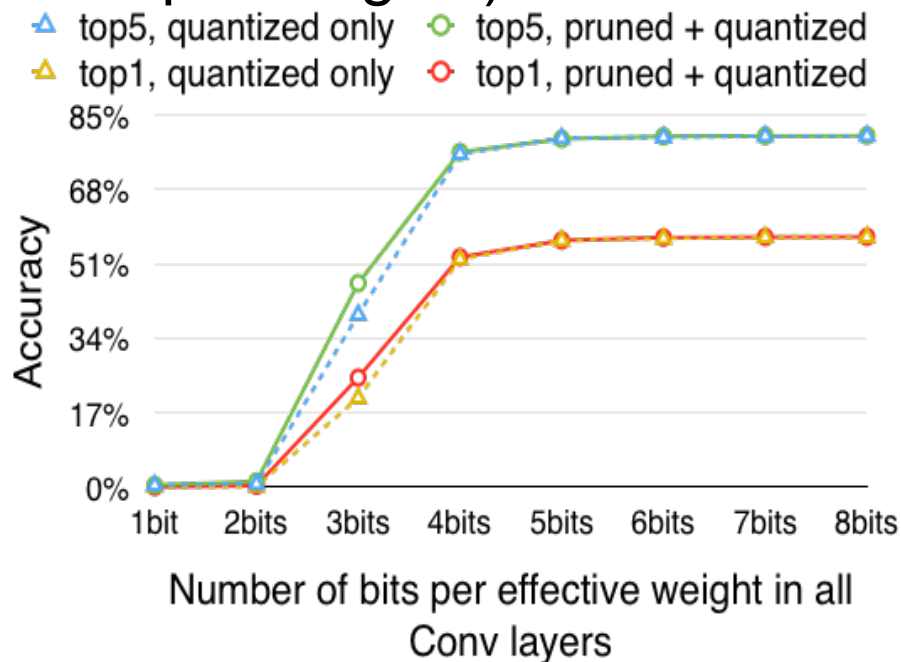
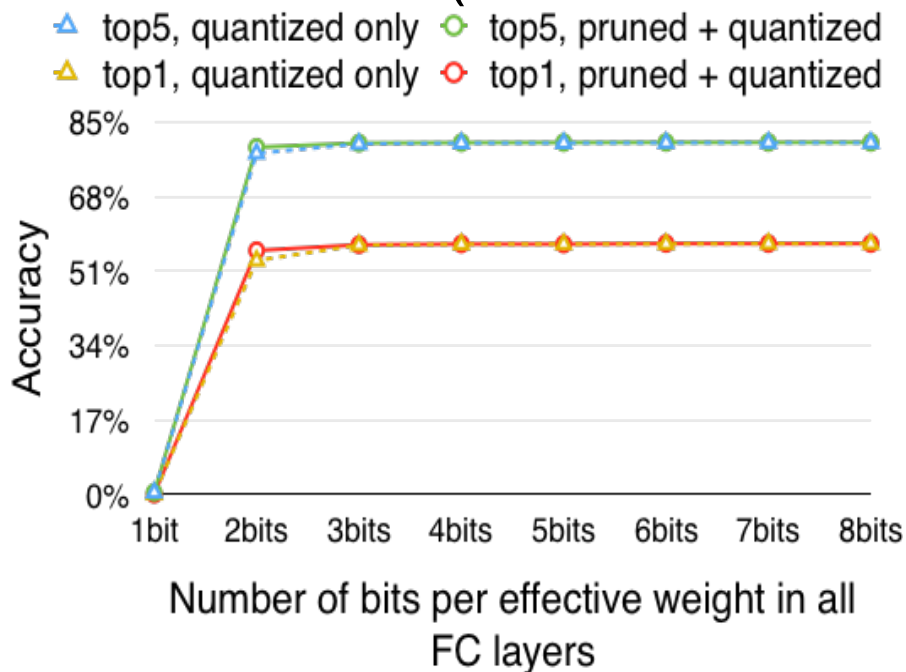
Low Rank Approximation on Phone

- Rank selection per Layer
- Tucker Decomposition (extension of SVD)
- Fine tuning

| Model | Top-5 | Weights | FLOPs | S6 | | Titan X |
|-----------------------------|------------------|-----------------|------------------|------------------|-------------------|-------------------|
| <i>AlexNet</i> | 80.03 | 61M | 725M | 117ms | 245mJ | 0.54ms |
| <i>AlexNet*</i> (imp.) | 78.33 (-1.70) | 11M (×5.46) | 272M (×2.67) | 43ms (×2.72) | 72mJ (×3.41) | 0.30ms (×1.81) |
| <i>VGG-S</i> | 84.60 | 103M | 2640M | 357ms | 825mJ | 1.86ms |
| <i>VGG-S*</i> (imp.) | 84.05 (-0.55) | 14M (×7.40) | 549M (×4.80) | 97ms (×3.68) | 193mJ (×4.26) | 0.92ms (×2.01) |
| <i>GoogLeNet</i> | 88.90 | 6.9M | 1566M | 273ms | 473mJ | 1.83ms |
| <i>GoogLeNet*</i> (imp.) | 88.66 (-0.24) | 4.7M (×1.28) | 760M (×2.06) | 192ms (×1.42) | 296mJ (×1.60) | 1.48ms (×1.23) |
| <i>VGG-16</i> | 89.90 | 138M | 15484M | 1926ms | 4757mJ | 10.67ms |
| <i>VGG-16*</i> (imp.) | 89.40 (-0.50) | 127M (×1.09) | 3139M (×4.93) | 576ms (×3.34) | 1346mJ (×3.53) | 4.58ms (×2.33) |

Weight Sharing + Vector Quantization

Trained Quantization: Weight Sharing via K-means clustering
(reduce number of unique weights)



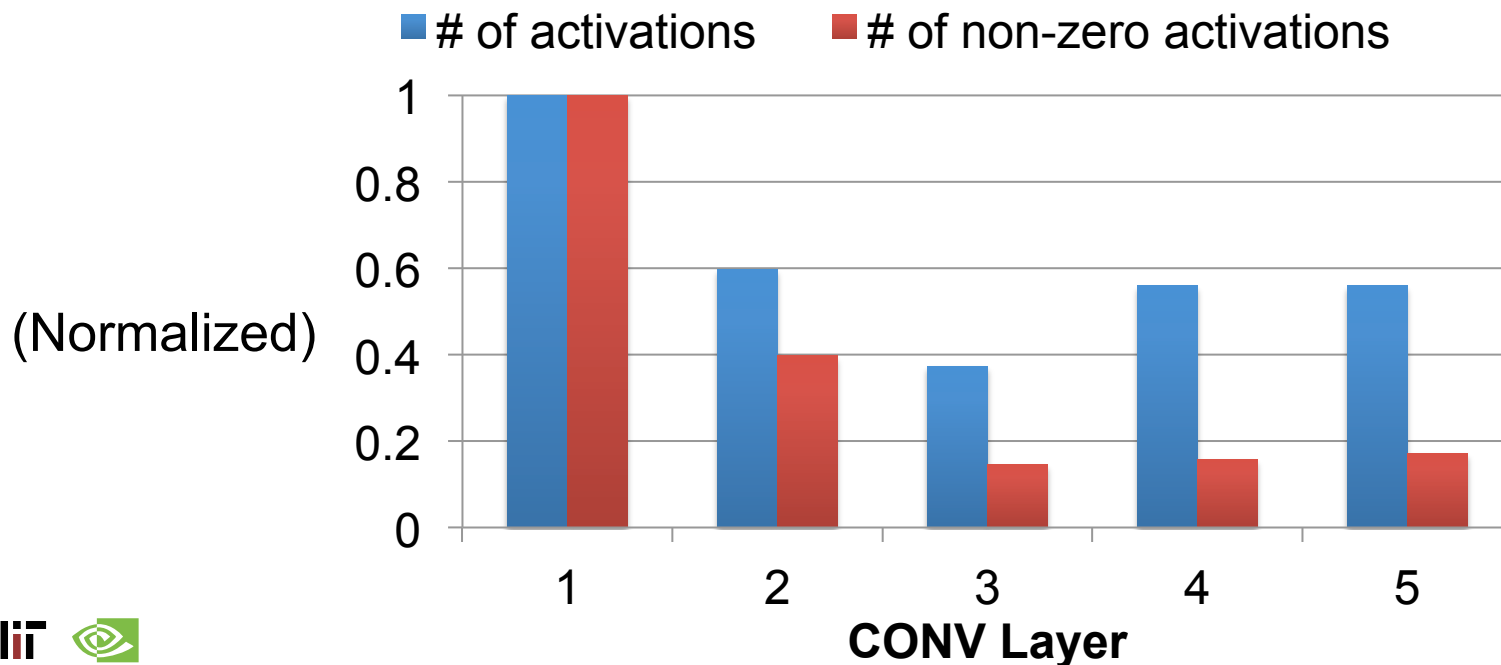
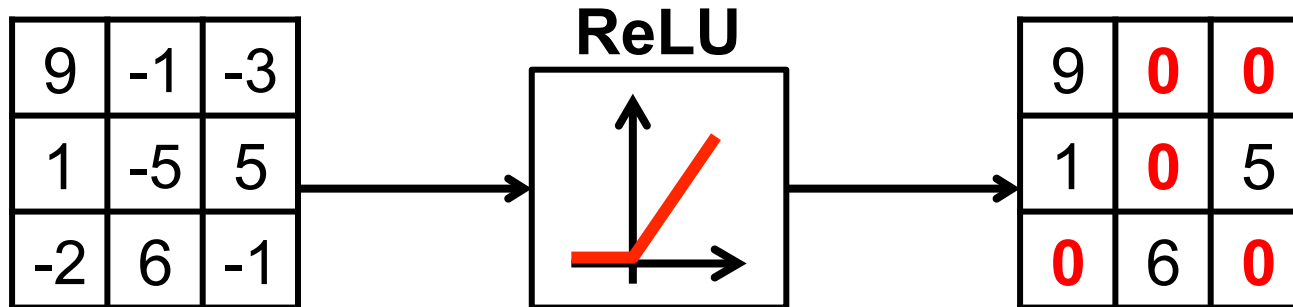
Reduce Bits for Storage (compute still 16-bits)



Exploit Data Statistics

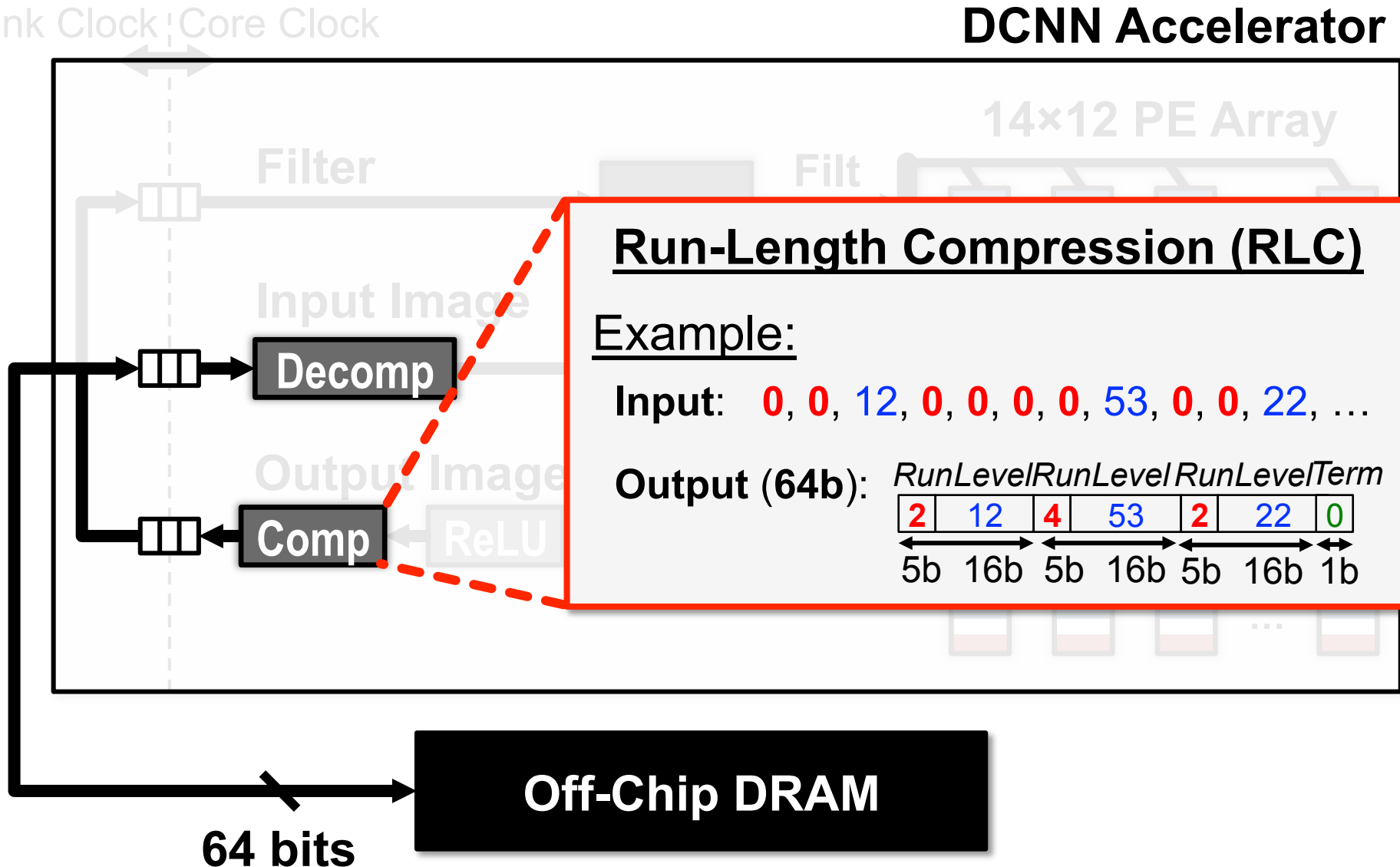
Sparsity in Fmaps

Many **zeros** in output fmaps after ReLU



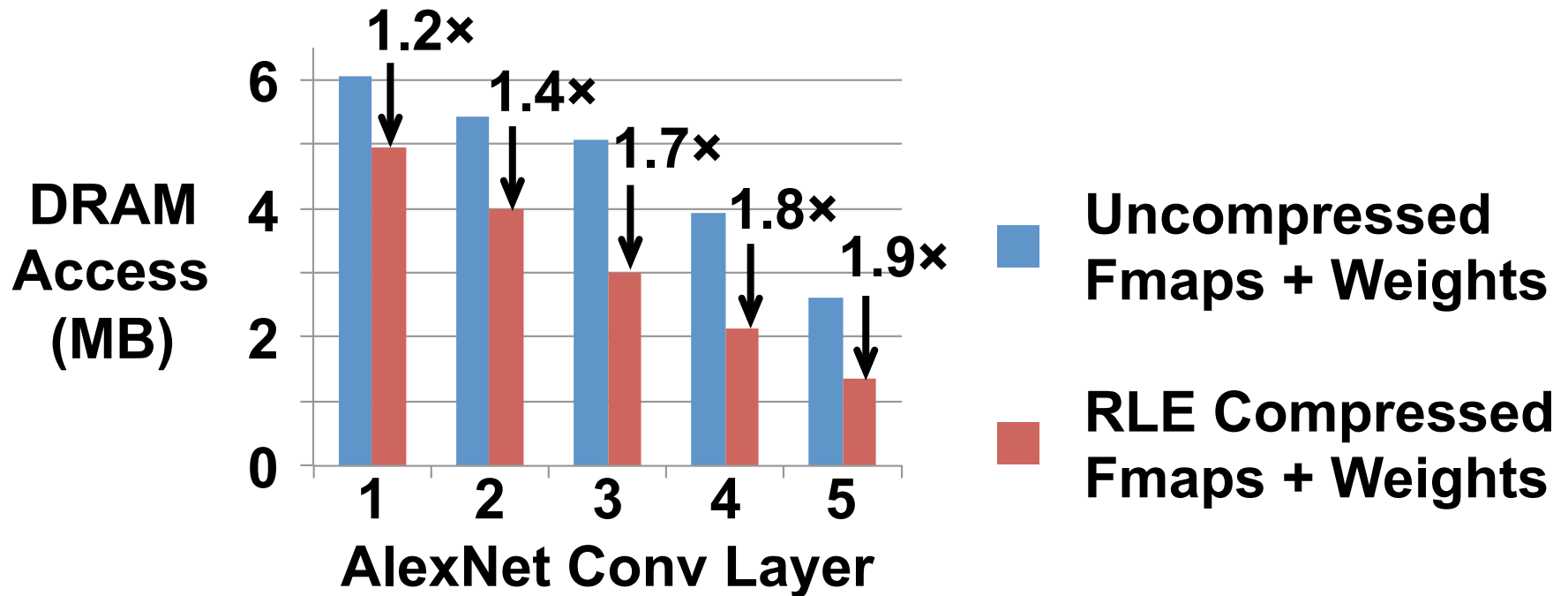
I/O Compression in Eyeriss

DCNN Accelerator



[Chen et al., ISSCC 2016]

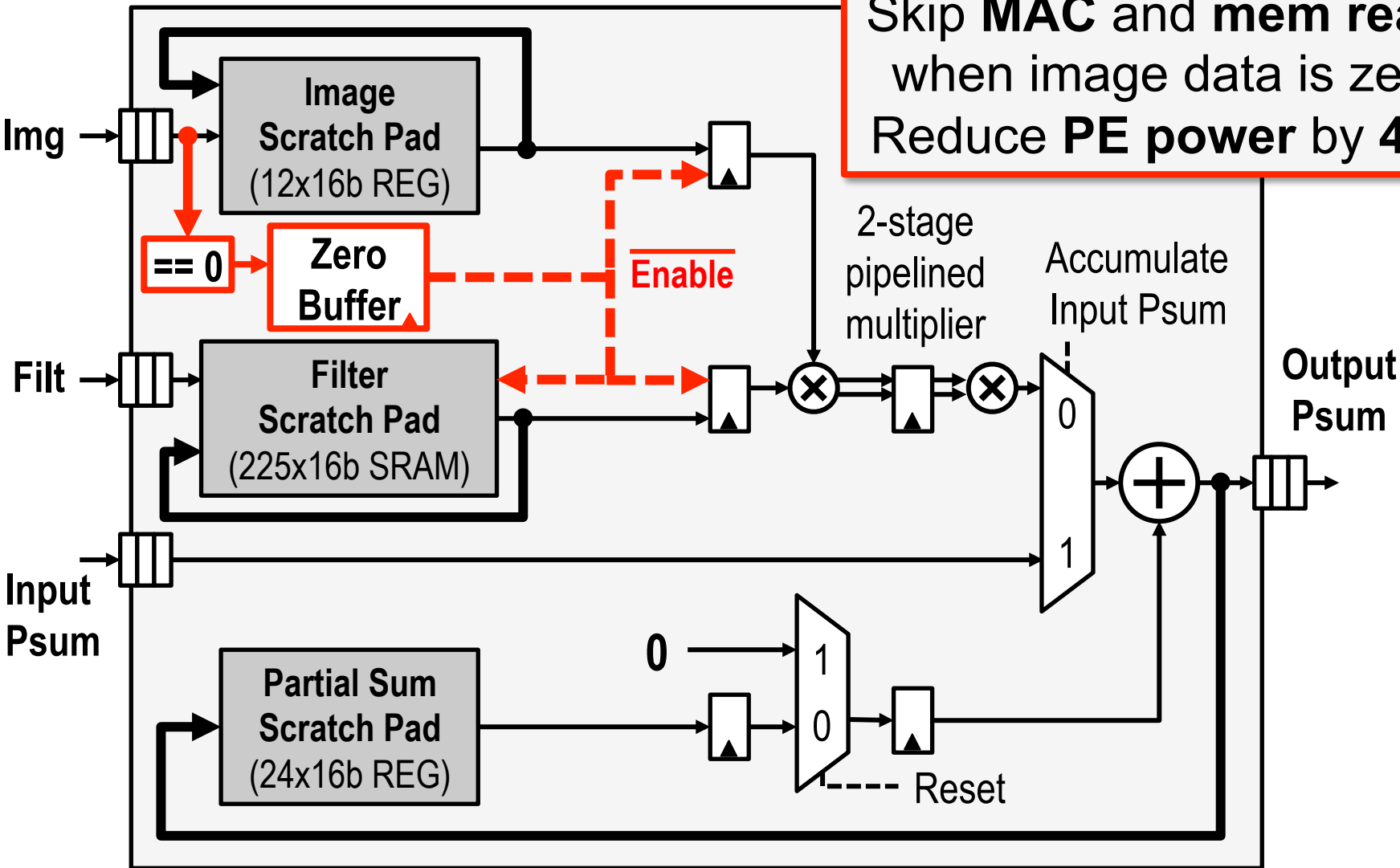
Compression Reduces DRAM BW



Simple RLC within 5% - 10% of theoretical entropy limit

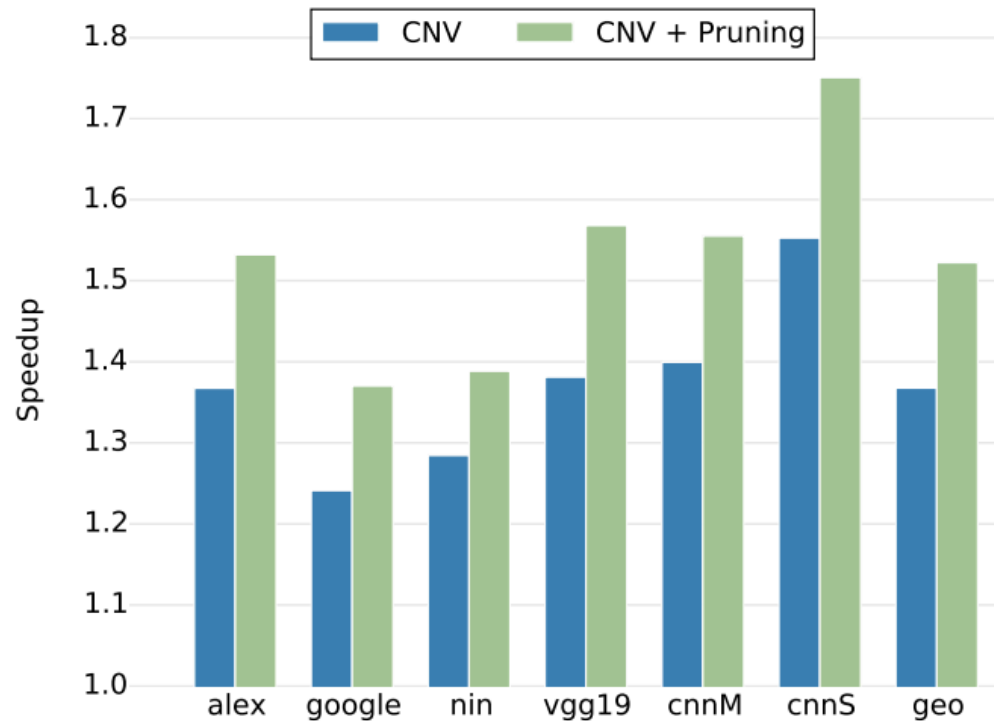
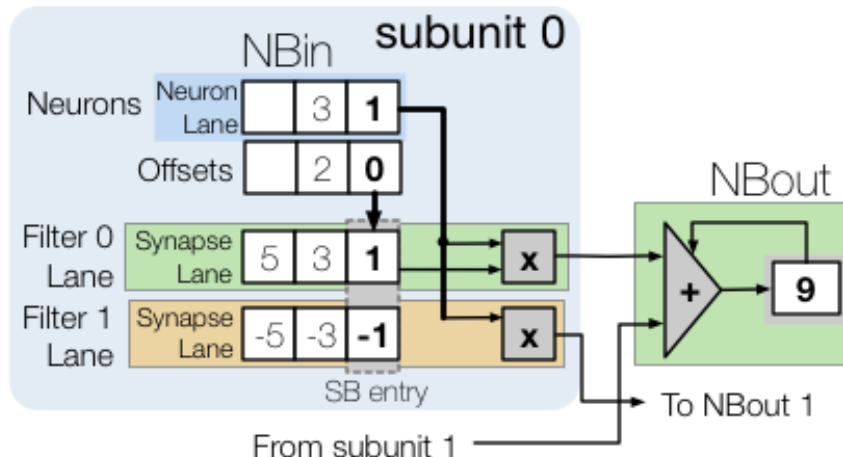
Data Gating / Zero Skipping in Eyeriss

Skip **MAC** and mem reads when image data is zero. Reduce **PE power by 45%**



Cnvlutin

- Process Convolution Layers
- Built on top of DaDianNao (4.49% area overhead)
- Speed up of 1.37x (1.52x with activation pruning)

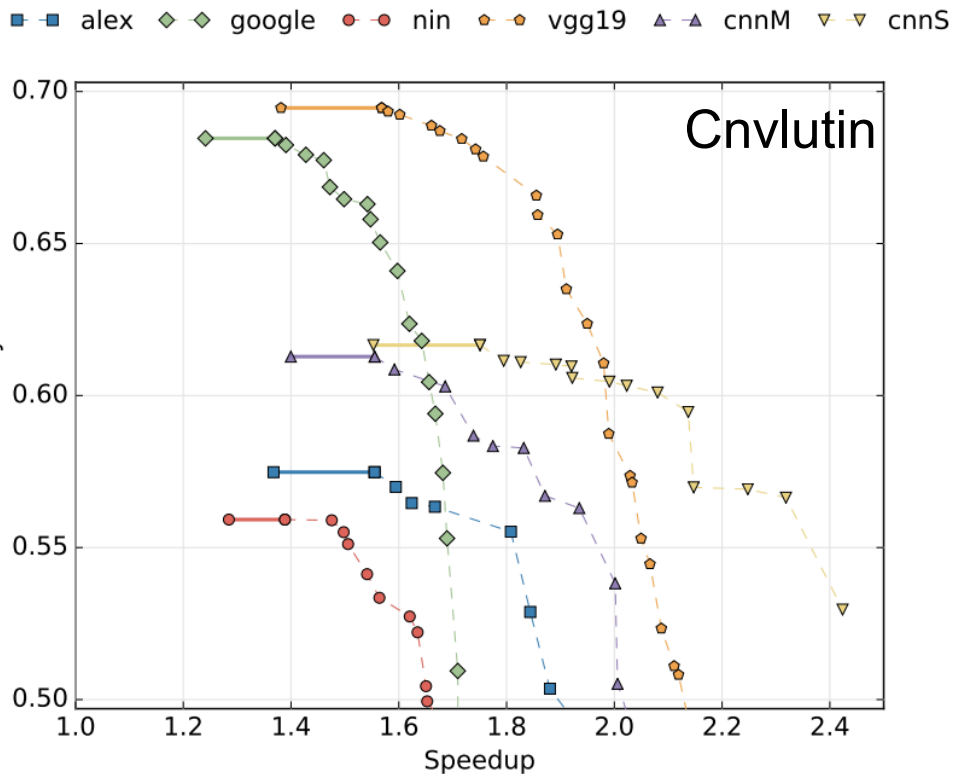


Pruning Activations

Remove small activation values

Speed up 11% (ImageNet)

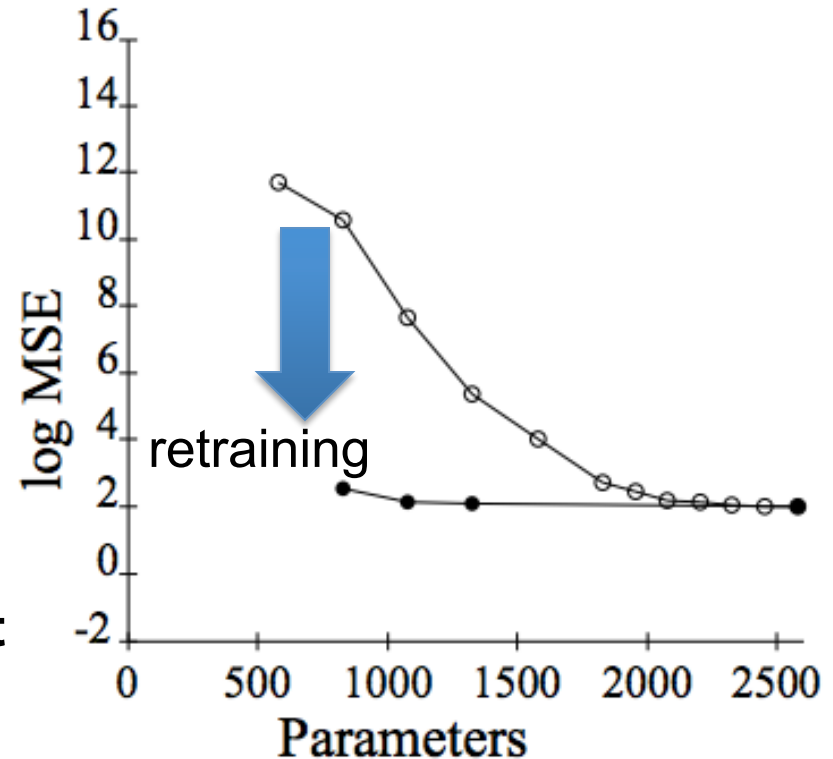
Reduce power 2x (MNIST)



Pruning – Make Weights Sparse

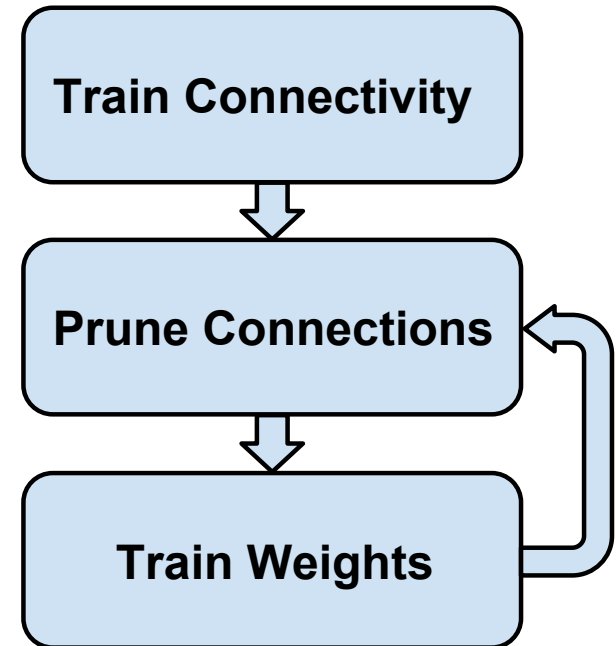
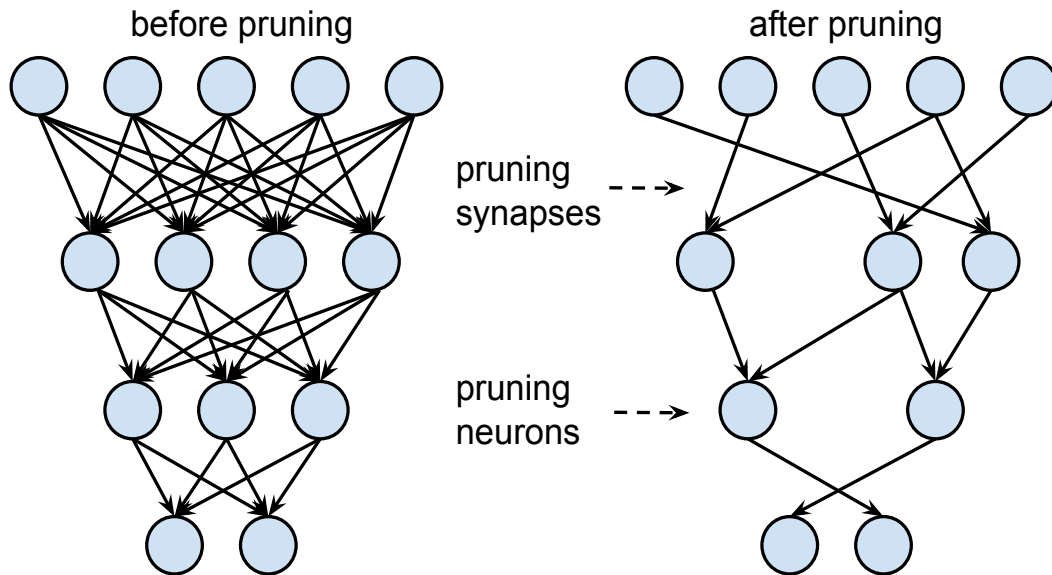
- **Optimal Brain Damage**

1. Choose a reasonable network architecture
2. Train network until reasonable solution obtained
3. Compute the second derivative for each weight
4. Compute saliencies (i.e. impact on training error) for each weight
5. Sort weights by saliency and delete low-saliency weights
6. Iterate to step 2



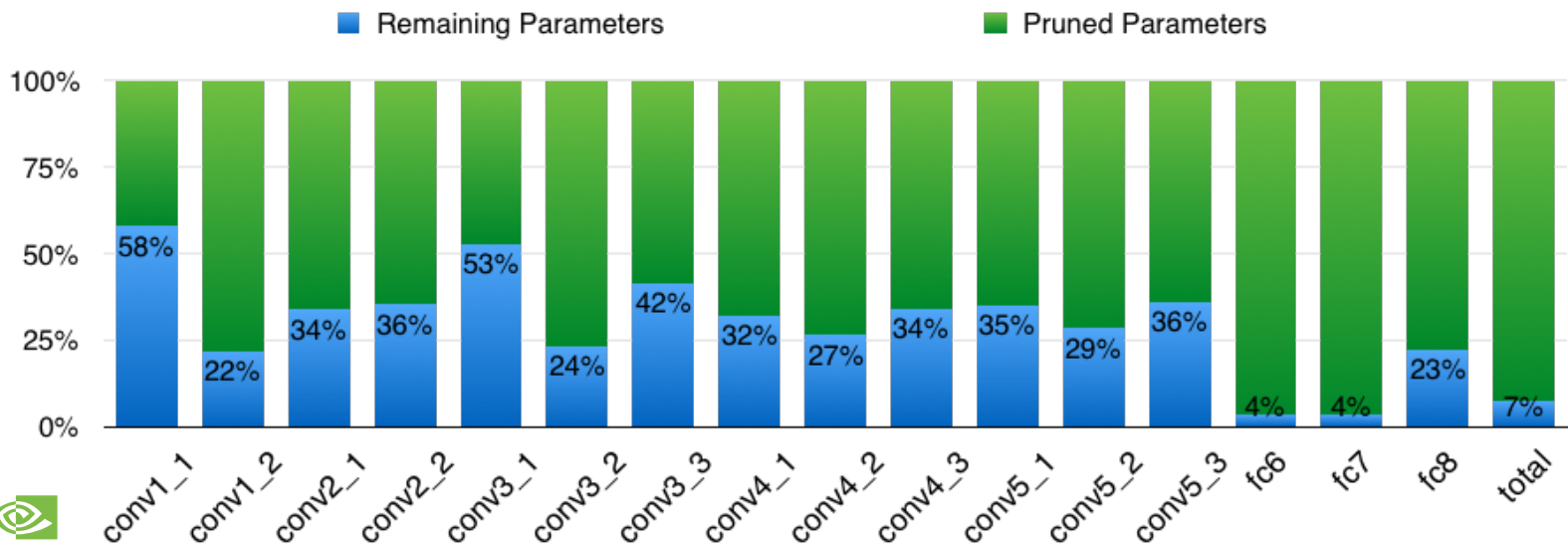
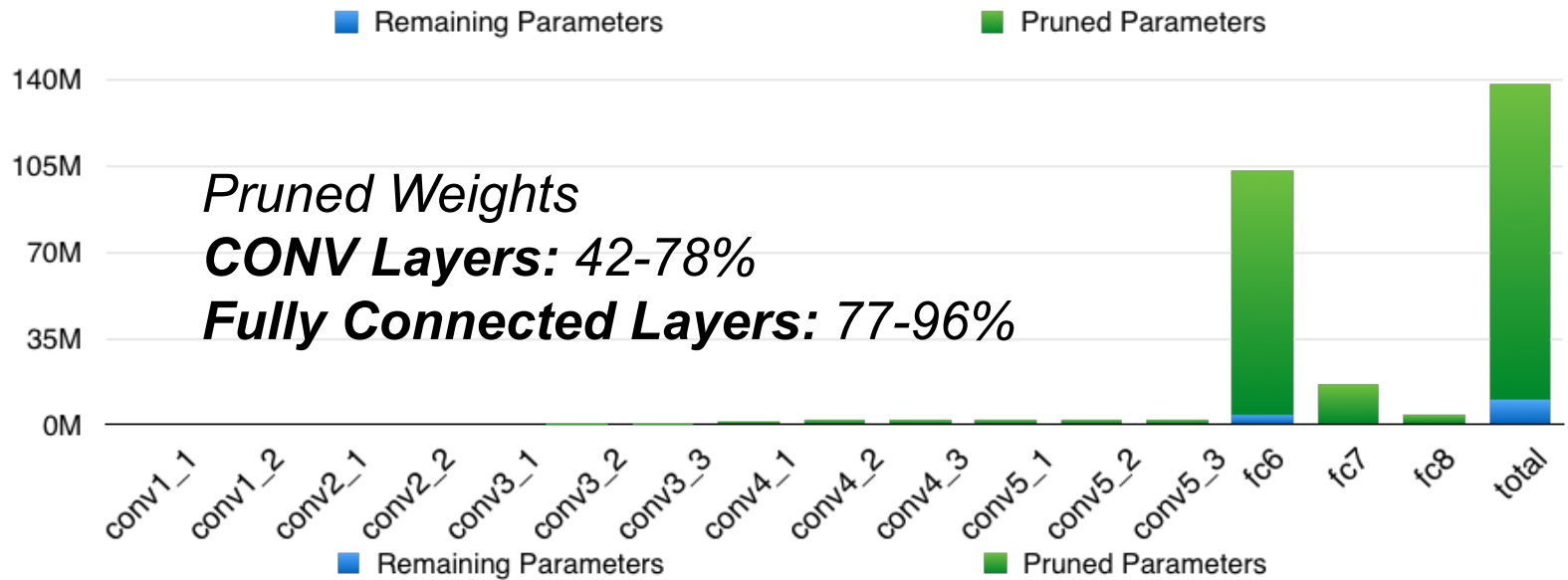
Pruning – Make Weights Sparse

Prune based on magnitude of weights



Pruning of VGG-16

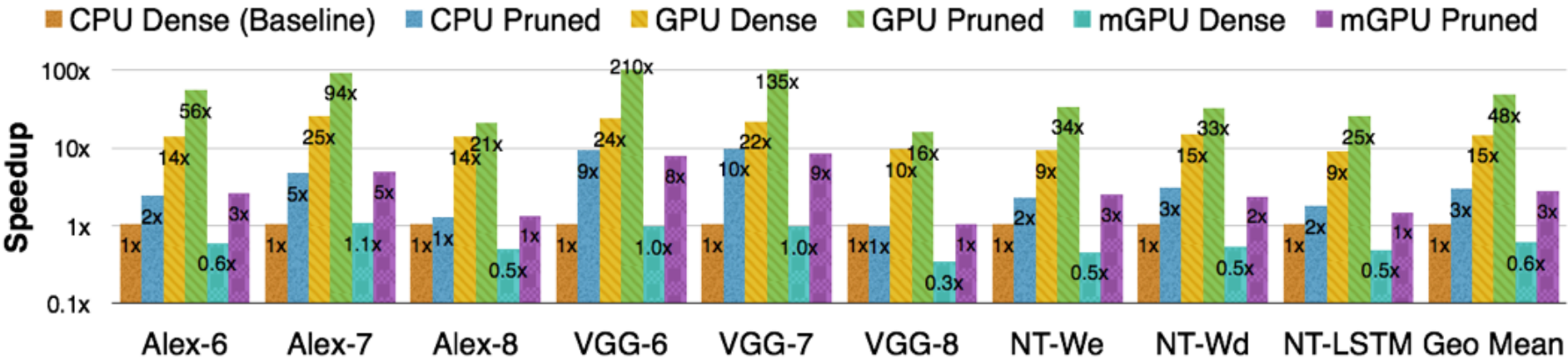
Pruning has most impact on Fully Connected Layers



Speed up of Weight Pruning on CPU/GPU

On Fully Connected Layers

Average Speed up of 3.2x on GPU, 3x on CPU, 5x on mGPU



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

Batch size = 1



Energy-Aware Pruning

- **# of Weights alone is not a good metric for energy**
 - **Example (AlexNet):**
 - **# of Weights (FC Layer) > # of Weights (CONV layer)**
 - **Energy (FC Layer) < Energy (CONV layer)**
- **Use energy evaluation method to estimate DNN energy**
 - **Account for data movement**
- **Prune based on energy rather than weights**
 - **Reduce overall energy (ALL layers) by 3.7x for AlexNet**
 - **1.8x more efficient than previous magnitude-based approach**
 - **1.6x energy reduction for GoogleNet**

Compression of Weights & Activations

- Compress weights and fmaps between DRAM and accelerator
- Variable Length / Huffman Coding

Example:

Value: **16'b0** → Compressed Code: {**1'b0**}

Value: **16'bx** → Compressed Code: {**1'b1**, **16'bx**}

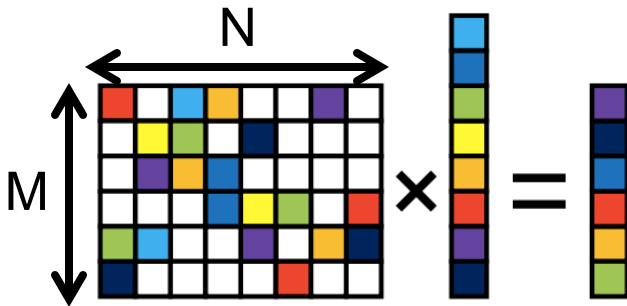
- Tested on AlexNet → **2× overall BW Reduction**

| Layer | Filter / Image bits (0%) | Filter / Image BW Reduc. | IO / HuffIO (MB/frame) | Voltage (V) | MMACs/ Frame | Power (mW) | Real (TOPS/W) |
|--------------|--------------------------|--------------------------|------------------------|-------------|--------------|------------|---------------|
| General CNN | 16 (0%) / 16 (0%) | 1.0x | | 1.1 | — | 288 | 0.3 |
| AlexNet 11 | 7 (21%) / 4 (29%) | 1.17x / 1.3x | 1 / 0.77 | 0.85 | 105 | 85 | 0.96 |
| AlexNet 12 | 7 (19%) / 7 (89%) | 1.15x / 5.8x | 3.2 / 1.1 | 0.9 | 224 | 55 | 1.4 |
| AlexNet 13 | 8 (11%) / 9 (82%) | 1.05x / 4.1x | 6.5 / 2.8 | 0.92 | 150 | 77 | 0.7 |
| AlexNet 14 | 9 (04%) / 8 (72%) | 1.00x / 2.9x | 5.4 / 3.2 | 0.92 | 112 | 95 | 0.56 |
| AlexNet 15 | 9 (04%) / 8 (72%) | 1.00x / 2.9x | 3.7 / 2.1 | 0.92 | 75 | 95 | 0.56 |
| Total / avg. | — | — | 19.8 / 10 | — | — | 76 | 0.94 |
| LeNet-5 11 | 3 (35%) / 1 (87%) | 1.40x / 5.2x | 0.005 / 0.001 | 0.7 | 0.3 | 25 | 1.07 |
| LeNet-5 12 | 4 (26%) / 6 (55%) | 1.25x / 1.9x | 0.050 / 0.042 | 0.8 | 1.6 | 35 | 1.75 |
| Total / avg. | — | — | 0.053 / 0.043 | — | — | 33 | 1.6 |

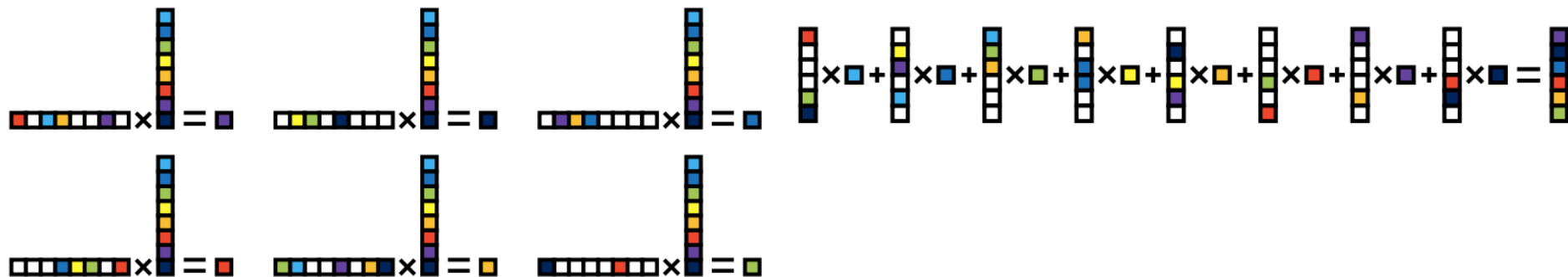
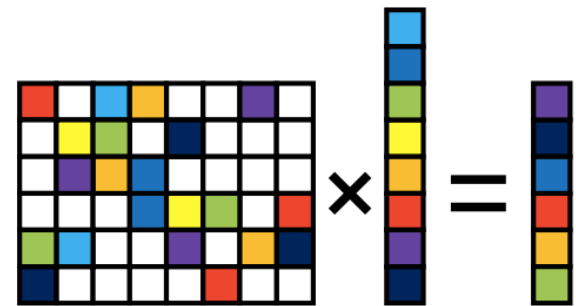
Sparse Matrix-Vector DSP

- Use **CSC** rather than **CSR** for **SpMxV**

Compressed Sparse Row (CSR)



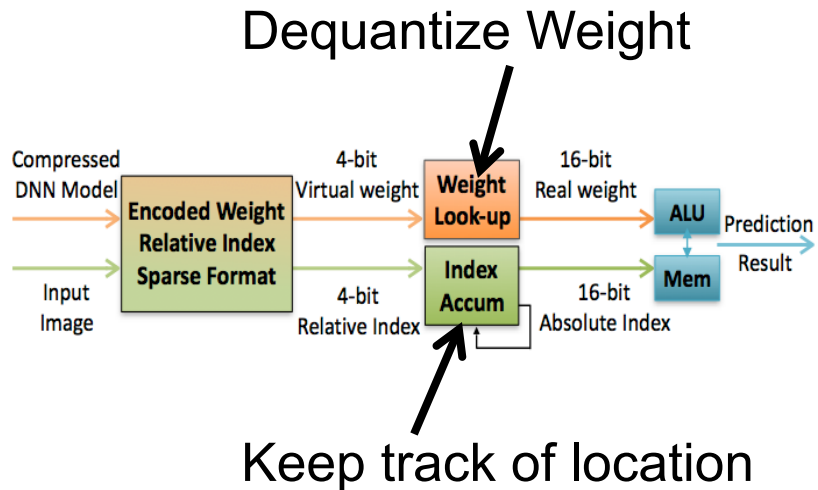
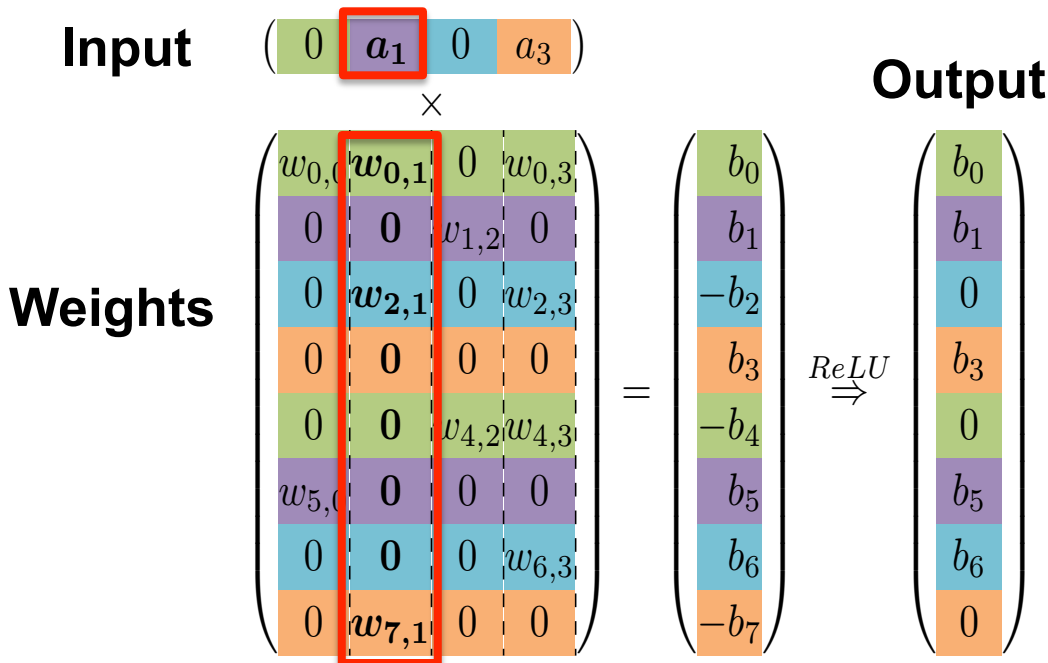
Compressed Sparse Column (CSC)



Reduce memory bandwidth by 2x (when not $M \gg N$)

EIE: A Sparse Linear Algebra Engine

- Process Fully Connected Layers (after Deep Compression)
- Store weights column-wise in Run Length format
 - Non-zero weights, Run-length of zeros
 - Start location of each column since variable length
- Read relative column when input is non-zero



Network Architecture

Reduce size and computation with 1x1 Filter

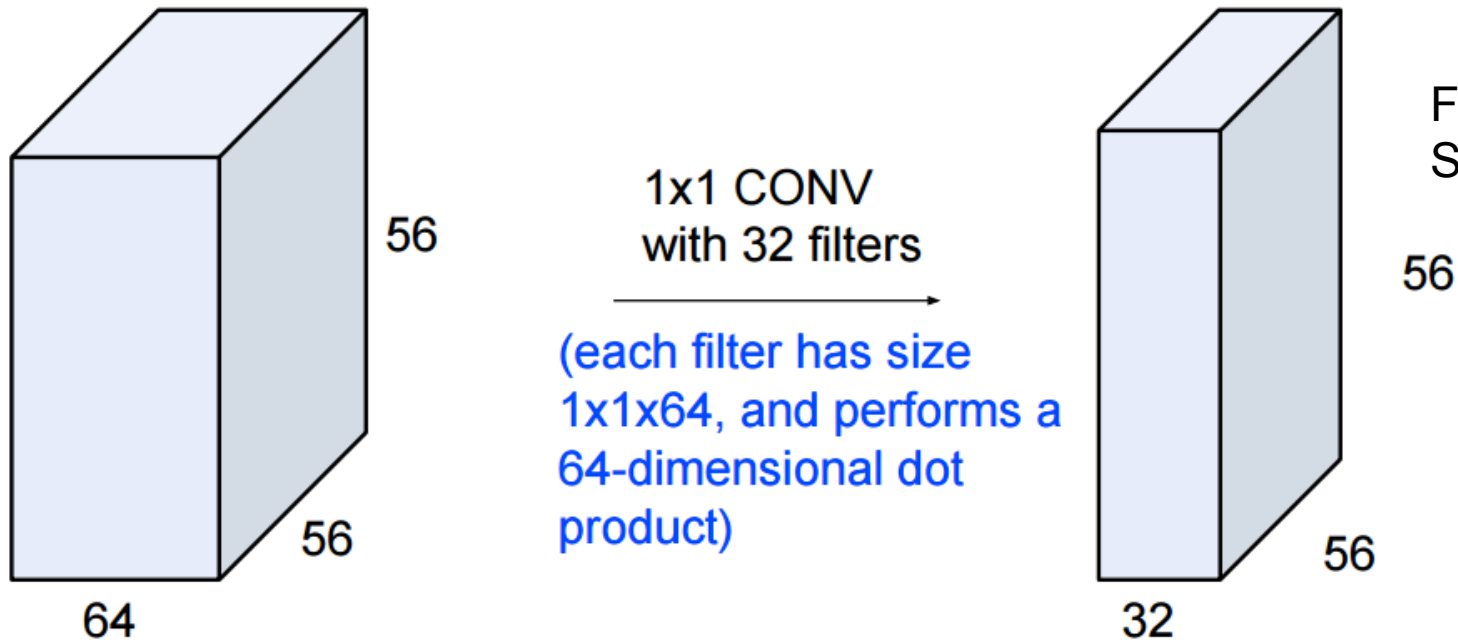


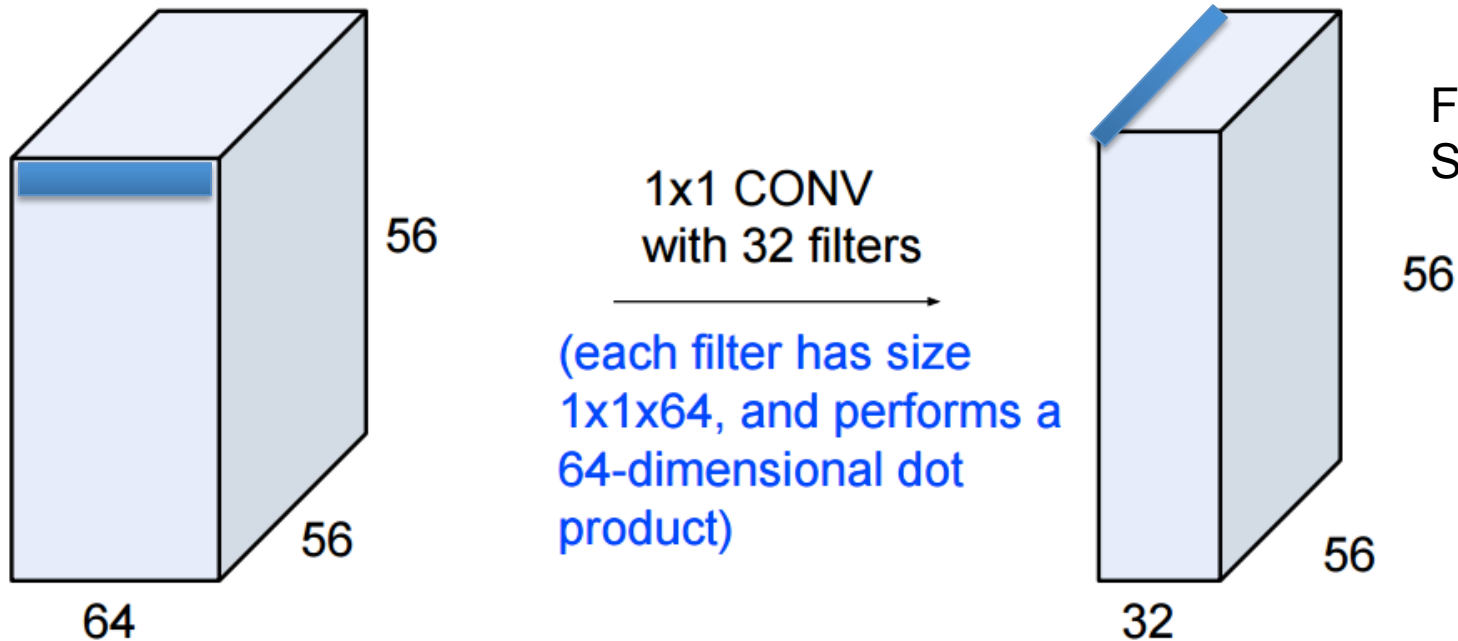
Figure Source:
Stanford cs231n

Used in Network In Network(NiN) and GoogLeNet

[Lin et al., ArXiv 2013 / ICLR 2014] [Szegedy et al., ArXiv 2014 / CVPR 2015]

Network Architecture

Reduce size and computation with 1x1 Filter



Used in Network In Network(NiN) and GoogLeNet

[Lin et al., ArXiv 2013 / ICLR 2014] [Szegedy et al., ArXiv 2014 / CVPR 2015]

Network Architecture

Reduce size and computation with 1x1 Filter

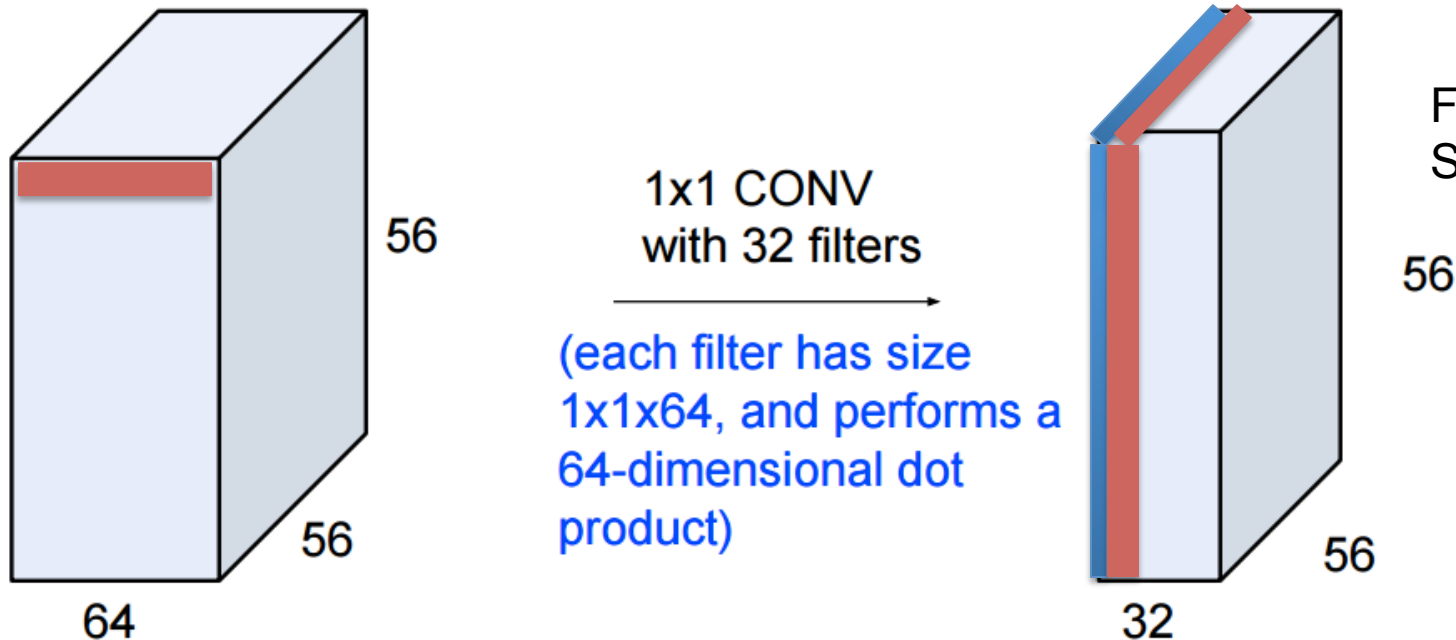


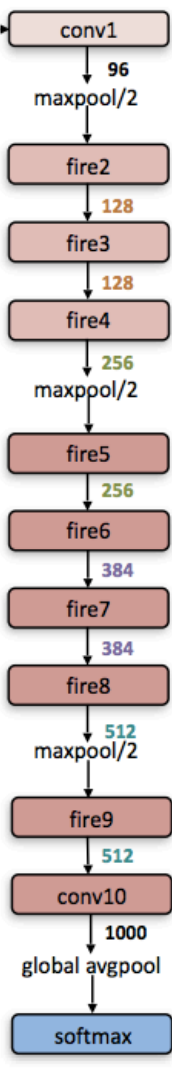
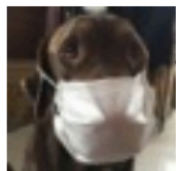
Figure Source:
Stanford cs231n

Used in Network In Network(NiN) and GoogLeNet

[Lin et al., ArXiv 2013 / ICLR 2014] [Szegedy et al., ArXiv 2014 / CVPR 2015]

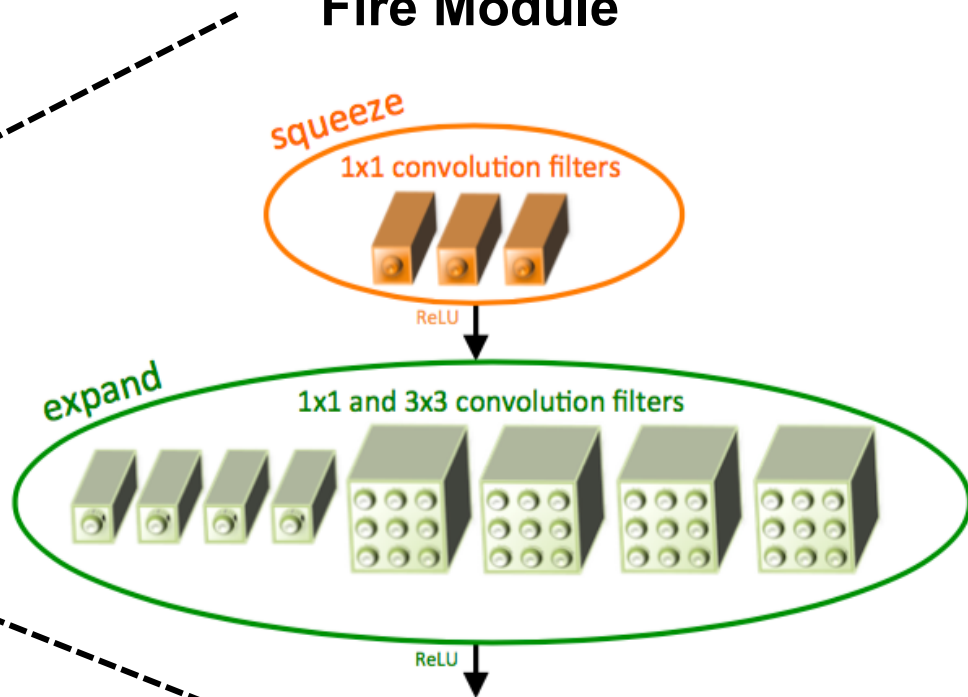
SqueezeNet

Reduce weights by reducing number of input channels by “squeezing” with 1x1
50x fewer weights than AlexNet



"labrador retriever dog"

Fire Module



Energy Consumption of Existing DNNs

- Maximally reducing # of weights does not necessarily result in optimized energy consumption
- Deeper CNNs with fewer weights (e.g. GoogLeNet, SqueezeNet), do not necessarily consume less energy than shallower CNNs with more weights (e.g. AlexNet)
- Reducing # of weights can provide equal or more reduction than reducing the bitwidth of weights (e.g. BWN)

