# DNN Model and Hardware Co-Design
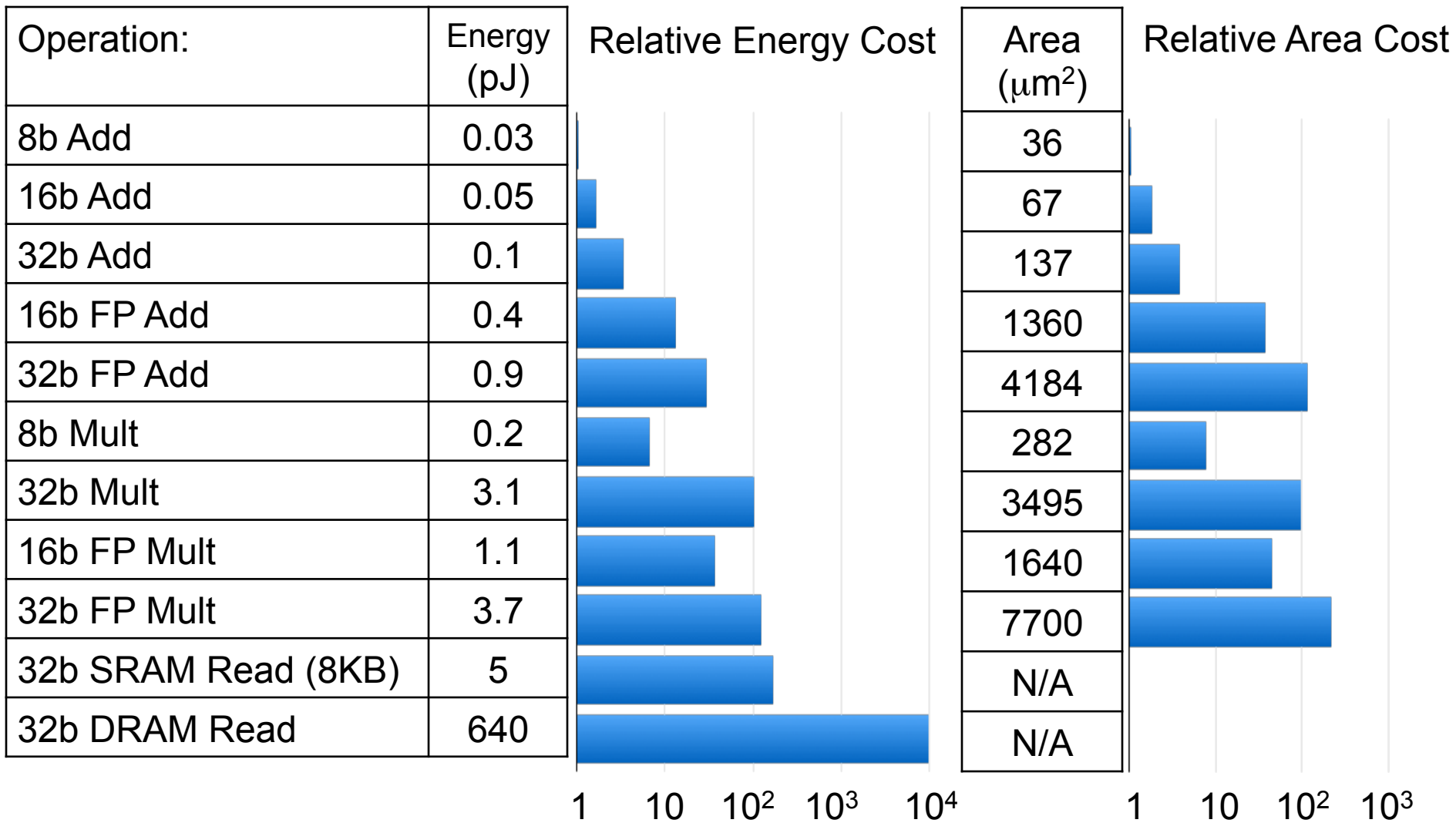
## CICS/MTL Tutorial (2017)

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

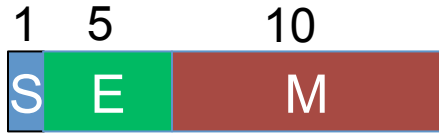Joel Emer, Vivienne Sze, Yu-Hsin Chen, Tien-Ju Yang

# Approaches

- **<u>Reduce size</u> of operands for storage/compute**

    - **Floating point -> Fixed point**

    - **Bit-width reduction**

    - **Non-linear quantization**

- **<u>Reduce number</u> of operations for storage/compute**

    - **Exploit Activation Statistics (Compression)**

    - **Network Pruning**

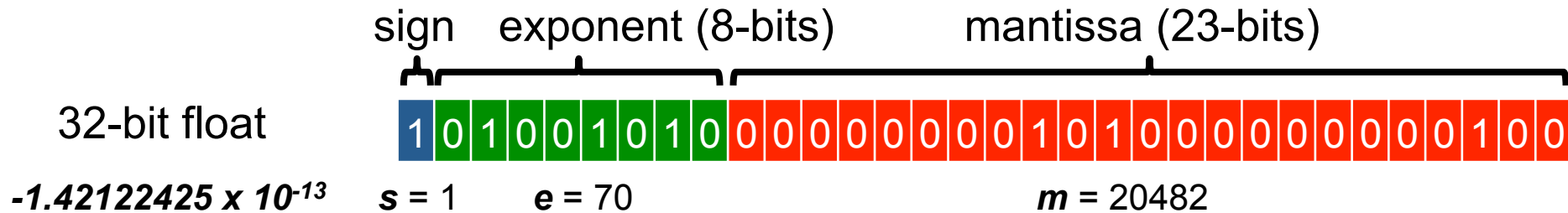    - **Compact Network Architectures**

# Cost of Operations

| Operation: | Energy (pJ) | Relative Energy Cost | Area (μm²) | Relative Area Cost |
|---|---|---|---|---|
| 8b Add | 0.03 | | 36 | |
| 16b Add | 0.05 | | 67 | |
| 32b Add | 0.1 | | 137 | |
| 16b FP Add | 0.4 | | 1360 | |
| 32b FP Add | 0.9 | | 4184 | |
| 8b Mult | 0.2 | | 282 | |
| 32b Mult | 3.1 | | 3495 | |
| 16b FP Mult | 1.1 | | 1640 | |
| 32b FP Mult | 3.7 | | 7700 | |
| 32b SRAM Read (8KB) | 5 | | N/A | |
| 32b DRAM Read | 640 | | N/A | |

Relative Energy Cost axis: $1 \quad 10 \quad 10^2 \quad 10^3 \quad 10^4$

Relative Area Cost axis: $1 \quad 10 \quad 10^2 \quad 10^3$

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

# Number Representation

| | | Range | Accuracy |
|---|---|---|---|
| FP32 | 1 / 8 / 23 → S E M | $10^{-38} - 10^{38}$ | .000006% |
| FP16 | 1 / 5 / 10 → S E M | $6 \times 10^{-5} - 6 \times 10^{4}$ | .05% |
| Int32 | 1 / 31 → S M | $0 - 2 \times 10^{9}$ | ½ |
| Int16 | 1 / 15 → S M | $0 - 6 \times 10^{4}$ | ½ |
| Int8 | 1 / 7 → S M | $0 - 127$ | ½ |

Image Source: B. Dally

# Floating Point → Fixed Point

## Floating Point

sign     exponent (8-bits)        mantissa (23-bits)

32-bit float

| 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |

$-1.42122425 \times 10^{-13}$    $s = 1$     $e = 70$        $m = 20482$

## Fixed Point

sign   mantissa (7-bits)

8-bit fixed

| 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 |

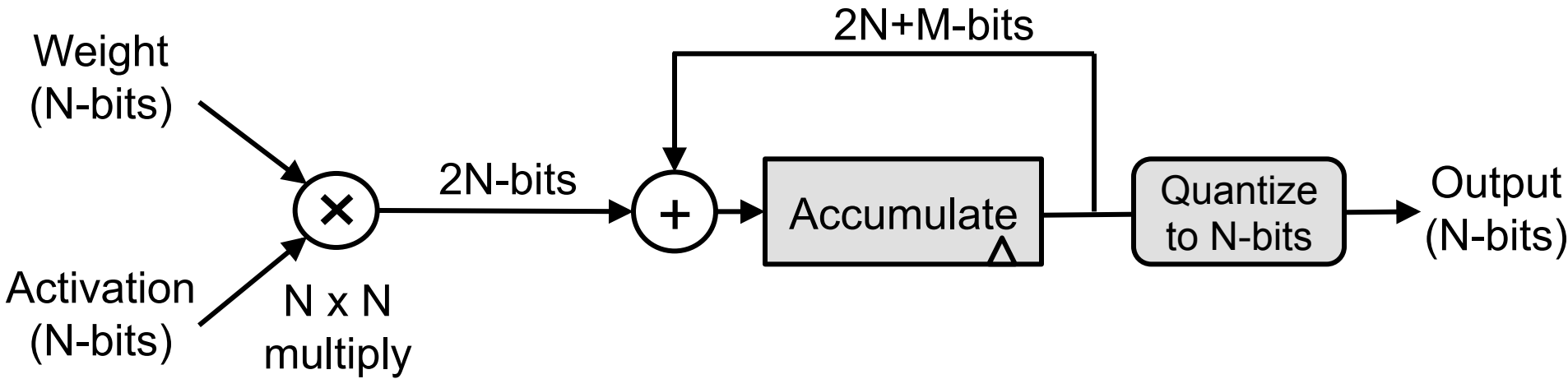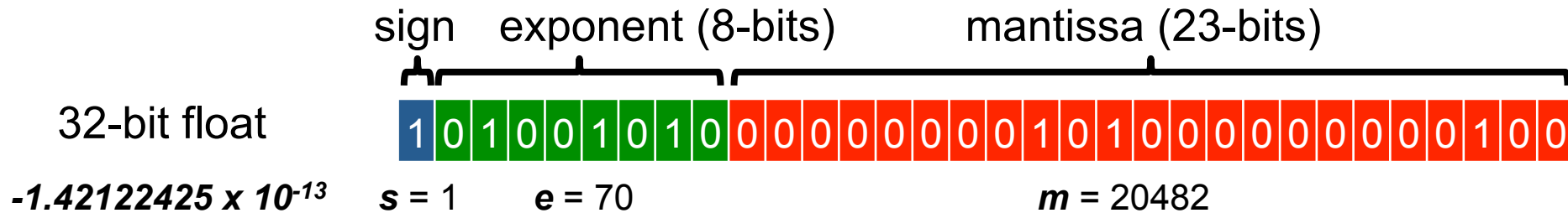integer (4-bits)     fractional (3-bits)

$12.75$      $s = 0$    $m = 102$

# N-bit Precision

For no loss in precision, **M** is determined based on largest filter size (in the range of 10 to 16 bits for popular DNNs)
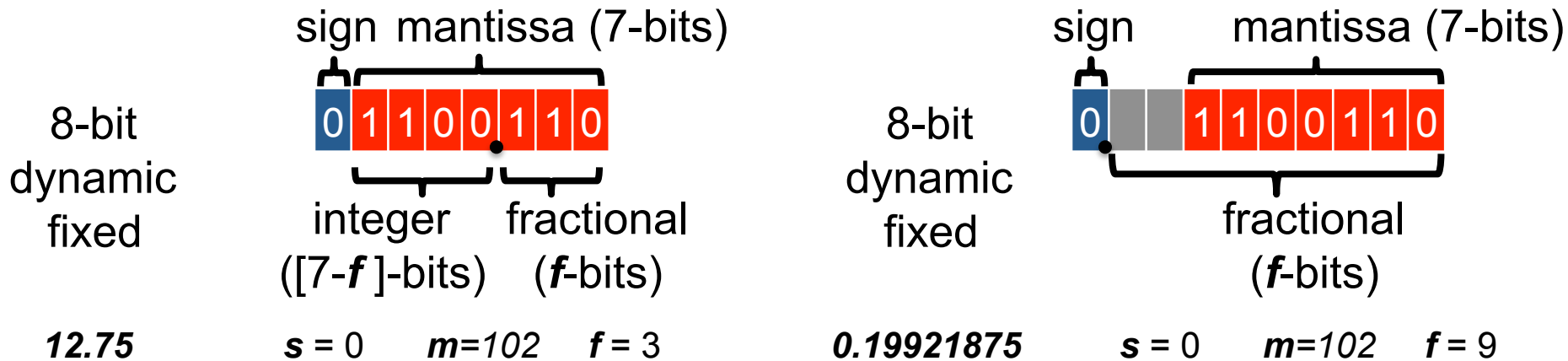
Weight (N-bits)

Activation (N-bits)
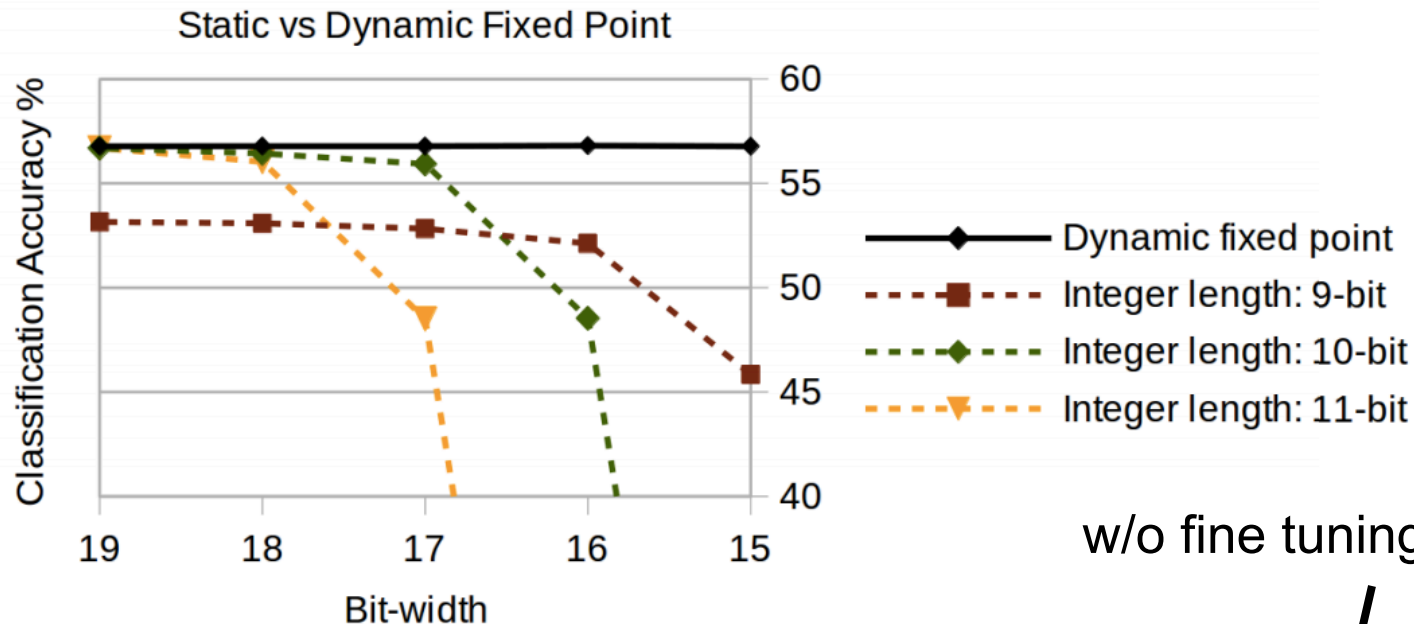
N x N multiply

2N-bits

2N+M-bits

Accumulate

Quantize to N-bits

Output (N-bits)

# Dynamic Fixed Point

## Floating Point

sign    exponent (8-bits)    mantissa (23-bits)

32-bit float

$1\ 0\ 1\ 0\ 0\ 1\ 0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0$

$-1.42122425 \times 10^{-13}$    $s = 1$    $e = 70$    $m = 20482$

## Fixed Point

sign  mantissa (7-bits)

8-bit dynamic fixed

$0\ 1\ 1\ 0\ 0\ 1\ 1\ 0$

integer ([7-$f$]-bits)    fractional ($f$-bits)

$12.75$    $s = 0$    $m=102$    $f = 3$

sign    mantissa (7-bits)

8-bit dynamic fixed

$0$ _ _ $1\ 1\ 0\ 0\ 1\ 1\ 0$

fractional ($f$-bits)

$0.19921875$    $s = 0$    $m=102$    $f = 9$

Allow $f$ to vary based on data type and layer

# Impact on Accuracy

Top-1 accuracy on of CaffeNet on ImageNet



Static vs Dynamic Fixed Point

- Dynamic fixed point
- Integer length: 9-bit
- Integer length: 10-bit
- Integer length: 11-bit

w/o fine tuning

|  | 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%) |

[Gysel et al., Ristretto, ICLR 2016]

8

# Avoiding Dynamic Fixed Point

Batch normalization 'centers' dynamic range

AlexNet
(Layer 6)

Image Source: Moons
et al, WACV 2016



'Centered' dynamic ranges might reduce need for
dynamic fixed point

# 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)

# 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)

# 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 |



[Judd et al., ArXiv 2016]    [Moons et al., WACV 2016]

# Bitwidth Scaling (Speed)

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

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



[Judd et al., Stripes, CAL 2016]

# Bitwidth Scaling (Power)

**Reduce Bit-width →
Shorter Critical Path
→ Reduce Voltage**



**33x gain
@ 1% RMSE**

Relative Power [-]

Root-Mean-Square Error [-]

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

274 mW

142

1.3X

107

1.9X

55

42      35      35      22

A      B      C      D

MEM + CTRL    MAC-array

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

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

[Moons et al., VLSI 2016]

# Binary Nets

*Binary Filters*

- ## Binary Connect (BC)

  – Weights {-1,1}, Activations 32-bit float

  – MAC → addition/subtraction

  – Accuracy loss: **19%** on AlexNet

  [Courbariaux, NIPS 2015]

- ## Binarized Neural Networks (BNN)

  – Weights {-1,1}, Activations {-1,1}

  – MAC → XNOR

  – Accuracy loss: **29.8%** on AlexNet

  [Courbariaux, arXiv 2016]

# Scale the Weights and Activations

- **Binary Weight Nets (BWN)**
  - Weights {-α, α} → except first and last layers are 32-bit float
  - Activations: 32-bit float
  - α determined by the $l_1$-norm of all weights in a layer
  - Accuracy loss: **0.8%** on AlexNet
- **XNOR-Net**
  - Weights {-α, α}
  - Activations {-$β_i$, $β_i$} → except first and last layers are 32-bit float
  - $β_i$ determined by the $l_1$-norm of all activations across channels *for given position i* of the input feature map
  - Accuracy loss: **11%** on AlexNet

Hardware needs to support both activation precisions

Scale factors (α, $β_i$) can change per layer or position in filter

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

# XNOR-Net

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

# Ternary Nets

- **Allow for weights to be zero**

  – Increase sparsity, but also increase number of bits (2-bits)

- **Ternary Weight Nets (TWN)** [Li et al., arXiv 2016]

  – Weights {-w, 0, w} → except first and last layers are 32-bit float

  – Activations: 32-bit float

  – Accuracy loss: <span style="color:red">3.7%</span> on AlexNet

- **Trained Ternary Quantization (TTQ)** [Zhu et al., ICLR 2017]

  – Weights {$-w_1$, 0, $w_2$} → except first and last layers are 32-bit float

  – Activations: 32-bit float

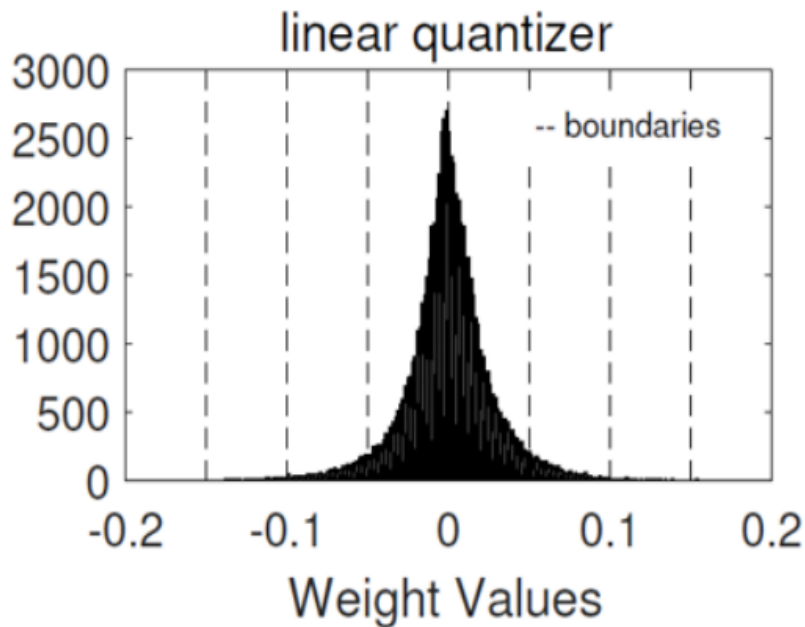  – Accuracy loss: <span style="color:red">0.6%</span> on AlexNet

# Non-Linear Quantization

- **Precision** refers to the **number of levels**
  - Number of bits = $\log_2$ (number of levels)

- **Quantization:** mapping data to a smaller set of **levels**
  - Linear, e.g., fixed-point
  - Non-linear
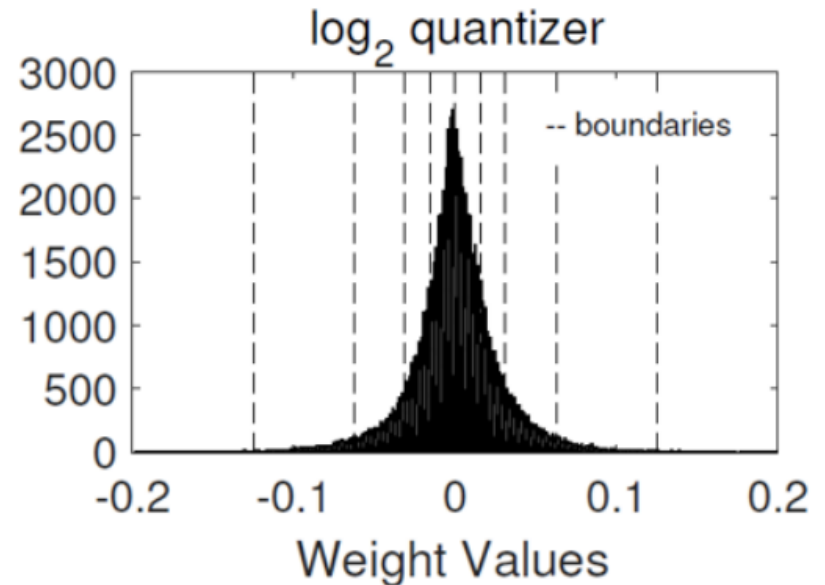    - Computed
    - Table lookup

Objective: Reduce size to improve speed and/or reduce energy while preserving accuracy

# Computed Non-linear Quantization

## Log Domain Quantization



linear quantizer

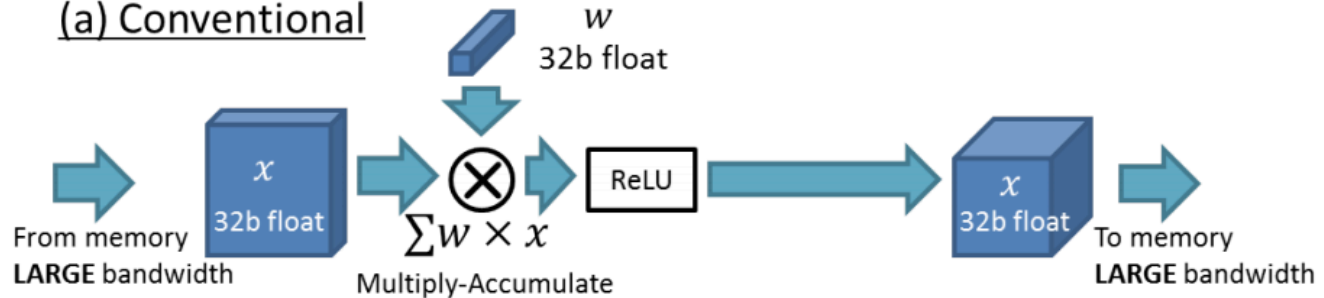log$_2$ quantizer

Product = X * W

Product = X << W

[Lee et al., LogNet, ICASSP 2017]

# Log Domain Computation



Only activation in log domain

Both weights and activations in log domain

max, bitshifts, adds/subs

[Miyashita et al., arXiv 2016]

# Log Domain Quantization

- **Weights: 5-bits for CONV, 4-bit for FC; Activations: 4-bits**

- Accuracy loss: 3.2% on AlexNet



[Miyashita et al., arXiv 2016],
[Lee et al., LogNet, ICASSP 2017]

# Reduce Precision Overview
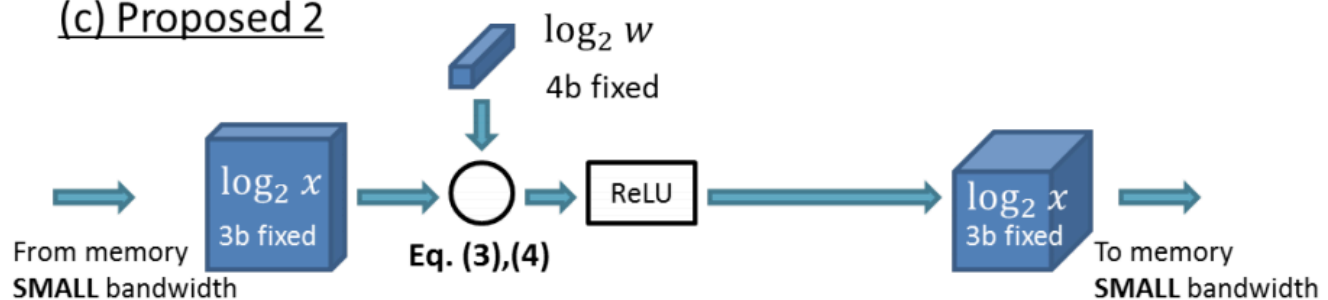
- **Learned mapping of data to quantization levels (e.g., k-means)**



*Implement with look up table*

[Han et al., ICLR 2016]

- **Additional Properties**
  - **Fixed or Variable (across data types, layers, channels, etc.)**

# Non-Linear Quantization Table Lookup

**Trained Quantization:** Find K weights via K-means clustering to reduce number of unique weights *per layer* (weight sharing)
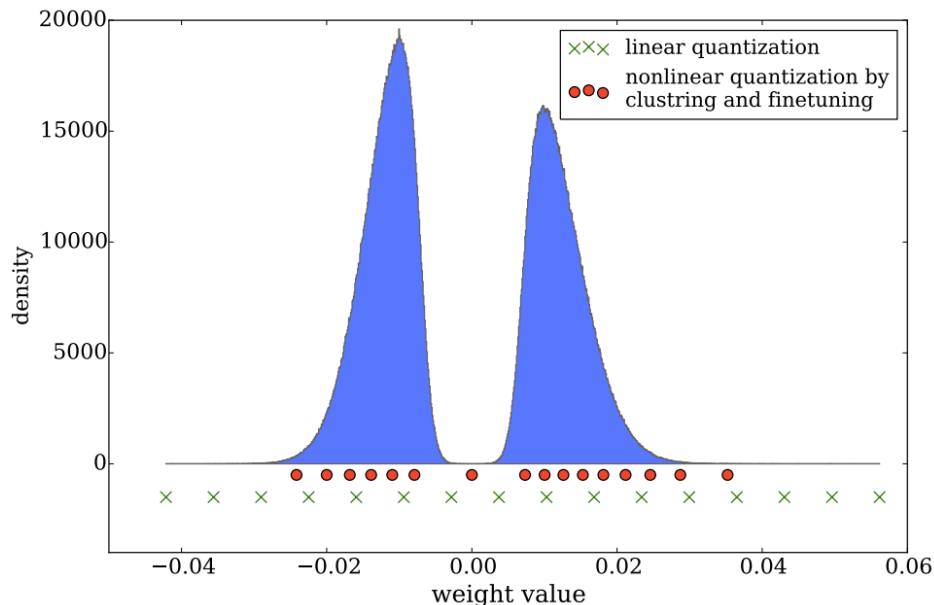
**Example:** AlexNet (no accuracy loss)
**256 unique weights** for CONV layer
**16 unique weights** for FC layer



*Smaller Weight Memory*

Weight Memory CRSM x $\log_2 U$-bits

Weight index ($\log_2 U$-bits)

*Overhead*

Weight Decoder/ Dequant U x 16b

**Weight (16-bits)**

*Does not reduce precision of MAC*

MAC

Input Activation (16-bits)

Output Activation (16-bits)

Consequences: Narrow weight memory and second access from (small) table

[Han et al., Deep Compression, ICLR 2016]

# Summary of Reduce Precision

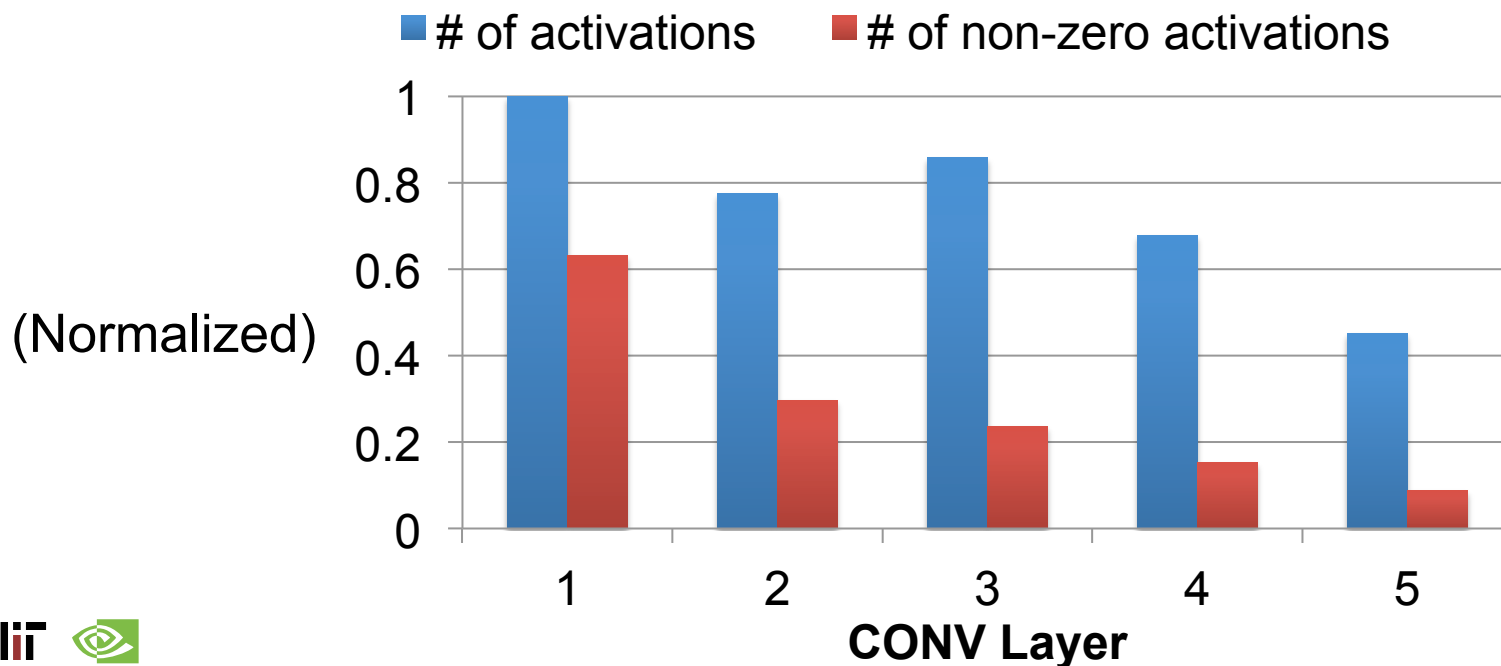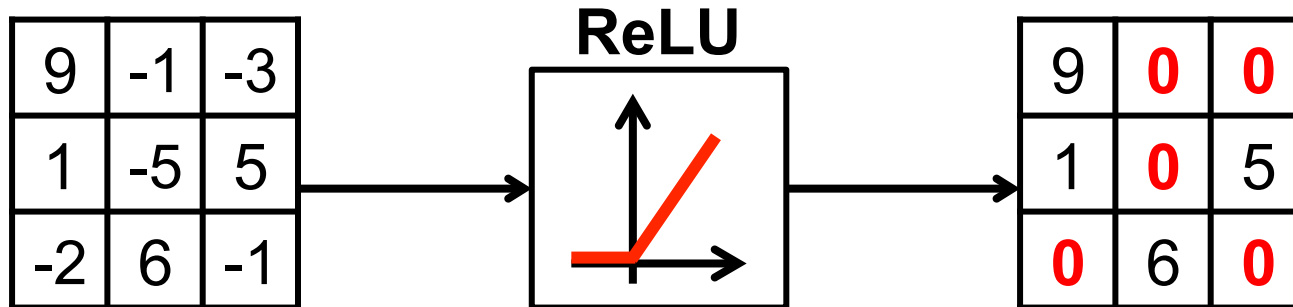| Category | Method | Weights (# of bits) | Activations (# of bits) | Accuracy Loss vs. 32-bit float (%) |
|---|---|---|---|---|
| Dynamic Fixed Point | w/o fine-tuning | 8 | 10 | 0.4 |
| | w/ fine-tuning | 8 | 8 | 0.6 |
| Reduce weight | Ternary weights Networks (TWN) | 2* | 32 | 3.7 |
| | Trained Ternary Quantization (TTQ) | 2* | 32 | 0.6 |
| | Binary Connect (BC) | 1 | 32 | 19.2 |
| | Binary Weight Net (BWN) | 1* | 32 | 0.8 |
| Reduce weight and activation | Binarized Neural Net (BNN) | 1 | 1 | 29.8 |
| | XNOR-Net | 1* | 1 | 11 |
| Non-Linear | LogNet | 5(conv), 4(fc) | 4 | 3.2 |
| | Weight Sharing | 8(conv), 4(fc) | 16 | 0 |

* first and last layers are 32-bit float

# Reduce Number of Ops and Weights

- **Exploit Activation Statistics**

- **Network Pruning**

- **Compact Network Architectures**

- **Knowledge Distillation**

# Sparsity in Fmaps

## Many **zeros** in **output fmaps** after **ReLU**

# I/O Compression in Eyeriss



[Chen et al., ISSCC 2016]

# Compression Reduces DRAM BW



DRAM Access (MB) vs AlexNet Conv Layer. Uncompressed Fmaps + Weights (blue) and RLE Compressed Fmaps + Weights (red). Compression ratios: 1.2×, 1.4×, 1.7×, 1.8×, 1.9×.

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

[Chen et al., ISSCC 2016]

# Data Gating / Zero Skipping in Eyeriss



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

[Chen et al., ISSCC 2016]

# Cnvlutin

- **Process Convolution Layers**

- **Built on top of DaDianNao (4.49% area overhead)**

- **Speed up of 1.37x (1.52x with activation pruning)**



[Albericio et al., ISCA 2016]

# Pruning Activations

## Remove small activation values

### *Speed up 11% (ImageNet)*



Cnvlutin

### *Reduce power 2x (MNIST)*



Minerva

[Albericio et al., ISCA 2016]        [Reagen et al., ISCA 2016]

# 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

retraining

[Lecun et al., NIPS 1989]

# Pruning – Make Weights Sparse

## Prune based on *magnitude* of weights



before pruning

after pruning

pruning synapses - - ->

pruning neurons - - ->

**Train Connectivity**

**Prune Connections**

**Train Weights**

***Example:*** *AlexNet*
***Weight Reduction:*** *CONV layers 2.7x, FC layers 9.9x*
*(Most reduction on fully connected layers)*
***Overall:*** *9x weight reduction, 3x MAC reduction*

[Han et al., NIPS 2015]

# 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 CSRMV
NVIDIA GeForce GTX Titan X: cuBLAS GEMV, cuSPARSE CSRMV
NVIDIA Tegra K1: cuBLAS GEMV, cuSPARSE CSRMV

Batch size = 1

[Han et al., NIPS 2015]

# Key Metrics for Embedded DNN

- **Accuracy → Measured on Dataset**

- **Speed → Number of MACs**

- **Storage Footprint → Number of Weights**

- **Energy → ?**

# 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**

[Yang et al., CVPR 2017]

# Energy-Evaluation Methodology



**CNN Shape Configuration**
**(# of channels, # of filters, etc.)**

**Hardware Energy Costs of each MAC and Memory Access**

Memory Accesses Optimization

# acc. at mem. level **1**
# acc. at mem. level **2**
⋮
# acc. at mem. level **n**

# of MACs Calculation

# of MACs

$E_{data}$

$E_{comp}$

**CNN Weights and Input Data**

[0.3, 0, -0.4, 0.7, 0, 0, 0.1, …]

Energy

L1 L2 L3 …

**CNN Energy Consumption**

# Key Observations

- Number of weights *alone* is not a good metric for energy

- **All data types** should be considered

**Energy Consumption
of GoogLeNet**



Computation
10%

Input Feature Map
25%

Weights
22%

Output Feature Map
43%

[Yang et al., CVPR 2017]

# Energy Consumption of Existing DNNs



Deeper CNNs with fewer weights do not necessarily consume less energy than shallower CNNs with more weights

[Yang et al., CVPR 2017]

# Magnitude-based Weight Pruning



Reduce number of weights by **removing small magnitude weights**

# Energy-Aware Pruning



Remove weights from layers **in order of highest to lowest energy**
**3.7x reduction in AlexNet / 1.6x reduction in GoogLeNet**

[Yang et al., CVPR 2017]

# Compression of Weights & Activations

- **Compress weights and activations 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 l1 | 7 (21%) / 4 (29%) | 1.17x / 1.3x | 1 / 0.77 | 0.85 | 105 | 85 | 0.96 |
| AlexNet l2 | 7 (19%) / 7 (89%) | 1.15x / **5.8x** | 3.2 / 1.1 | 0.9 | 224 | 55 | 1.4 |
| AlexNet l3 | 8 (11%) / 9 (82%) | 1.05x / 4.1x | 6.5 / 2.8 | 0.92 | 150 | 77 | 0.7 |
| AlexNet l4 | 9 (04%) / 8 (72%) | 1.00x / 2.9x | 5.4 / 3.2 | 0.92 | 112 | 95 | 0.56 |
| AlexNet l5 | 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 l1 | 3 (35%) / 1 (87%) | 1.40x / 5.2x | 0.002 / 0.001 | 0.7 | 0.3 | 25 | 1.07 |
| LeNet-5 l2 | 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** |

[Moons et al., VLSI 2016; Han et al., ICLR 2016]

# 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** >> **N**)
*For DNN, **M** = # of filters, **N** = # of weights per filter*

[Dorrance et al., FPGA 2014]

# 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**



**Input**

**Output**

**Weights**

Dequantize Weight

Keep track of location

[Han et al., ISCA 2016]

# Compact Network Architectures

- **Break large convolutional layers into a series of smaller convolutional layers**
  - **Fewer weights, but same effective receptive field**

- **Before Training: Network Architecture Design**

- **After Training: Decompose Trained Filters**

# Network Architecture Design

## Build Network with series of Small Filters

**GoogleNet/Inception v3**

5x5 filter

decompose →

5x1 filter

1x5 filter

*separable filters*

Apply sequentially

**VGG-16**

5x5 filter

decompose →

Two 3x3 filters

Apply sequentially

# Network Architecture Design

Reduce size and computation with 1x1 Filter (**bottleneck**)



1x1 CONV
with 32 filters

(each filter has size
1x1x64, and performs a
64-dimensional dot
product)

56

64

56

Figure Source:
Stanford cs231n

56

56

32

Used in Network In Network(NiN) and GoogLeNet

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

# Network Architecture Design

Reduce size and computation with 1x1 Filter (**bottleneck**)



1x1 CONV
with 32 filters

(each filter has size
1x1x64, and performs a
64-dimensional dot
product)

56

64

56

56

56

32

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 Design

Reduce size and computation with 1x1 Filter (**bottleneck**)



1x1 CONV
with 32 filters

(each filter has size
1x1x64, and performs a
64-dimensional dot
product)

Figure Source:
Stanford cs231n

56
56
64

56
56
32
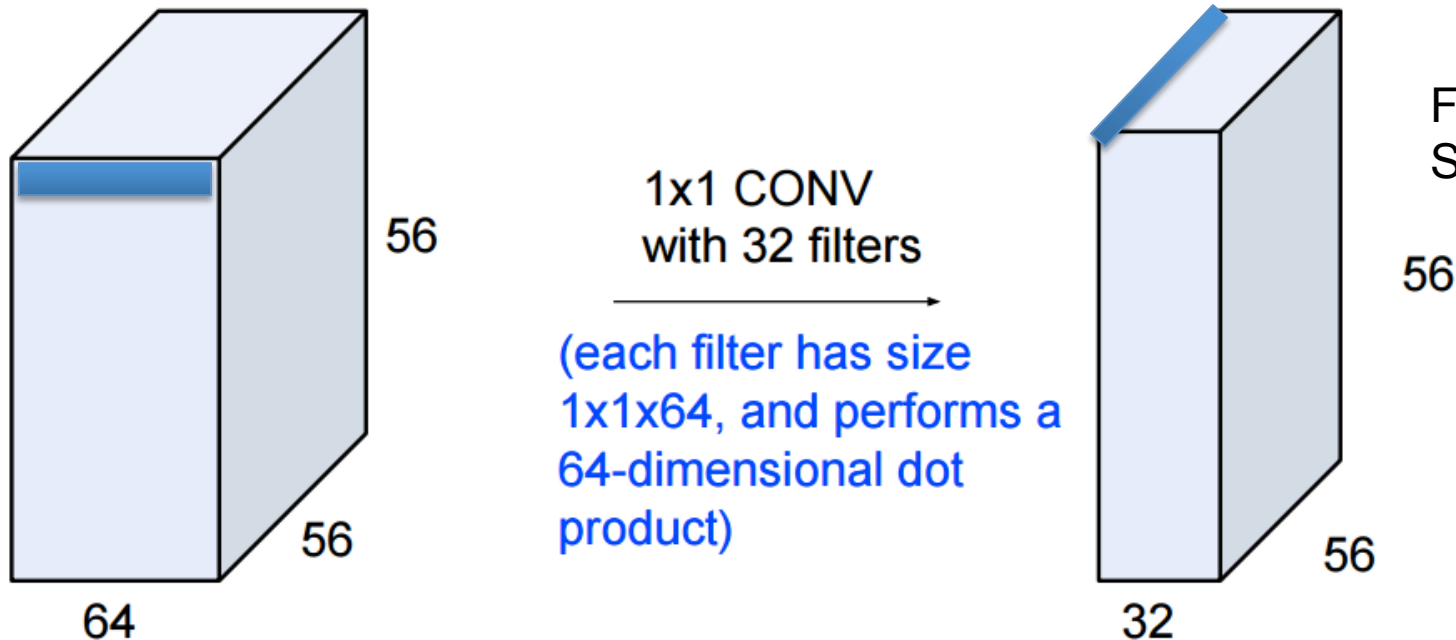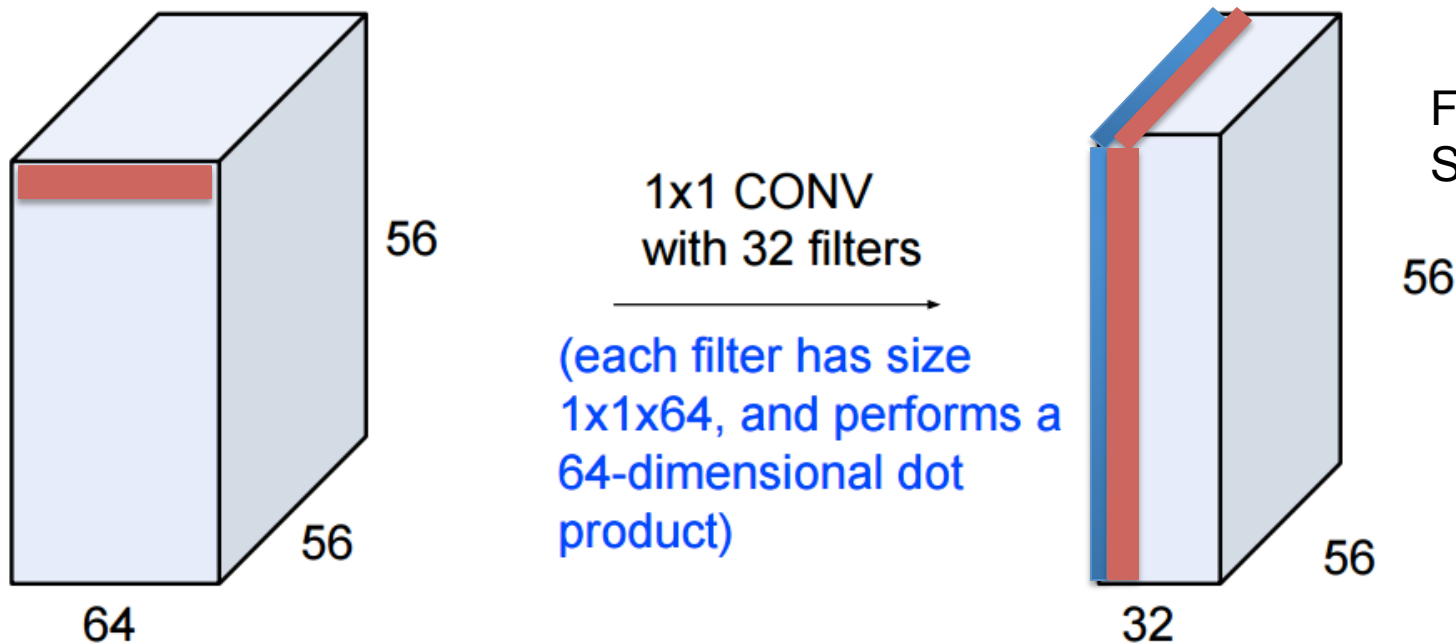
Used in Network In Network(NiN) and GoogLeNet

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

# Bottleneck in Popular DNN models

# SqueezeNet

Reduce weights by reducing number of input channels by "squeezing" with 1x1
**50x fewer weights than AlexNet (no accuracy loss)**



**Fire Module**

[F.N. Iandola et al., ArXiv, 2016]

# Energy Consumption of Existing DNNs



Deeper CNNs with fewer weights do not necessarily consume less energy than shallower CNNs with more weights

[Yang et al., CVPR 2017]

# Decompose Trained Filters

**After training,** perform **low-rank approximation by applying tensor decomposition** to weight kernel; then **fine-tune** weights for accuracy



(a) Full convolution

(b) Two-component decomposition (Jaderberg et al., 2014a)

(c) CP-decomposition

**R** = canonical rank

[Lebedev et al., ICLR 2015]

# Decompose Trained Filters

## Visualization of Filters

Original  Approx.



- **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

[Denton et al., NIPS 2014]

# Decompose Trained Filters on Phone

## Tucker Decomposition



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

[Kim et al., ICLR 2016]

# Knowledge Distillation



[Bucilu et al., KDD 2006],[Hinton et al., arXiv 2015]

# Metrics to Compare DNN Models

- **How can we compare different models?**

- **Accuracy**

- **Network Architecture**
  - **# Layers, filter size, # of filters, # of channels**

- **# of Weights (storage capacity)**
  - **Number of non-zero (NZ) weights**

- **# of MACs (operations)**
  - **Number of non-zero (NZ) MACS**

# Metrics of DNN Models

| Metrics | AlexNet | VGG-16 | GoogLeNet (v1) | ResNet-50 |
|---|---|---|---|---|
| Accuracy (top-5 error)* | 19.8 | 8.80 | 10.7 | 7.02 |
| Input | 227x227 | 224x224 | 224x224 | 224x224 |
| **# of CONV Layers** | **5** | **16** | **21** | **49** |
| Filter Sizes | 3, 5, 11 | 3 | 1, 3, 5, 7 | 1, 3, 7 |
| # of Channels | 3 - 256 | 3 - 512 | 3 - 1024 | 3 - 2048 |
| # of Filters | 96 - 384 | 64 - 512 | 64 - 384 | 64 - 2048 |
| Stride | 1, 4 | 1 | 1, 2 | 1, 2 |
| # of Weights | 2.3M | 14.7M | 6.0M | 23.5M |
| # of MACs | 666M | 15.3G | 1.43G | 3.86G |
| **# of FC layers** | **3** | **3** | **1** | **1** |
| # of Weights | 58.6M | 124M | 1M | 2M |
| # of MACs | 58.6M | 124M | 1M | 2M |
| **Total Weights** | **61M** | **138M** | **7M** | **25.5M** |
| **Total MACs** | **724M** | **15.5G** | **1.43G** | **3.9G** |

*Single crop results: https://github.com/jcjohnson/cnn-benchmarks

# Metrics of DNN Models

| Metrics | AlexNet | VGG-16 | GoogLeNet (v1) | ResNet-50 |
|---|---|---|---|---|
| Accuracy (top-5 error)* | 19.8 | 8.80 | 10.7 | 7.02 |
| # of CONV Layers | 5 | 16 | 21 | 49 |
| # of Weights | 2.3M | 14.7M | 6.0M | 23.5M |
| # of MACs | 666M | 15.3G | 1.43G | 3.86G |
| # of NZ MACs** | 394M | 7.3G | 806M | 1.5G |
| # of FC layers | 3 | 3 | 1 | 1 |
| # of Weights | 58.6M | 124M | 1M | 2M |
| # of MACs | 58.6M | 124M | 1M | 2M |
| # of NZ MACs** | 14.4M | 17.7M | 639k | 1.8M |
| Total Weights | 61M | 138M | 7M | 25.5M |
| Total MACs | 724M | 15.5G | 1.43G | 3.9G |
| # of NZ MACs** | 409M | 7.3G | 806M | 1.5G |

*Single crop results: https://github.com/jcjohnson/cnn-benchmarks

**# of NZ MACs based on 50k ImageNet validation images

# Metrics of DNN Algorithms

| Metrics | AlexNet | AlexNet (sparse) |
|---|---|---|
| Accuracy (top-5 error) | 19.8 | 19.8 |
| **# of Conv Layers** | **5** | **5** |
| # of Weights | 2.3M | 2.3M |
| # of MACs | 666M | 666M |
| # of NZ weights | **2.3M** | **863k** |
| # of NZ MACs | **394M** | **207M** |
| **# of FC layers** | **3** | **3** |
| # of Weights | 58.6M | 58.6M |
| # of MACs | 58.6M | 58.6M |
| # of NZ weights | **58.6M** | **5.9M** |
| # of NZ MACs | **14.4M** | **2.1M** |
| **Total Weights** | **61M** | **61M** |
| **Total MACs** | **724M** | **724M** |
| **# of NZ weights** | **61M** | **6.8M** |
| **# of NZ MACs** | **409M** | **209M** |

# of NZ MACs based on 50k ImageNet validation images

# Tutorial Summary

- **DNNs are a critical component in the AI revolution**, delivering record breaking accuracy on many important AI tasks for a wide range of applications; however, it comes at the cost of **high computational complexity**

- **Efficient processing of DNNs** is an important area of research with many promising opportunities for innovation at **various levels of hardware design, including algorithm co-design**

- When considering different DNN solutions it is important to **evaluate with the appropriate workload** in term of both input and model, and recognize that they are **evolving rapidly**.

- It's important to consider a **comprehensive set of metrics** when evaluating different DNN solutions: **accuracy, speed, energy, and cost**

# Resources

- **Eyeriss Project: http://eyeriss.mit.edu**
  - Tutorial Slides
  - Benchmarking
  - Energy modeling
  - Mailing List for updates    Follow @eems_mit
    - **http://mailman.mit.edu/mailman/listinfo/eems-news**
  - **Paper based on today's tutorial:**
    - V. Sze, Y.-H. Chen, T-J. Yang, J. Emer, "*Efficient Processing of Deep Neural Networks: A Tutorial and Survey*", arXiv, 2017