Lecture Series on
Hardware for Deep Learning

# Part 4:
# Reducing the Complexity

Dr. Adam Teman

EnICS Labs, Bar-Ilan University

28 April 2020

EnICS
Emerging Nanoscaled
Integrated Circuits and Systems Labs

Tradition of Excellence

Bar-Ilan University
אוניברסיטת בר-אילן

# Outline

Motivation

# Models are Getting Larger



Dally, NIPS'2016 workshop on Efficient Methods for Deep Neural Networks © Adam Teman, 2020

# Explosion in size, complexity, energy

## Deep neural networks are energy hungry and growing fast

AI is being powered by the explosive growth of deep neural networks

Number of parameters in neural net (~energy consumption)

1940  1950  1960  1970  1980  1990  2000  2010  2020  2030

*Source: Qualcomm*

| Network | Model size (MB) | GFLOPS |
|---|---|---|
| AlexNet* | 233 | 0.7 |
| VGG-16* | 528 | 15.5 |
| VGG-19* | 548 | 19.6 |
| ResNet-50* | 98 | 3.9 |
| ResNet-101* | 170 | 7.6 |
| ResNet-152* | 230 | 11.3 |
| GoogleNet# | 27 | 1.6 |
| InceptionV3# | 89 | 6 |
| MobileNet# | 38 | 0.58 |
| SequeezeNet# | 30 | 0.84 |

*: Characterization and Benchmarking of Deep Learning, Natalia Vassilieva
#: https://github.com/albanie/convnet-burden

# Big Three Challenges

- **First Challenge: Model Size**
  - Hard to distribute large models through over-the-air update

- **Second Challenge: Speed**
  - Such long training time limits ML researcher's productivity

- **Third Challenge: Energy Efficiency**
  - AlphaGo: 1920 CPUs and 280 GPUs, $3000 electric bill per game
  - On mobile: drains battery
  - On data-center: increases TCO



This item is over 100MB.
Microsoft Excel will not download until you connect to Wi-Fi.

Cancel    OK

App icon is in the public domain
Phone image is licensed under CC-BY 2.0

This image is licensed under CC-BY 2.0

|            | Error rate | Training time |
|------------|------------|---------------|
| ResNet18:  | 10.76%     | 2.5 days      |
| ResNet50:  | 7.02%      | 5 days        |
| ResNet101: | 6.21%      | 1 week        |
| ResNet152: | 6.16%      | 1.5 weeks     |

*Source: Han*

# Where is the Energy Consumed?

- **Larger model**

- **More memory references**

- **More energy**

- **How can we make our models more energy efficient?**

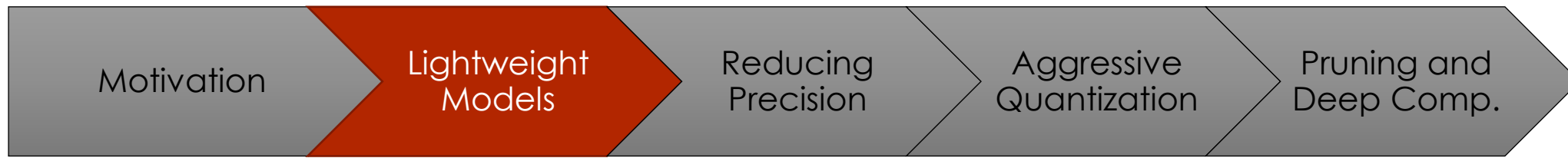| Operation | Energy [pJ] |
|---|---|
| 32 bit int ADD | 0.1 |
| 32 bit float ADD | 0.9 |
| 32 bit Register File | 1 |
| 32 bit int MULT | 3.1 |
| 32 bit float MULT | 3.7 |
| 32 bit SRAM Cache | 5 |
| **32 bit DRAM Memory** | **640** |

Relative Energy Cost

$1 \times = 1000 \times +$

This image is in the public domain

*Source: Han*

# Lightweight Models

# Reminder: Standard Convolution

- **Layer sizes:**
  - Input fmap: $HxWxC$
  - Filter size: $RxSxC$
  - Output size: $ExFxM$
- **A bit simplified:**
  - Assume: $H=W=E=F$
  - Assume: $R=S=k$
- **Cost of convolution:**
  - $M$ output maps of size $H^2$.
  - Each one requires $k^2*C$ MACs
  - Total MACs: $M*H^2*k^2*C$
  - Total Weights: $M*k^2*C$

Proportional to number of Input Channels: $C$

Proportional to size of convolutional kernel: $k^2$

Proportional to spatial size of output map: $H^2$

Proportional to number of Output Channels: $M$

Standard | Group | Pointwise | Depthwise | Factorized

# Spatial and Channel Connectivity

- **To visualize the connectivity complexity, we can use a pair of illustrations**
  - For a 3x3 kernel, looking at one spatial dimension (e.g., one row), the connectivity between the input activation and output fmap looks as follows:

  - And across channels, each input channel is connected to each output channel, so we get:

- **So we see that for convolutions:**
  - **Spatially**, the inputs and outputs are connected *locally*.
  - Across **channels**, the inputs and outputs are *fully connected*.



*Source: Yusuke Uchida*

# Group Convolutions

- **Observation:**
  - The more filters in a layer ($M$), the more *intermediate features* we learn.
- **Problem:**
  - This leads to a lot of operations (Total MACs: $M*H^2*k^2*C$)
- **Grouped Convolutions:**
  - Reduce the number of operations by dividing the input into several groups.
  - Essentially, we can learn different features through different routes.
  - First used by **AlexNet** to split a network onto two GPUs.



*Source: Krizhevsky 2012*

Standard  Group  Pointwise  Depthwise  Factorized

# Group Convolutions

- **So now we have:**
  - $G$ groups of $M/G$ filters
  - $G$ output fmaps of $M/G$ depth
  - Total MACs: $G*(M/G*H^2*k^2*C/G)$
  - That's a reduction of $1/G$.
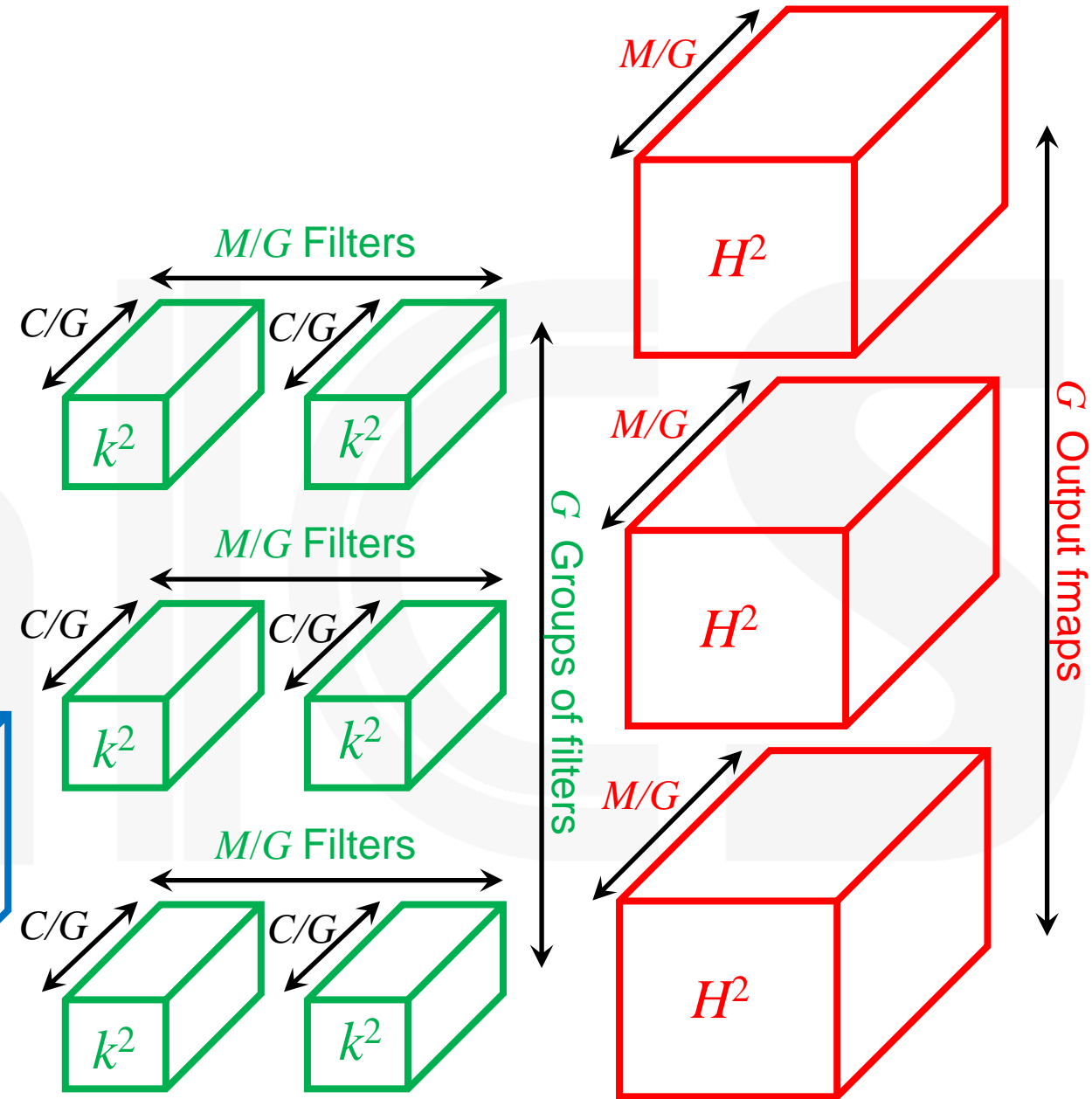- **Visualization: Gconv 3x3**



*Source: Yusuke Uchida*

Standard  Group  Pointwise  Depthwise  Factorized

# Pointwise (1x1) Convolution

- **Problem:**
  - Convolving a large filter over many input channels is expensive ($k^2 * C$)
- **Solution:**
  - Merge channels with a $1x1xC$ filter
  - Use $M$ filters to get the desired input channel depth
  - Total cost: $M * H^2 * C$.
- **This "blends" information across channels:**



*Source: Chi-Feng Wang*

Spatial          Channel

Input

Output

*Source: Yusuke Uchida*

Standard   Group   Pointwise   Depthwise   Factorized

# Example: Inception (GoogLeNet)

- **GoogLeNet was intended to solve three problems**
  - Previous models kept going deeper
    → computationally expensive
  - Variation in location of information
    →Need several filter sizes for each feature
  - Deep networks are prone to overfitting

- **Solution: Go *Wider***
  - Use an **"Inception Layer"** to split activations into several routes with different filter sizes

- **But this is computationally expensive**
  - So reduce dimensionality with 1x1 convolution and then stack a larger filter on top



*Source: Google, Inception v*

Standard  Group  Pointwise  Depthwise  Factorized

# Example: SqueezeNet

- **The "Fire Module" of SqueezeNet:**
  - Uses 1x1 convolutions to reduce channel depth
  - Uses 1x1 and 3x3 convolutions to expand it back

**Two other interesting concepts in SqueezeNet:**

- **Downsampling**
  - Use pooling with a stride of ½ late in the network.
  - This provides late convolution layers with many parameters
- **No fully connected layers**
  - Finish with $N$ channels for $N$ classification categories
  - Use average pooling on a channel for a classification score



*Source: Song Han*

Standard | Group | Pointwise | Depthwise | Factorized

# Depthwise Convolutions

- **A popular way of doing low cost convolutions is to combine *Group Convolutions* with *Pointwise Convolutions*.**

- **Let's start by looking at a standard convolution:**

  - Starting with an input of $H$x$W$x$C$ we want to arrive at an output of $E$x$F$x$M$.
  - The standard approach is to use $M$ filters with a depth of $C$.
  - For example, a 7x7x3 input to a 5x5x128 output needs 128 3x3x3 filters.
  - Total MACs: 86,400
  - Weights: 3,456



Source: Kunlun Bai

x 128

7 7 3

3 3 3

5 5 128

Standard  Group  Pointwise  **Depthwise**  Factorized

# Depthwise Convolutions

- **Instead let's make a <span style="color:purple">group convolution</span> with $C$ groups:**
  - $C$ filters of $k$x$k$x$1$.
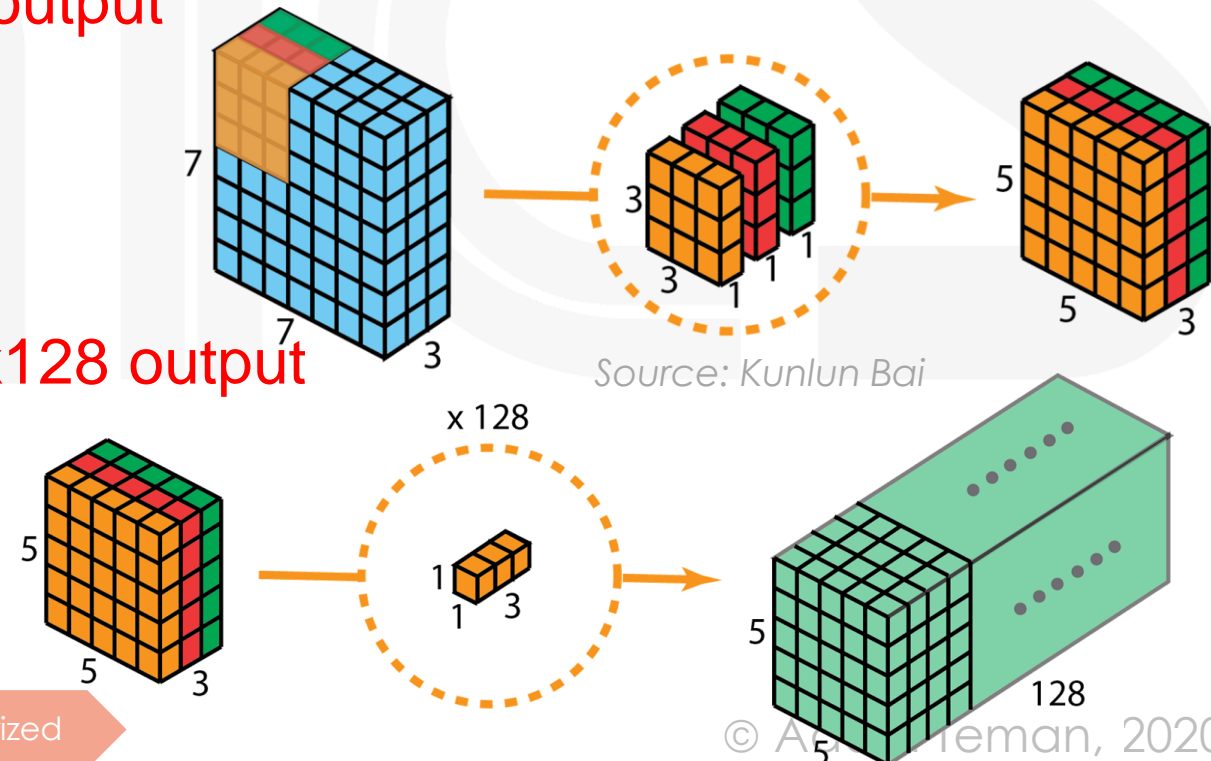  - Each filter is applied to <span style="color:blue">one input channel</span>, providing <span style="color:red">one output fmap</span>.
  - Concatenating these we get an <span style="color:red">output of $E$x$F$x$C$</span>.
  - In our example, <span style="color:green">3 3x3x1 filters</span>, <span style="color:red">5x5x3 output</span>
- **Now use a pointwise (1x1) convolution:**
  - $M$ filters of $1$x$1$x$C$.
  - Provides the desired <span style="color:red">output of $E$x$F$x$M$</span>.
  - In our example, <span style="color:green">128 1x1x3 filters</span>, <span style="color:red">5x5x128 output</span>
- **How much did it cost?**
  - Total MACs: <span style="color:purple">16,675</span> (-80%)
  - Total weights: <span style="color:purple">411</span> (-90%)

*Source: Kunlun Bai*

x 128

Standard → Group → Pointwise → Depthwise → Factorized

# Example: MobileNet

- **Introduced by Google in 2017**
  - Applies Batch Normalization and ReLU after each Depthwise Convolution
  - Better accuracy than **VGG-16** with 97% fewer weights and 97% fewer MACs

Table 8. MobileNet Comparison to Popular Models

| Model | ImageNet Accuracy | Million Mult-Adds | Million Parameters |
|---|---|---|---|
| 1.0 MobileNet-224 | 70.6% | 569 | 4.2 |
| GoogleNet | 69.8% | 1550 | 6.8 |
| VGG 16 | 71.5% | 15300 | 138 |

http://blog.csdn.net/u011995719

3x3 Depthwise Conv → BN → ReLU → 1x1 Conv → BN → ReLU

Table 1. MobileNet Body Architecture

| Type / Stride | Filter Shape | Input Size |
|---|---|---|
| Conv / s2 | $3 \times 3 \times 3 \times 32$ | $224 \times 224 \times 3$ |
| Conv dw / s1 | $3 \times 3 \times 32$ dw | $112 \times 112 \times 32$ |
| Conv / s1 | $1 \times 1 \times 32 \times 64$ | $112 \times 112 \times 32$ |
| Conv dw / s2 | $3 \times 3 \times 64$ dw | $112 \times 112 \times 64$ |
| Conv / s1 | $1 \times 1 \times 64 \times 128$ | $56 \times 56 \times 64$ |
| Conv dw / s1 | $3 \times 3 \times 128$ dw | $56 \times 56 \times 128$ |
| Conv / s1 | $1 \times 1 \times 128 \times 128$ | $56 \times 56 \times 128$ |
| Conv dw / s2 | $3 \times 3 \times 128$ dw | $56 \times 56 \times 128$ |
| Conv / s1 | $1 \times 1 \times 128 \times 256$ | $28 \times 28 \times 128$ |
| Conv dw / s1 | $3 \times 3 \times 256$ dw | $28 \times 28 \times 256$ |
| Conv / s1 | $1 \times 1 \times 256 \times 256$ | $28 \times 28 \times 256$ |
| Conv dw / s2 | $3 \times 3 \times 256$ dw | $28 \times 28 \times 256$ |
| Conv / s1 | $1 \times 1 \times 256 \times 512$ | $14 \times 14 \times 256$ |
| $5\times$   Conv dw / s1 | $3 \times 3 \times 512$ dw | $14 \times 14 \times 512$ |
|   Conv / s1 | $1 \times 1 \times 512 \times 512$ | $14 \times 14 \times 512$ |
| Conv dw / s2 | $3 \times 3 \times 512$ dw | $14 \times 14 \times 512$ |
| Conv / s1 | $1 \times 1 \times 512 \times 1024$ | $7 \times 7 \times 512$ |
| Conv dw / s2 | $3 \times 3 \times 1024$ dw | $7 \times 7 \times 1024$ |
| Conv / s1 | $1 \times 1 \times 1024 \times 1024$ | $7 \times 7 \times 1024$ |
| Avg Pool / s1 | Pool $7 \times 7$ | $7 \times 7 \times 1024$ |
| FC / s1 | $1024 \times 1000$ | $1 \times 1 \times 1024$ |
| Softmax / s1 | Classifier | $1 \times 1 \times 1000$ |

Standard | Group | Pointwise | Depthwise | Factorized
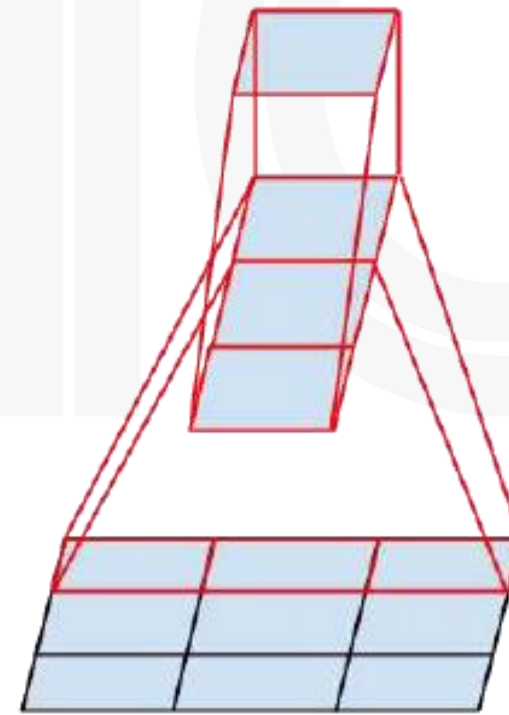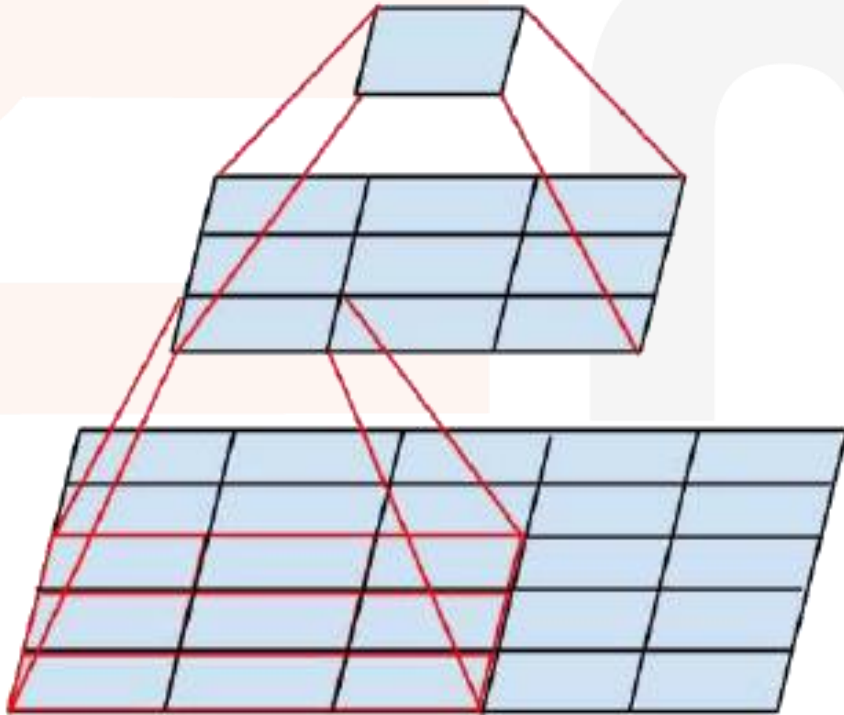
*Source: Google, MobileNets*  © Adam Teman, 2020

# Example: ShuffleNet

- **Apply a "Channel Shuffle"**
  - 1x1 Group Convolution and shuffle the outputs
- **Also use Depthwise Convolutions and Residuals**
- **Outperforms MobileNet**



*Source: Zhang, et al, ShuffleNet*

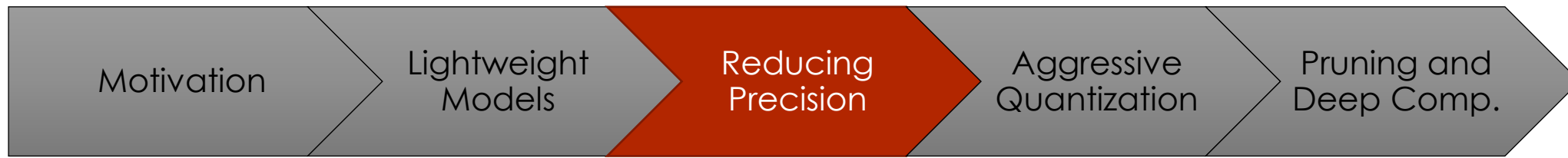Standard | Group | Pointwise | **Depthwise** | Factorized

# Factorized (Stacked) Convolutions

- **Reduce the number of weights using two smaller filters:**
  - **VGG**: two 3x3 filters (18 weights) replace one 5x5 filter (25 weights)
  - **Inception v2**: $1xn$ and $nx1$ filters ($2n$ weights) replace $nxn$ filter ($n^2$ weights)
    For example: 3x1 and 1x3 filters (6 weights) replace 3x3 filter (9 weights)
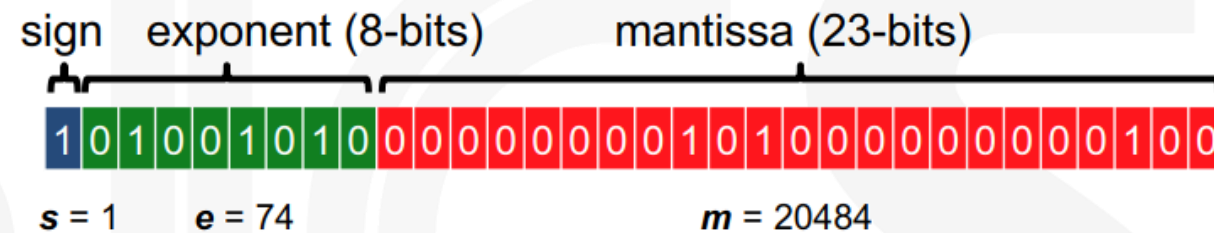
*Source: Sik-Ho Tsang*

Standard ⟩ Group ⟩ Pointwise ⟩ Depthwise ⟩ **Factorized**

# Reducing Precision

Emerging Nanoscaled
Integrated Circuits and Systems Labs

Bar-Ilan University
אוניברסיטת בר-אילן

# Taxonomy

- **Precision refers to the number of levels**
  - *Number of bits = $\log_2$ (number of levels)*
  - Normal Precision: FP32
  - Low Precision: FP16, INT8
- **Mixed Precision**
  - Utilizing several precisions (e.g., FP32 and FP16) in model.
- **Quantization:** mapping data to a smaller set of levels
  - Linear, e.g., fixed-point (e.g., INT8, binary)
  - Non-linear
    - Computed (e.g., floating point, log-domain)
    - Table lookup (e.g., learned)

## Floating Point (FP32):

$$-1.112934 \times 10^{-16}$$



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

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

$s = 1$    $e = 74$    $m = 20484$

## Fixed Point (INT8):

$$12.75$$



sign | mantissa (7-bits)

0 1 1 0 0 1 1 0

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

$s = 0$    $m = 102$

# Number Representation

## Dynamic Fixed Point



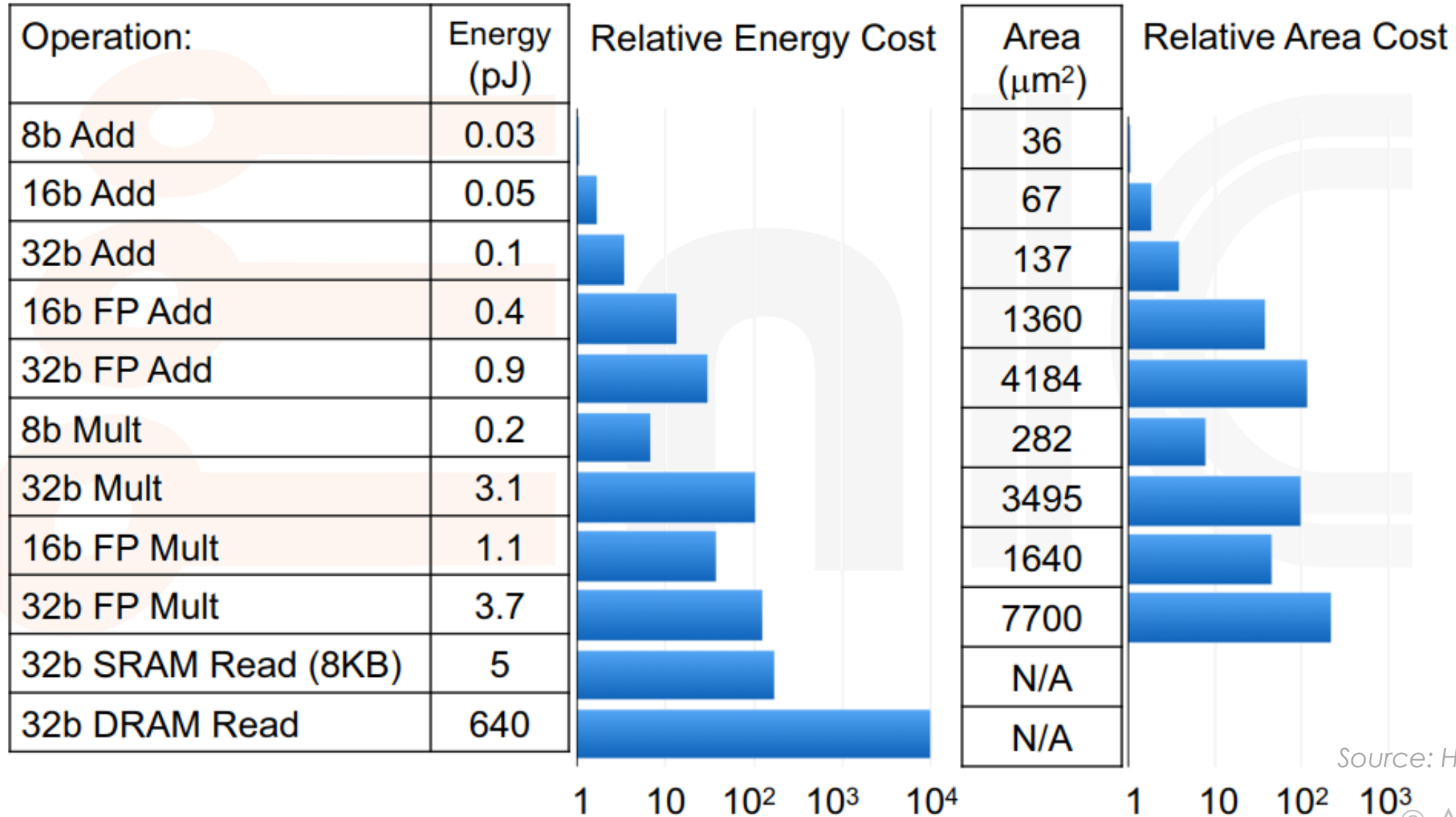| | Range | Accuracy |
|---|---|---|
| **FP32** S E M (1, 8, 23) | $10^{-38} - 10^{38}$ | .000006% |
| **FP16** S E M (1, 5, 10) | $6 \times 10^{-5} - 6 \times 10^{4}$ | .05% |
| **Int32** S M (1, 31) | $0 - 2 \times 10^{9}$ | ½ |
| **Int16** S M (1, 15) | $0 - 6 \times 10^{4}$ | ½ |
| **Int8** S M (1,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

*Source: B. Dally*

- Several scaling factors
- Different range for activations, bias, weights and gradients.

*Sources: Courbariaux, Montreal*
*UC Davis, Ristretto Project*

- Same dynamic range as FP32
- Easier for training and debugging than FP16
- Supported by Google TPU, Intel Xeon and Nirvana, others

## bfloat16: Brain Floating Point Format

Exponent: 8 bits    Mantissa (Significand): 7 bits

S E E E E E E E E M M M M M M M

**Range: ~$1e^{-38}$ to ~$3e^{38}$**

*Source: Patterson, GoogleAI*

# 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 | |

*Source: Horowitz, ISSCC 2014*

© Adam Teman, 2020

# Mixed Precision

- **Mixed Precision refers to using both full and reduced precision in a model:**
  - Identify the steps that require FP32, and use lower precision (e.g., FP16) everywhere else.
  - Has been shown to provide 2-4X speedup.
- **Low precision is supported by hardware and software platforms**
  - Google TPUs support a mix of FP32 and bfloat16
  - Nvidia Tensor Cores accelerate FP16 matrix multiplications and convolutions
  - Keras provides a mixed precision API in TensorFlow

| Model | Speedup |
|---|---|
| BERT Q&A | 3.3X speedup |
| GNMT | 1.7X speedup |
| NCF | 2.6X speedup |
| ResNet-50-v1.5 | 3.3X speedup |
| SSD-RN50-FPN-640 | 2.5X speedup |



$$D = \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix}$$

FP16 or FP32     FP16     FP16     FP16 or FP32

*Source: NVIDIA*

# Quantization

- **Quantization is mapping to a smaller set of levels**
  - e.g., floating point (FP32) to integer (INT8)
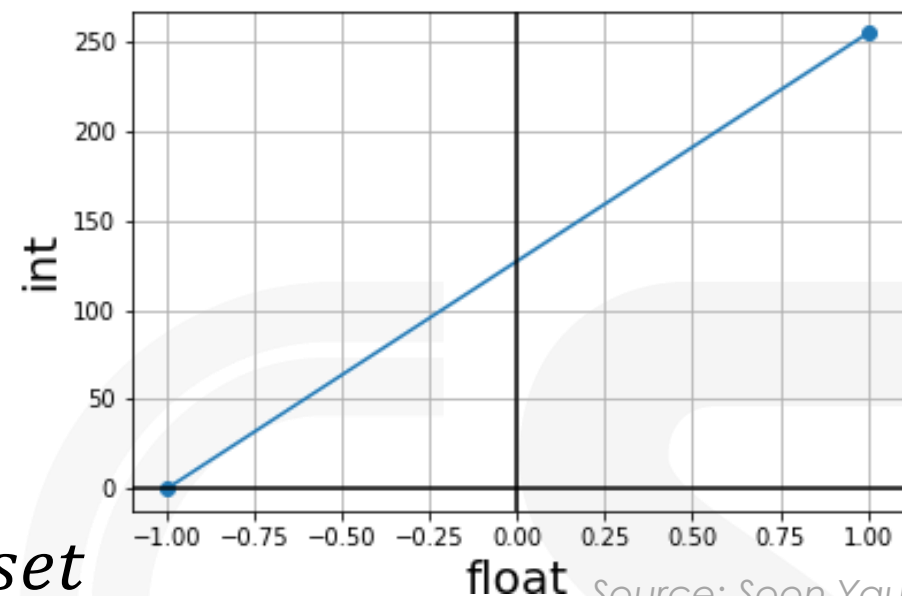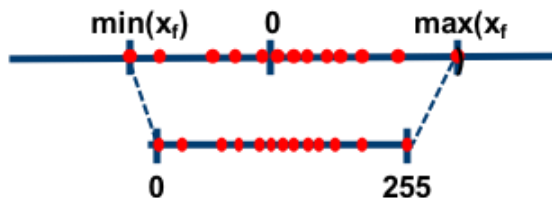- **How is it done?**
  - Well, there are a lot of tips and tricks, but basically we just need to scale and offset:
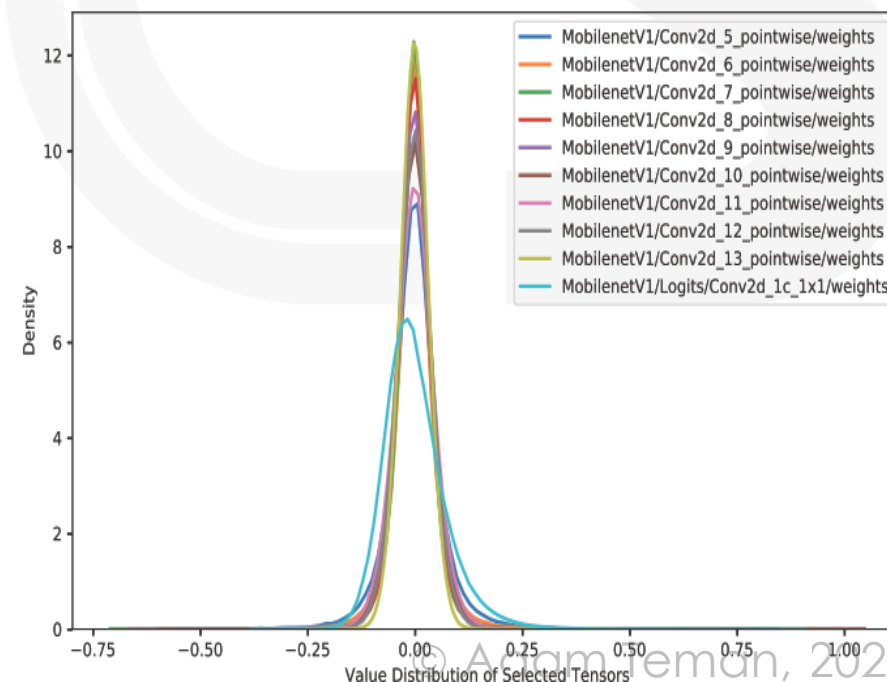
$$x_q = \frac{x_f}{scale} + offset$$

  - The scaling factor is dependent on the range of the floating point values

$$scale = \frac{\max x_f - \min x_f}{\max x_q - \min x_q}$$

  - The tighter the distribution, the better the accuracy
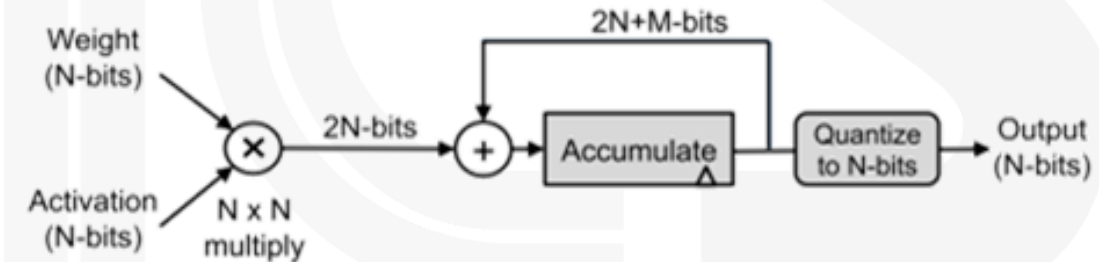  - Luckily, weights tend to have a tight distribution

*Source: Soon Yau*

© Adam Teman, 2020

# Uniform Quantization

- **Uniform quantization is straightforward quantization of floating point to integer**
  - **INT8 add**: 30X less energy, 116X less area than FP32
  - **INT8 multiply**: 18.5X less energy, 27.5X less area than FP32
- **Precision of internal values of MAC is higher than weights and activations**
  - Given $N$-bit weights and inputs →Need $N$x$N$ multiplier→$2N$-bit output product
  - Accumulator: $(2N+M)$-bit

$$M = \log_2 CSR$$

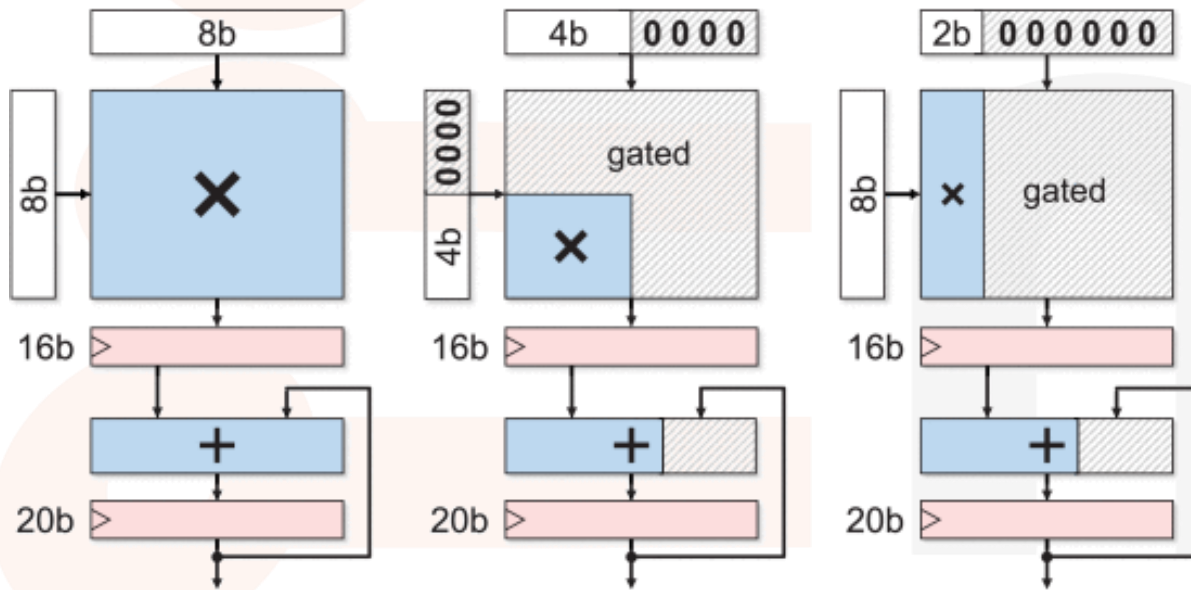  - Final output activation reduced to $N$-bits



- **No significant impact on accuracy if the distribution of weights and activations is centered near zero.**
  - 8-bit arithmetic used in Google's TPU, Nvidia's PASCAL, Intel's NNP-L
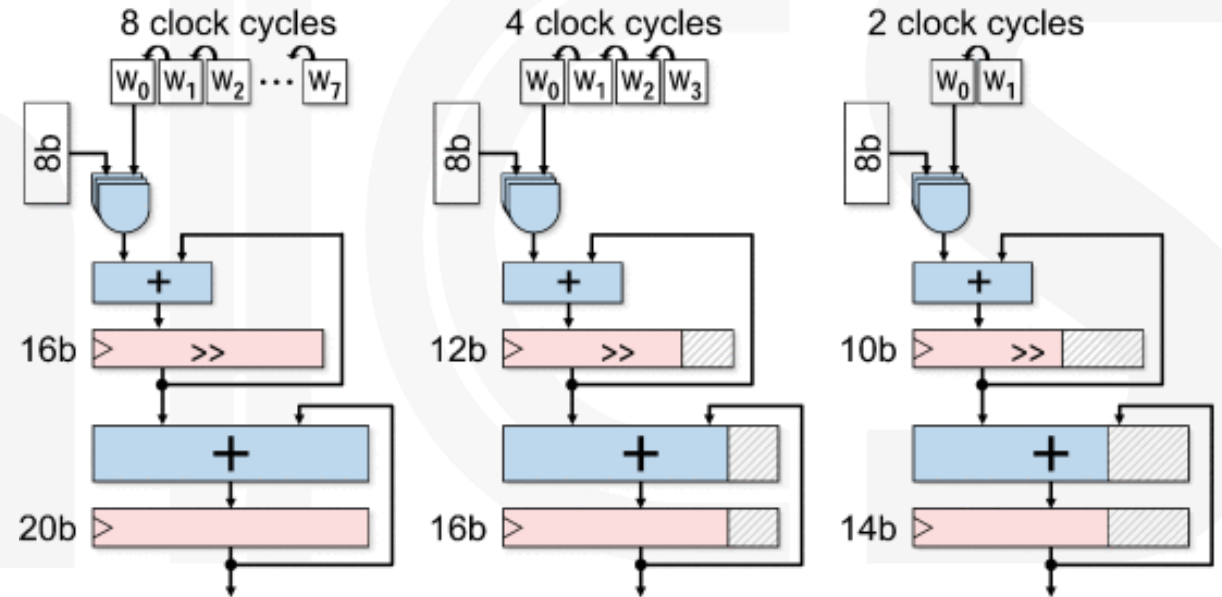
# Configurable MACs for Mixed Precision

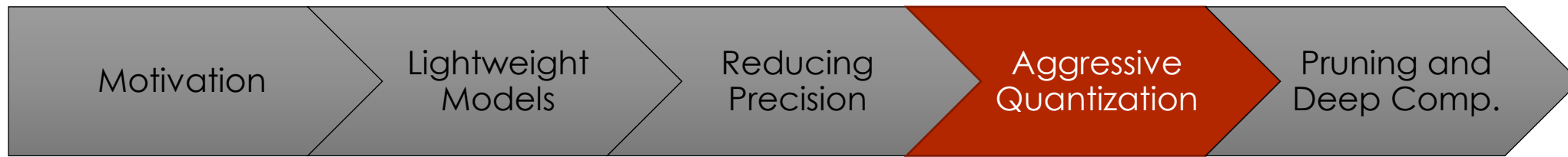- **Use precision-scalable arithmetic for power savings**

**Data Gated MACs**
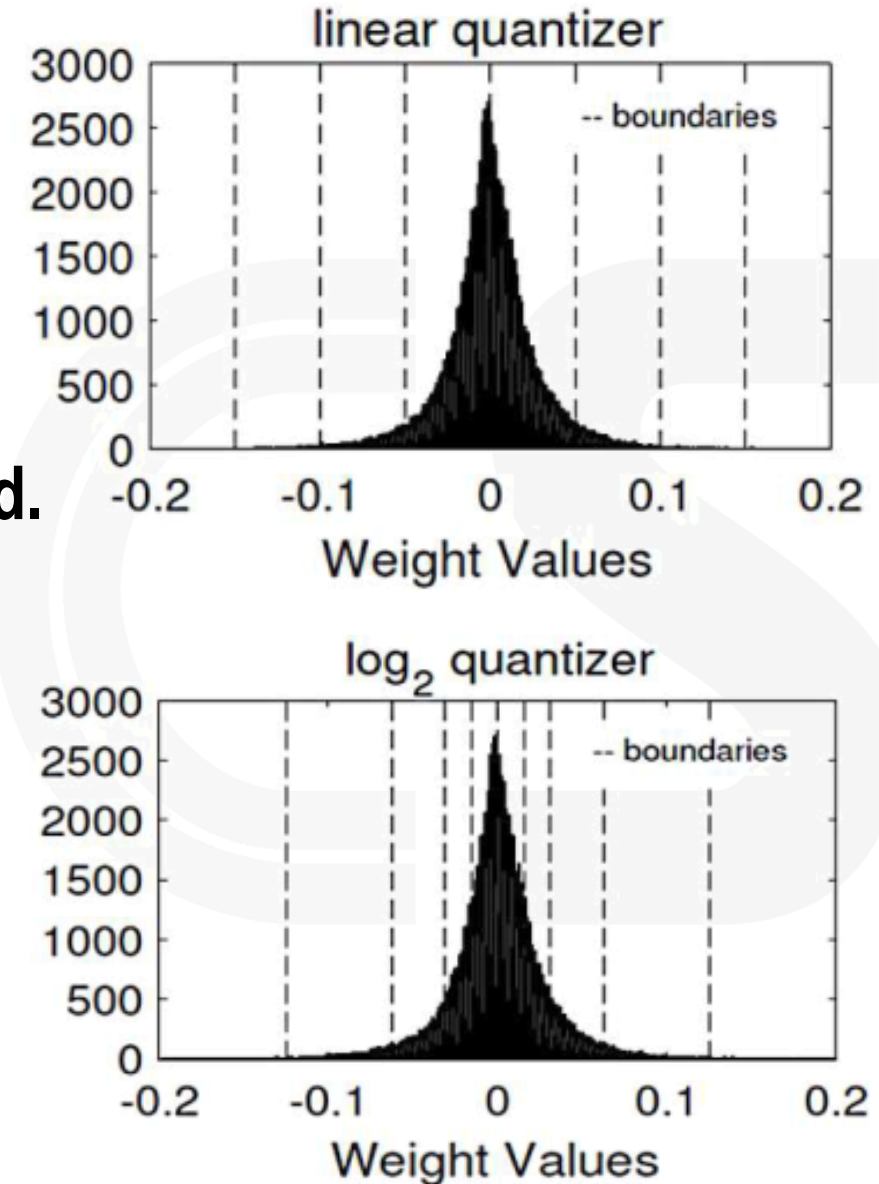
**Bit Serial Designs**

Source: Camus, JETCAS 2019

Source: Camus, JETCAS 2019

- **However, many approaches have overhead that reduce benefits.**

# Aggressive Quantization

EnICS
Emerging Nanoscaled
Integrated Circuits and Systems Labs

Bar-Ilan University
אוניברסיטת בר-אילן

# Non-Uniform Quantization

- **In standard uniform quantization, values are equally spaced out**

- **However, computing a quantization that better fits the distribution, better accuracy can be achieved.**

  - e.g. with 4-bit *log-domain quantization*, VGG-16 shows only a 5% loss (vs. 28% with uniform quantization)

- **Log-domain quantization further allows replacing multiplication with bit-shift**

- **Weight sharing, for example through *learned quantization*, can provide an even better solution.**
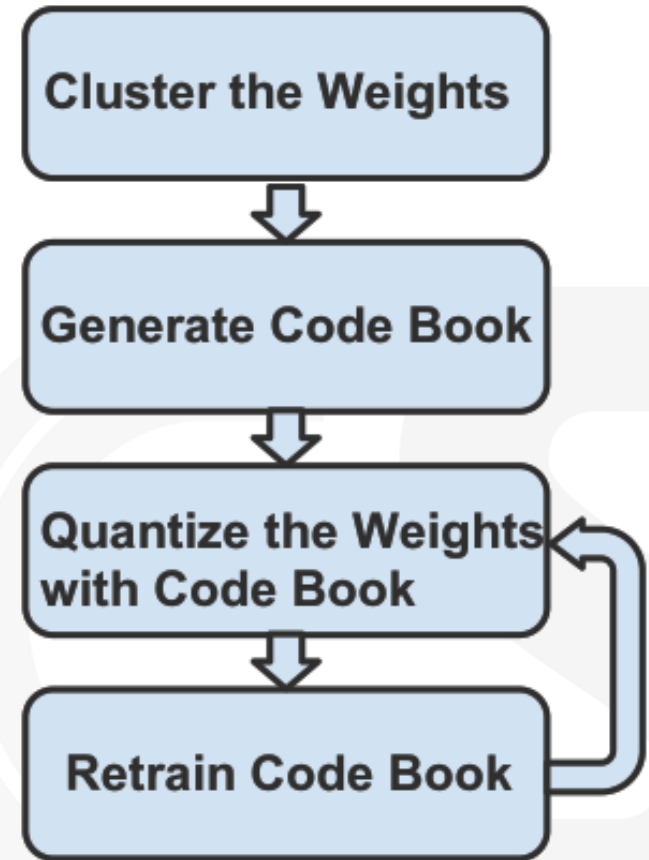


*Source: Camus, Lee, ICASSP 2017*

# Trained Quantization

2.09, 2.12, 1.92, 1.87

↓

2.0

Cluster the Weights

↓

Generate Code Book

↓

Quantize the Weights with Code Book

↓

Retrain Code Book

**32 bit**

**4bit** 8x less memory footprint

[Han et al. ICLR'16]

# Trained Quantization

**Weights (FP32)**

| | | | |
|------|-------|------|-------|
| 2.09 | -0.98 | 1.48 | 0.09 |
| 0.05 | -0.14 | -1.08 | 2.12 |
| -0.91 | 1.92 | 0 | -1.03 |
| 1.87 | 0 | 1.53 | 1.49 |

**cluster** →

| | | | |
|---|---|---|---|
| 3 | 0 | 2 | 1 |
| 1 | 1 | 0 | 3 |
| 0 | 3 | 1 | 0 |
| 3 | 1 | 2 | 2 |

**centroid** →

| |
|------|
| 2.00 |
| 1.50 |
| 0.00 |
| -1.00 |

→ ( **-** ) →

| |
|------|
| 1.96 |
| 1.48 |
| -0.04 |
| -0.97 |

**x learning rate**

**Gradient (FP32)**

| | | | |
|-------|-------|-------|-------|
| -0.03 | 0.01 | 0.03 | 0.02 |
| -0.01 | 0.01 | -0.02 | 0.12 |
| -0.01 | 0.02 | 0.04 | 0.01 |
| -0.07 | -0.02 | 0.01 | -0.02 |

**group** →

| | | | |
|-------|------|------|-------|
| -0.03 | 0.12 | 0.02 | -0.07 |

| | | |
|-------|------|-------|
| -0.03 | 0.01 | -0.02 |

| | | | | |
|------|-------|------|------|-------|
| 0.02 | -0.01 | 0.01 | 0.04 | -0.02 |

| | | | |
|-------|-------|-------|------|
| -0.01 | -0.02 | -0.01 | 0.01 |

**reduce** →

| |
|-------|
| 0.04 |
| 0.02 |
| 0.04 |
| -0.03 |

33

# Trained Quantization



- **AlexNet:**
  - 8-bit quantization on CONV layers, 5-bit quantization on FC layers without any loss of accuracy
  - Only 2% loss of accuracy for 4-bit CONV and 2-bit FC layer quantization

- **Need "cookbook" for index translation**
  - See "Deep Compression" later on in the lecture.

*Source: Han*
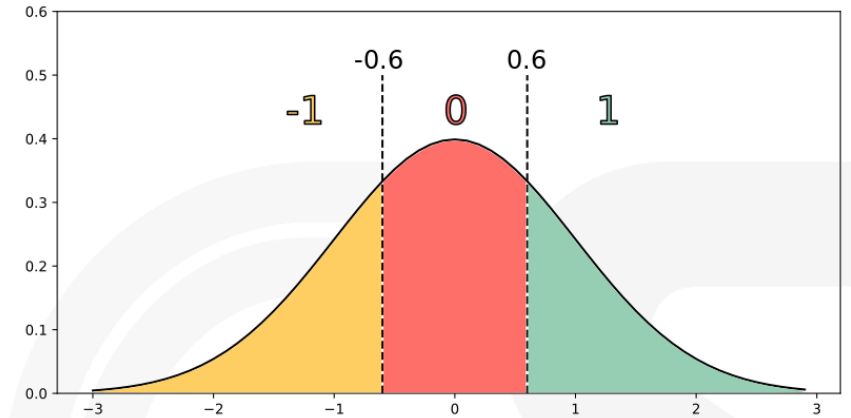
# More aggressive quantization

- **Ternary Connect (2014)**
  - Train with real valued weights
  - Ternarize the weights to $W_B \in \{-H, 0, H\}$

- **Binary Connect (2015)**
  - Binary weights ($W_B \in \{-1, 1\}$), full precision activations
  - Simple multipliers, full precision accumulation
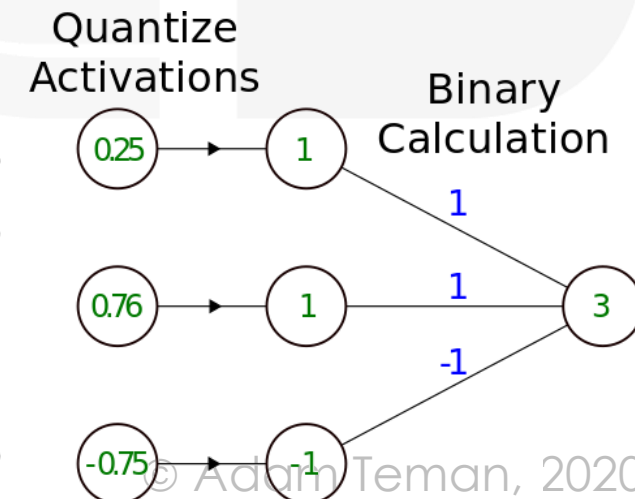  - Training (backprop updates) uses real valued weights ($W_R$) clipped at -1, 1.

- **BinaryNet, Binarized Neural Networks, XNOR-Net (2016)**
  - Binary Weights and Activations
  - Use XNOR for multiplication "popcount" for accumulation
  - Keep first and last layers at full precision

$$W_B = \text{sign}(W_R)$$

| Encoding (Value) | | XNOR (Multiply) |
|---|---|---|
| 0 (−1) | 0 (−1) | 1 (+1) |
| 0 (−1) | 1 (+1) | 0 (−1) |
| 1 (+1) | 0 (−1) | 0 (−1) |
| 1 (+1) | 1 (+1) | 1 (+1) |

# Summary

| 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

# Pruning and Deep Compression

Emerging Nanoscaled
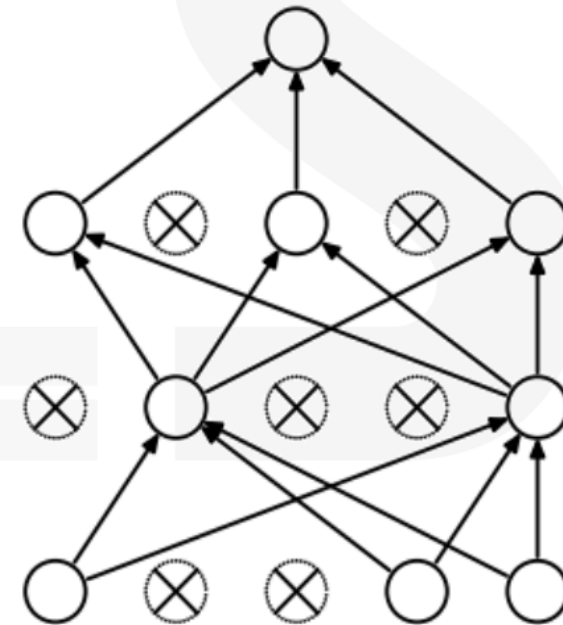Integrated Circuits and Systems Labs

Bar-Ilan University
אוניברסיטת בר-אילן

# Precursor: Dropout

- **A well-known technique for eliminating overfitting is called "Dropout"**
  - During each iteration of training,
    zero out a random fraction of nodes in fully connected layers
  - During inference, use all connections
- **But regularization through batch normalization has almost made this unnecessary.**

- **However, it raises the question:**
  **"Do we actually need all synapses?"**

(a) Standard Neural Net

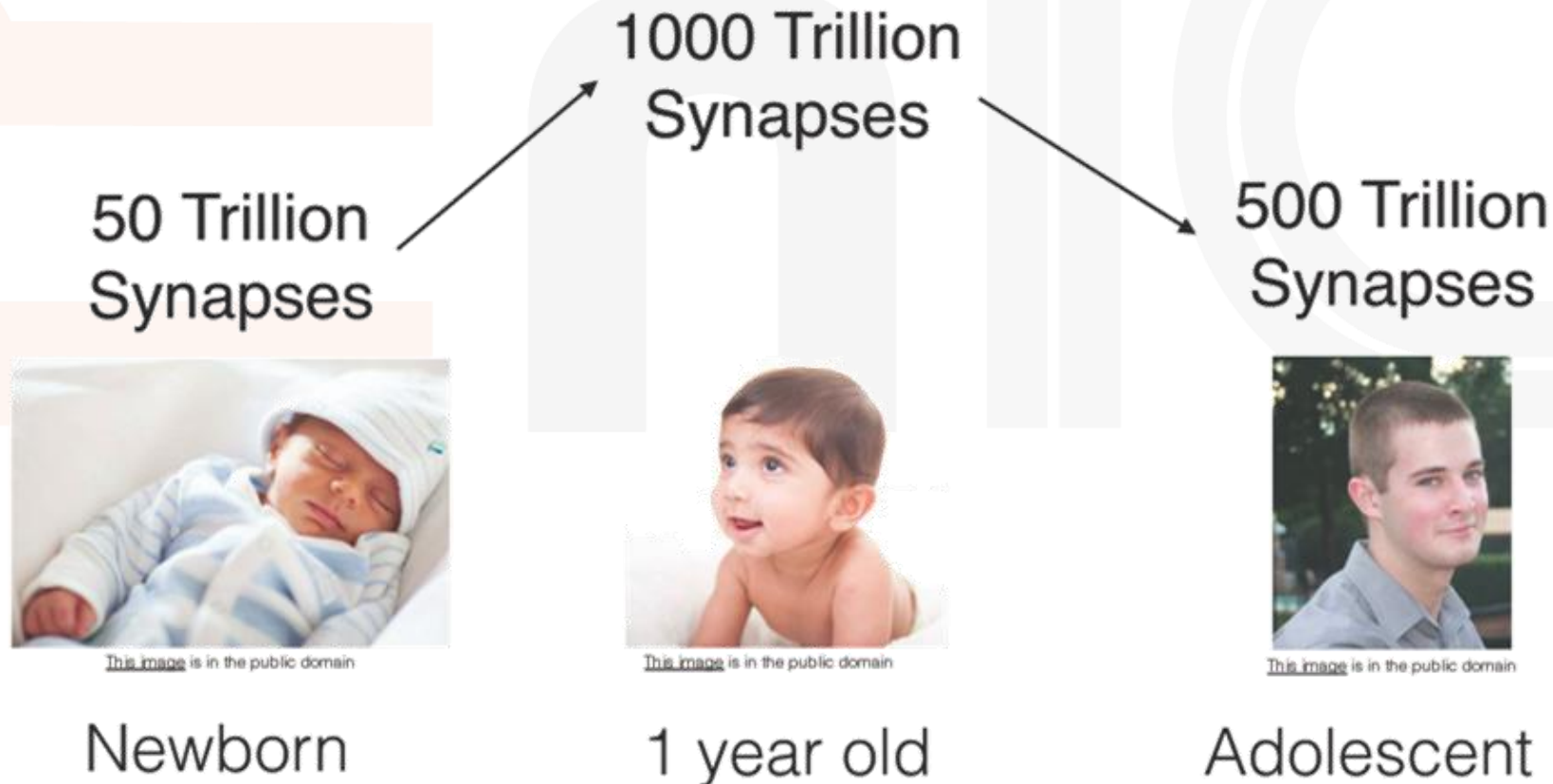(b) After applying dropout.

*Source: Srivastava, et al.*
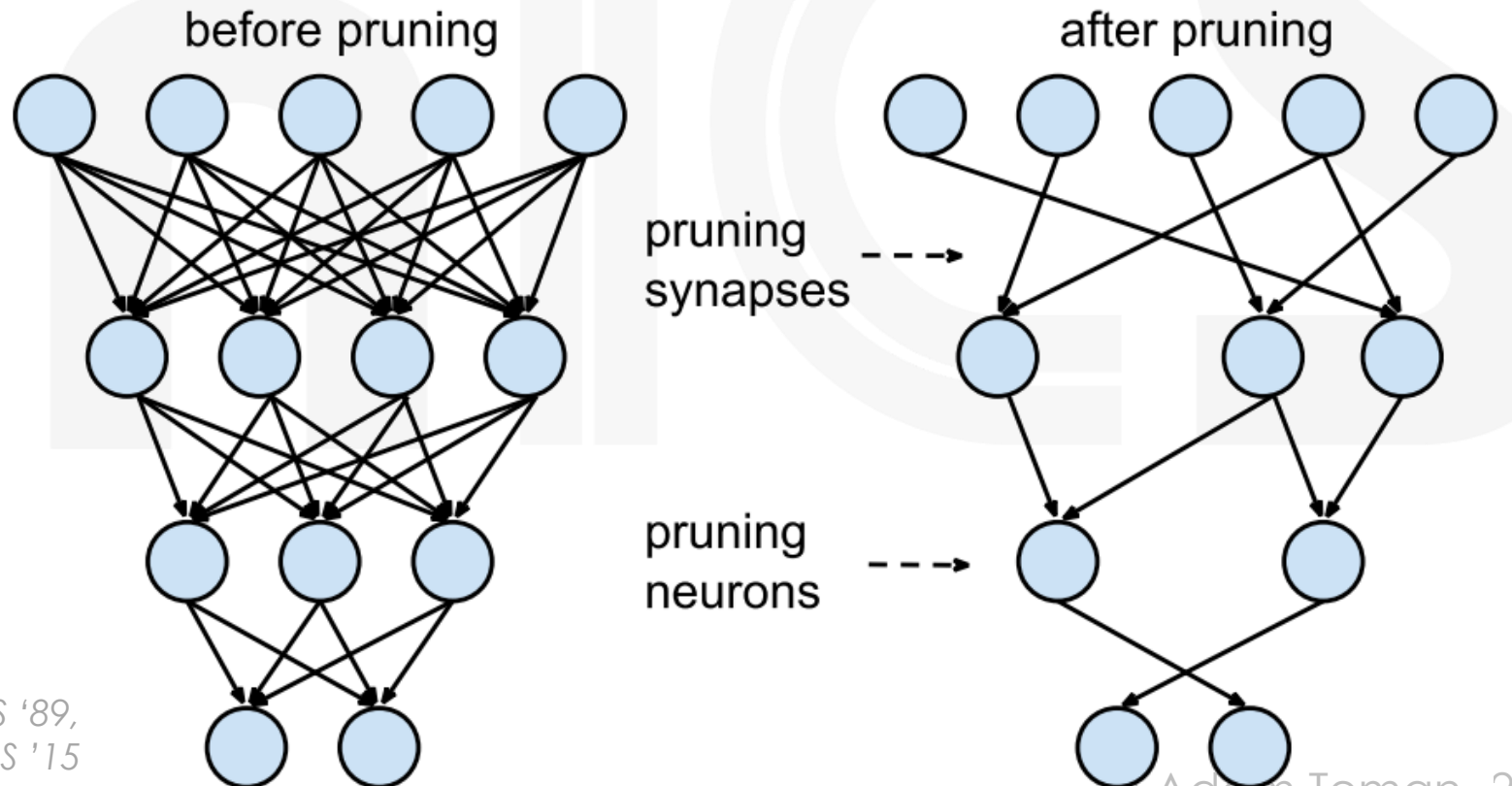
# Synaptic Pruning

- **The human (all mammals) body prunes synapses**
  - Axons and dentrites completely decay and die off during lifetime
  - Starts near birth and continues into the mid-20s



1000 Trillion Synapses

50 Trillion Synapses

500 Trillion Synapses

This image is in the public domain

Newborn

1 year old

Adolescent

*Source: Walsh, Nature 2013*
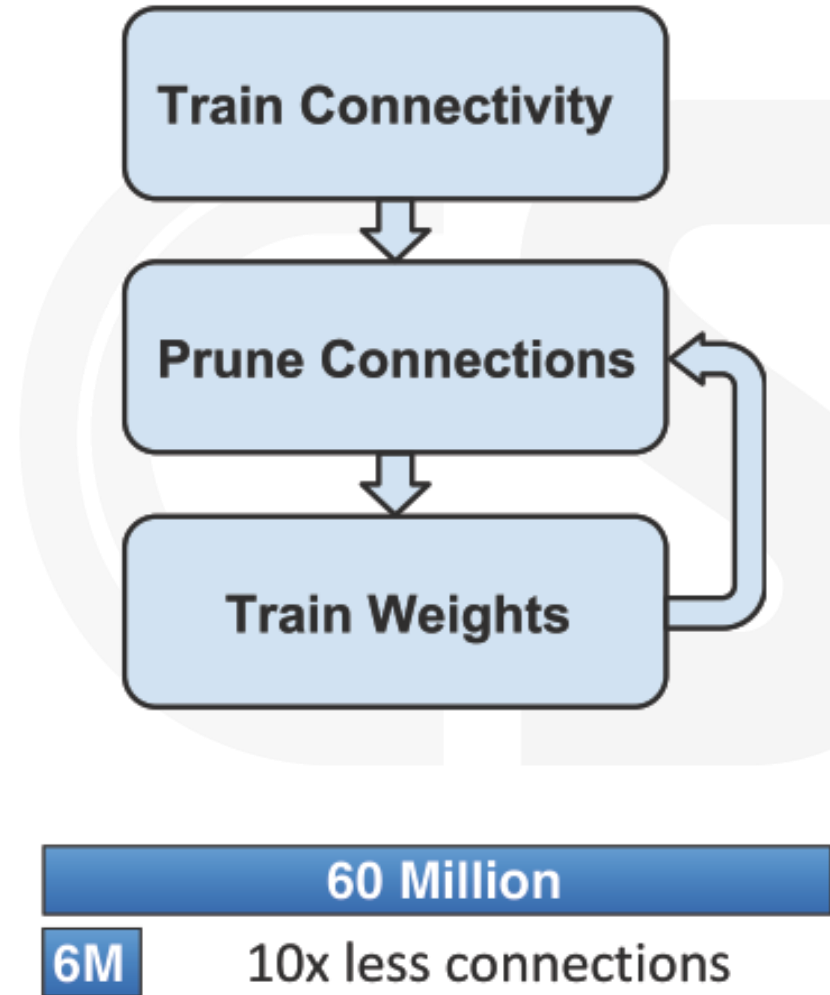
# Optimal Brain Damage

- **In 1989, Yann Lecun suggested pruning neural networks**
  - Compute the impact of each weight on the training loss = weight saliency
  - Remove low-saliency weights and fine tune remaining weights

- **Unlike in "Dropout", pruned synapses are removed for good.**



*Source: Lecun, NIPS '89,*
*Han, NIPS '15*
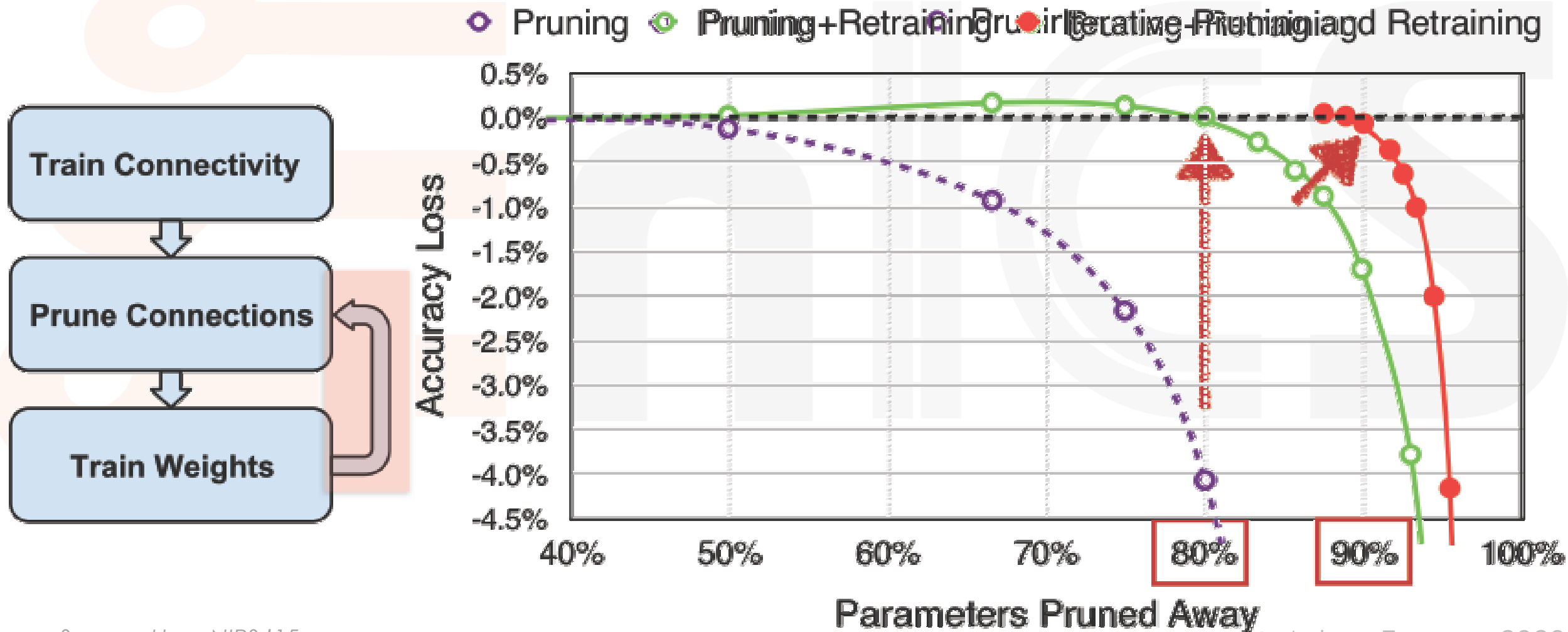
# Pruning Deep Neural Networks

- **Pruning DNNs leads to sparsity**
  - Easier to compress
  - Skip multiplications by zero

- **Han, et al., showed that 90% of the connections in AlexNet can be pruned without incurring accuracy loss!**
  - Weights were pruned below a threshold

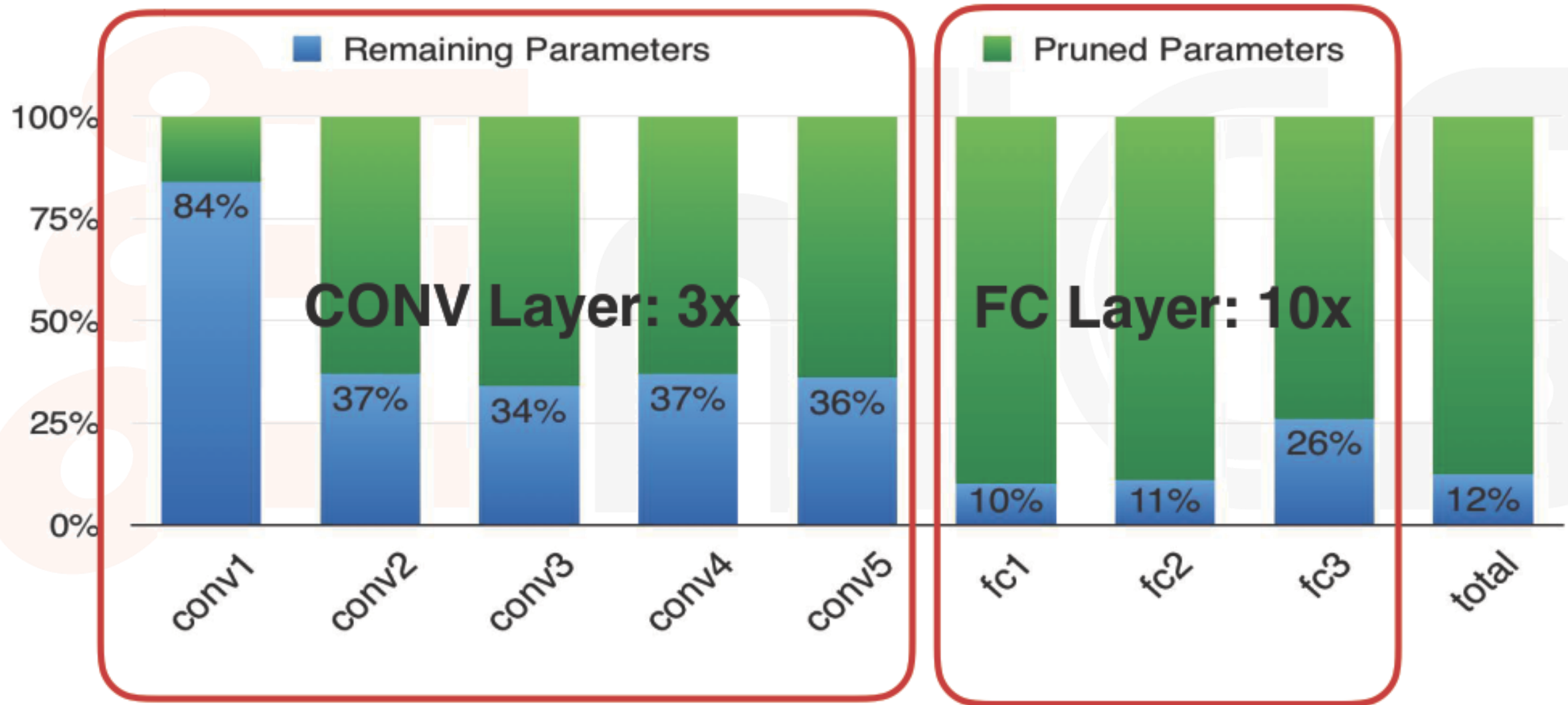- **The Train-Prune-Retrain pipeline was used**



Train Connectivity

Prune Connections

Train Weights

**60 Million**

**6M** 10x less connections

*Source: Han, NIPS '15*

# Pruning Deep Neural Networks

• **Iteratively Retrain to recover accuracy**



*Source: Han, NIPS '15* © Adam Teman, 2020

# Pruning AlexNet



*Source: Han, NIPS '15* © Adam Teman, 2020

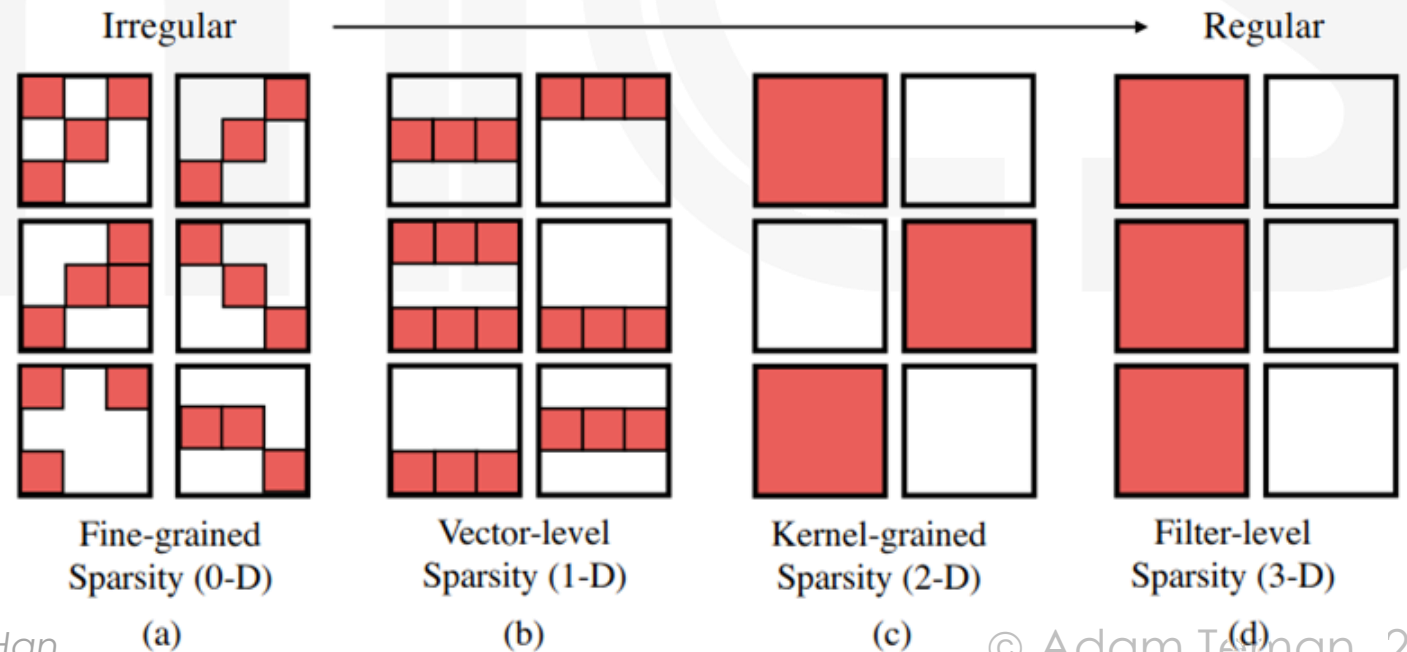# Pruning Changes Weight Distribution



**Before Pruning**     **After Pruning**     **After Retraining**

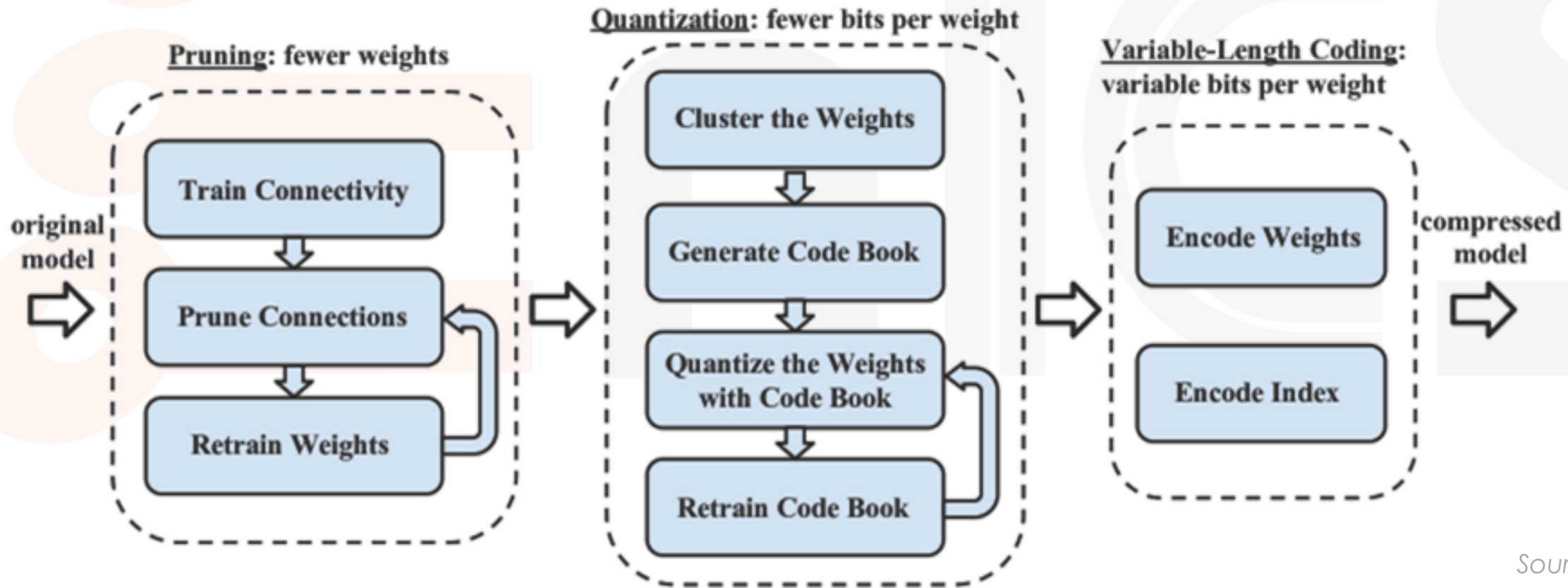Conv5 layer of Alexnet. Representative for other network layers as well.

# Hardware Efficiency Considerations in Pruning

- **Pruning leads to irregularity, which is difficult to parallelize in hardware**

- **Load-balance aware pruning**
  - Sort the weights in every sub-matrix and prune the same amount in each, such that each PE works on the same number of non-zero weights
  - Need to index every non-zero weight

- **Pruning with structure**
  - Prune by rows/columns, kernels, or whole filters
  - Can index a larger space
  - For example, prune a column according to L2 norm



Irregular &rarr; Regular

Fine-grained Sparsity (0-D) (a)    Vector-level Sparsity (1-D) (b)    Kernel-grained Sparsity (2-D) (c)    Filter-level Sparsity (3-D) (d)

*Source: Han*

# Deep Compression

- **Deep Compression combines pruning, trained quantization and variable length coding in a pipeline:**



*Source: Han*

# Storing the Meta Data

- **How do we store the index and weight?**
  - For each non-zero weight store the weight and the index
  - Instead of the actual index, store the distance from the previous non-zero index
  - Select a small bit-width for the index representation – if the span is larger, then pad with zeros.
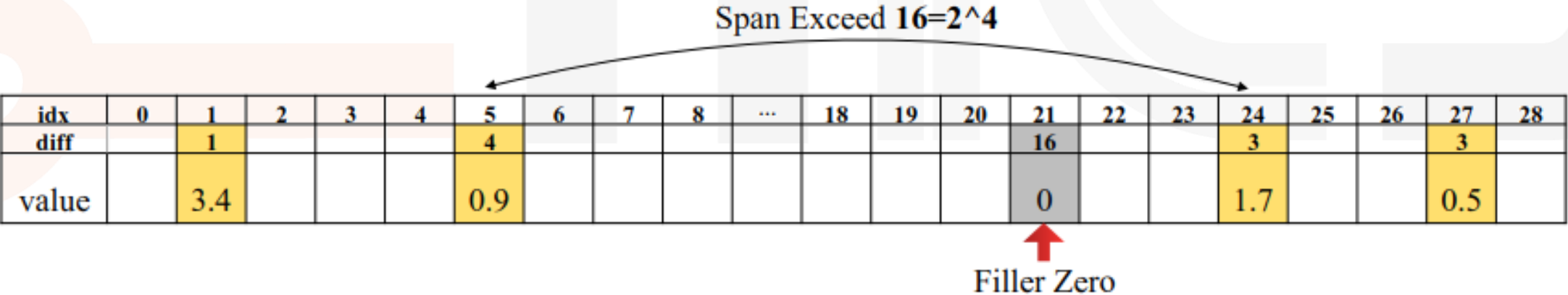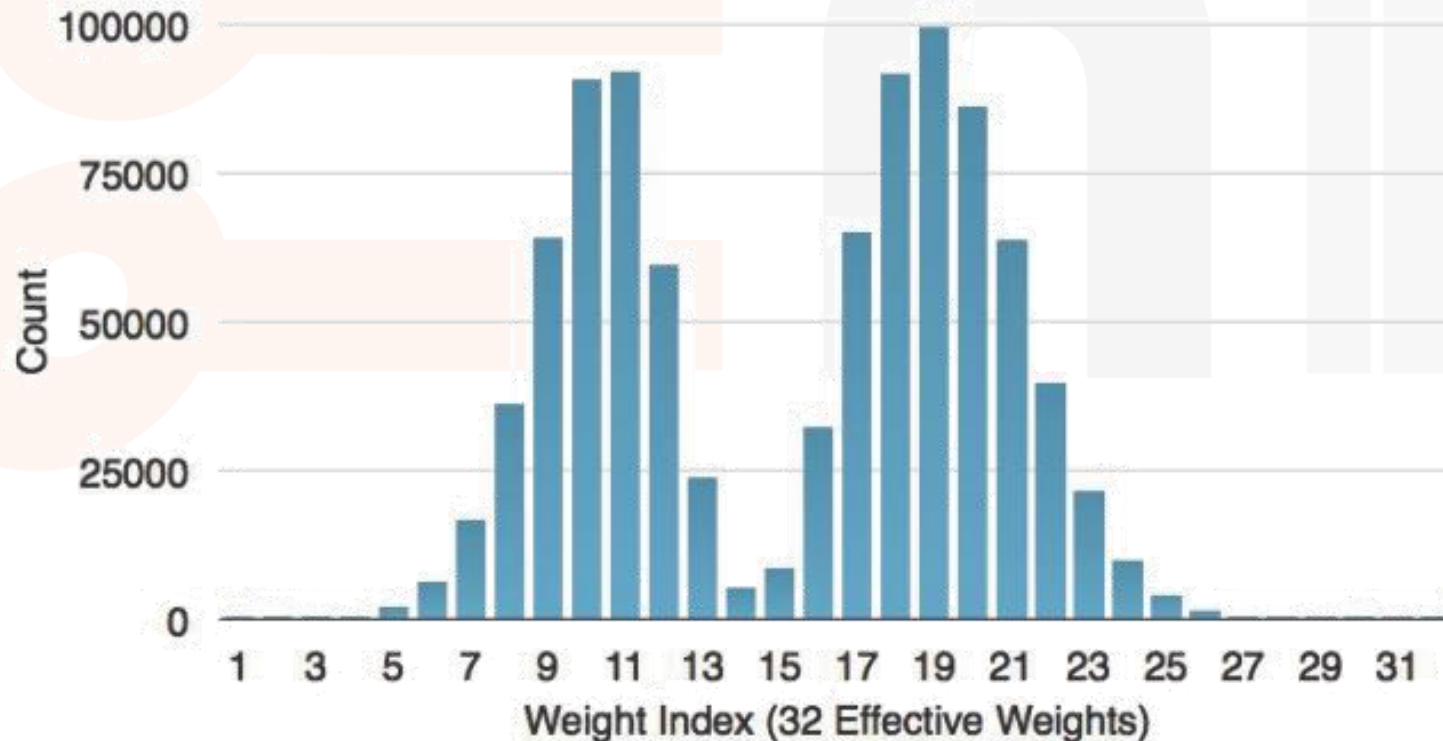- **A separate codebook is stored for each layer**

Span Exceed **16=2^4**

| idx | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | ... | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| diff | | 1 | | | | 4 | | | | | | | | 16 | | | 3 | | | 3 | |
| value | | 3.4 | | | | 0.9 | | | | | | | | 0 | | | 1.7 | | | 0.5 | |

Filler Zero

Figure 4.5: Pad a filler zero to handle overflow when representing a sparse vector with relative index.

# Variable-Length Coding

- **The idea is:**
  - Infrequent weights: use more bits to represent
  - Frequent weights: use less bits to represent
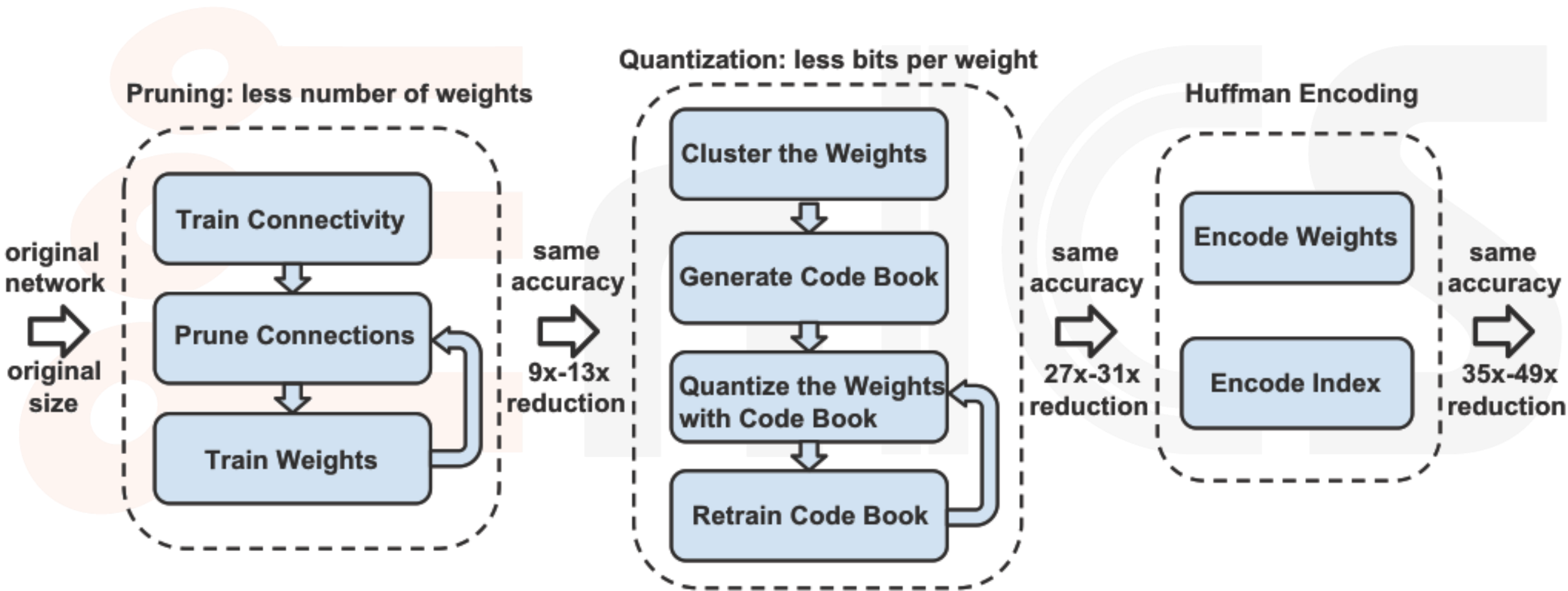- **Huffman coding is used for Deep Compression.**



**Huffman Encoding**

Encode Weights

Encode Index

*Source: Han*

# Summary of Deep Compression
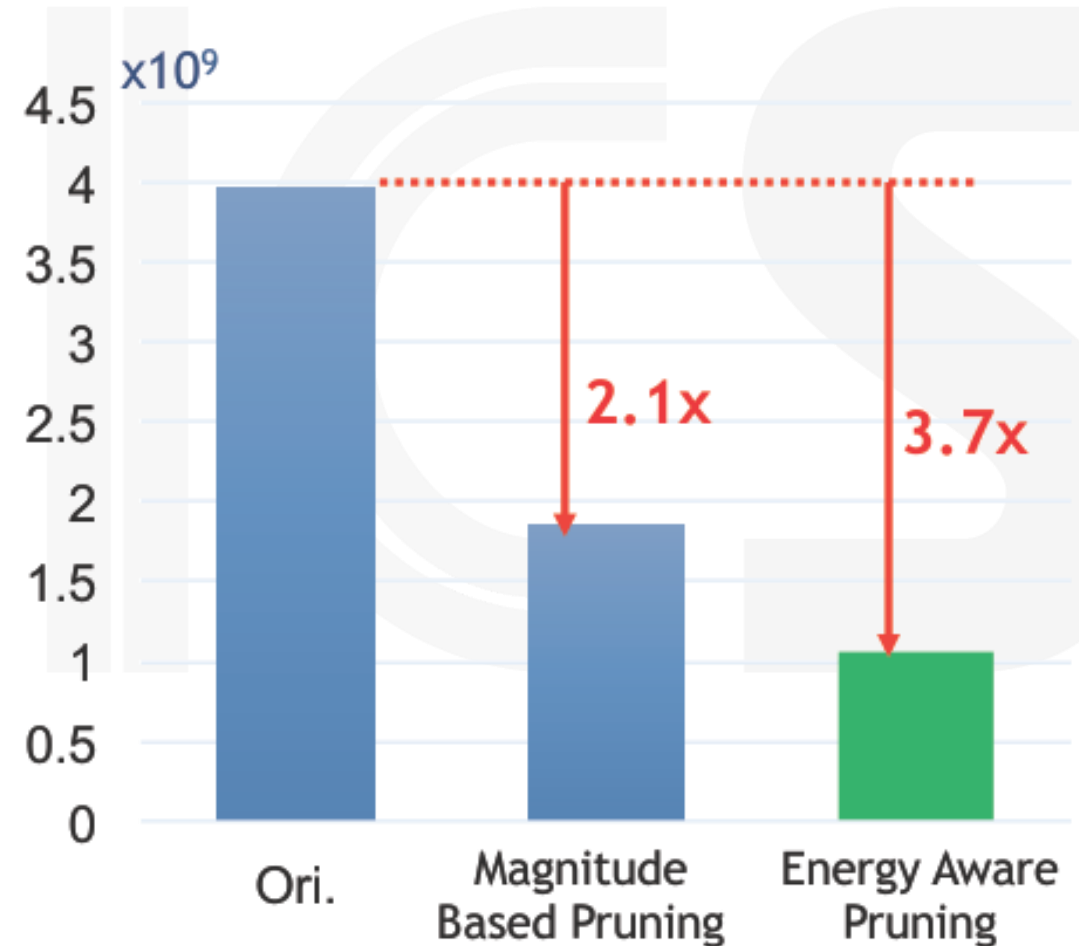
# Results: Compression Ratio

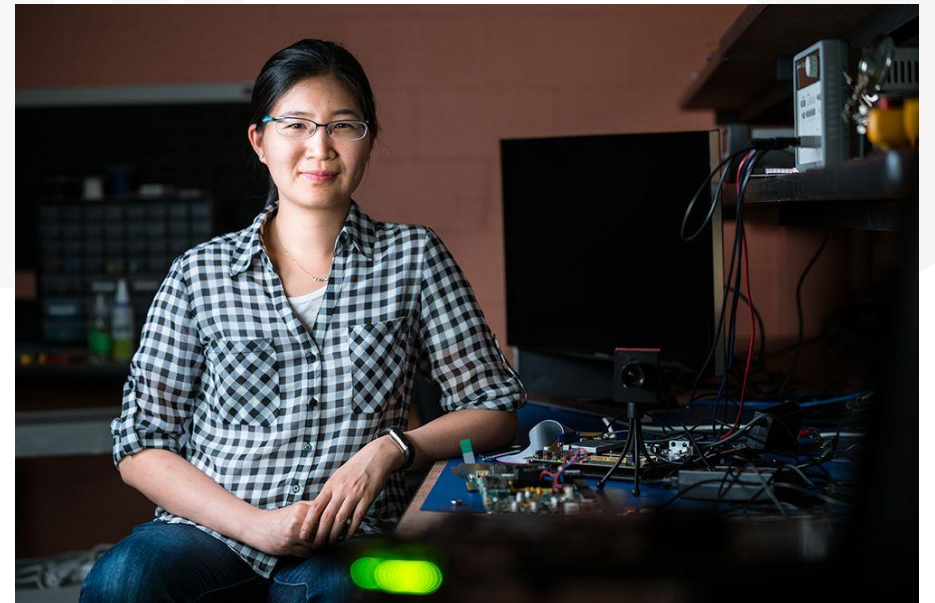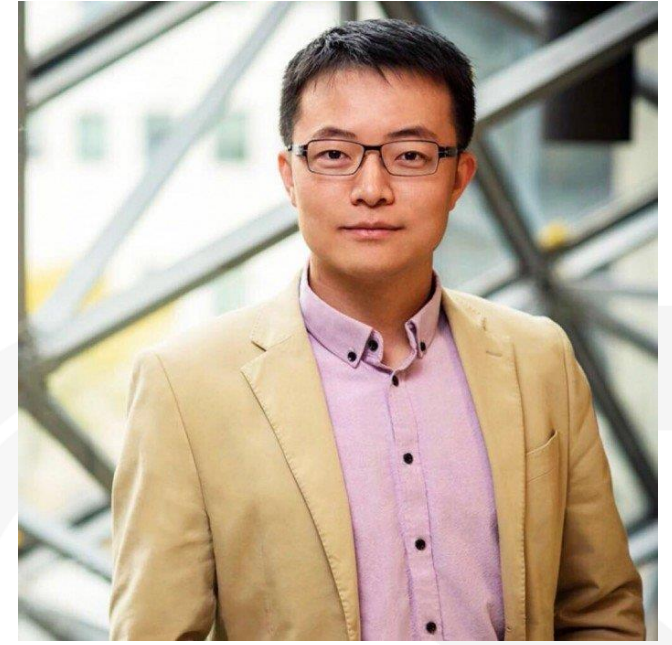| Network | Original Size | Compressed Size | Compression Ratio | Original Accuracy | Compressed Accuracy |
|---|---|---|---|---|---|
| LeNet-300 | 1070KB → | 27KB | **40x** | 98.36% → | 98.42% |
| LeNet-5 | 1720KB → | 44KB | **39x** | 99.20% → | 99.26% |
| AlexNet | 240MB → | 6.9MB | **35x** | 80.27% → | 80.30% |
| VGGNet | 550MB → | 11.3MB | **49x** | 88.68% → | 89.09% |
| GoogleNet | 28MB → | 2.8MB | **10x** | 88.90% → | 88.92% |
| ResNet-18 | 44.6MB → | 4.0MB | **11x** | 89.24% → | 89.28% |

*Source: Han*

# Energy-Aware Pruning

- **The value of weights alone is not a good metric for energy**
  - Instead prune according to energy.
- **Sort layers based on energy and prune layers that consume the most energy first**
- **Energy-aware pruning reduces AlexNet energy by 3.7x and outperforms the previous work that uses magnitude-based pruning by 1.7x**

*Source: Sze, MIT*          © Adam Teman, 2020

# Main References

- **Song Han, various talks**

- **Vivienne Sze, various talks**

- **Bill Dally, various talks**

- **Towards Data Science:**
  - Bharath Raj
  - Yusuke Uchida
  - Arthur Douillard
  - Sik-Ho Tsang
  - Chi-Feng Wang
  - Ranjeet Singh
  - others

*Source: MIT*