# CMSC5743

## L06: Binary/Ternary Network

**Bei Yu**

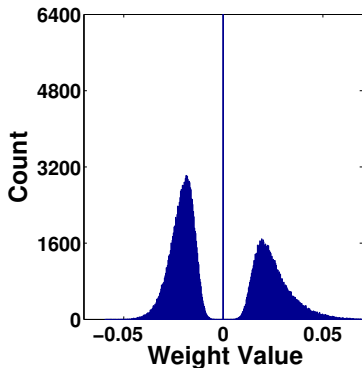(Latest update: November 2, 2020)

Fall 2020

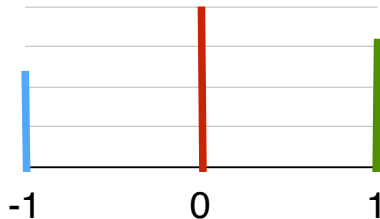**These slides contain/adapt materials developed by**

- ▶ Ritchie Zhao et al. (2017). "Accelerating binarized convolutional neural networks with software-programmable FPGAs". In: *Proc. FPGA*, pp. 15–24

- ▶ Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542
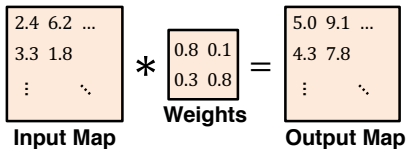
# Binary / Ternary Net: Motivation

# Binarized Neural Networks (BNN)

**CNN**



$$\begin{bmatrix} 2.4 & 6.2 & ... \\ 3.3 & 1.8 & \\ \vdots & & \ddots \end{bmatrix} * \begin{bmatrix} 0.8 & 0.1 \\ 0.3 & 0.8 \end{bmatrix} = \begin{bmatrix} 5.0 & 9.1 & ... \\ 4.3 & 7.8 & \\ \vdots & & \ddots \end{bmatrix}$$
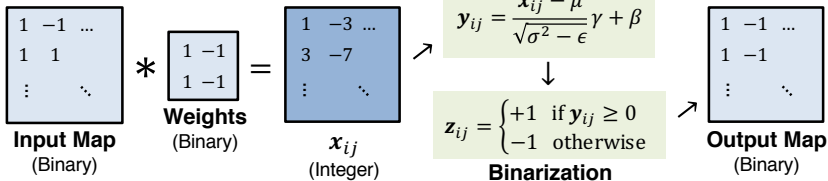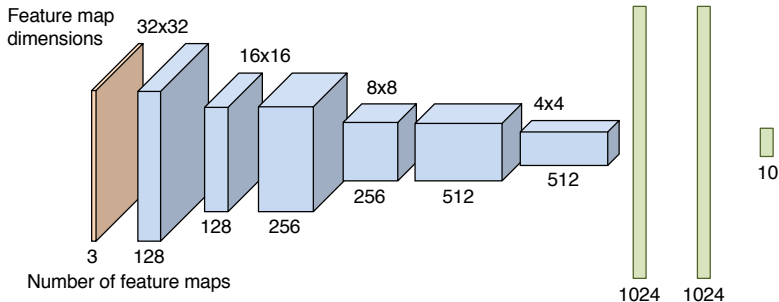
**Input Map**     **Weights**     **Output Map**

## Key Differences

1. Inputs are binarized (−1 or +1)
2. Weights are binarized (−1 or +1)
3. Results are binarized after **batch normalization**

**BNN**

$$\begin{bmatrix} 1 & -1 & ... \\ 1 & 1 & \\ \vdots & & \ddots \end{bmatrix} * \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} 1 & -3 & ... \\ 3 & -7 & \\ \vdots & & \ddots \end{bmatrix}$$

**Input Map**   (Binary)    **Weights** (Binary)    $x_{ij}$ (Integer)

**Batch Normalization**

$$y_{ij} = \frac{x_{ij} - \mu}{\sqrt{\sigma^2 - \epsilon}}\gamma + \beta$$

↓

$$z_{ij} = \begin{cases} +1 & \text{if } y_{ij} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

**Binarization**

$$\begin{bmatrix} 1 & -1 & ... \\ 1 & -1 & \\ \vdots & & \ddots \end{bmatrix}$$

**Output Map** (Binary)

6

# BNN CIFAR-10 Architecture [2]



Feature map dimensions

32x32
16x16
8x8
4x4

3  128  128  256  256  512  512

Number of feature maps

1024  1024

10

- 6 conv layers, 3 dense layers, 3 max pooling layers
- All conv filters are 3x3
- First conv layer takes in floating-point input
- **13.4 Mbits total model size** (after hardware optimizations)

[2] M. Courbariaux et al. **Binarized Neural Networks: Training Deep Neural Networks with Weights and Activations Constrained to +1 or -1**. *arXiv:1602.02830*, Feb 2016.

7

# Advantages of BNN

## 1. Floating point ops replaced with binary logic ops

| $b_1$ | $b_2$ | $b_1 \times b_2$ |
|----|----|-------|
| +1 | +1 | +1 |
| +1 | −1 | −1 |
| −1 | +1 | −1 |
| −1 | −1 | +1 |

| $b_1$ | $b_2$ | $b_1$ XOR $b_2$ |
|----|----|-------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

– Encode {+1,−1} as {0,1} → multiplies become XORs
– Conv/dense layers do dot products → XOR and popcount
– Operations can map to LUT fabric as opposed to DSPs

## 2. Binarized weights may reduce total model size
– Fewer bits per weight may be offset by having more weights

8

# BNN vs CNN Parameter Efficiency

| Architecture | Depth | Param Bits (Float) | Param Bits (Fixed-Point) | Error Rate (%) |
|---|---|---|---|---|
| ResNet [3] (CIFAR-10) | 164 | 51.9M | 13.0M* | 11.26 |
| BNN [2] | 9 | - | 13.4M | 11.40 |

\* Assuming each float param can be quantized to 8-bit fixed-point

▶ **Comparison:**
 – Conservative assumption: ResNet can use 8-bit weights
 – BNN is based on VGG (less advanced architecture)
 – BNN seems to hold promise!

[2] M. Courbariaux et al. **Binarized Neural Networks: Training Deep Neural Networks with Weights and Activations Constrained to +1 or -1**. *arXiv:1602.02830*, Feb 2016.
[3] K. He, X. Zhang, S. Ren, and J. Sun. **Identity Mappings in Deep Residual Networks.** *ECCV 2016.*

9

Minimize the Quantization Error

Reduce the Gradient Error

# Overview

Minimize the Quantization Error

Reduce the Gradient Error

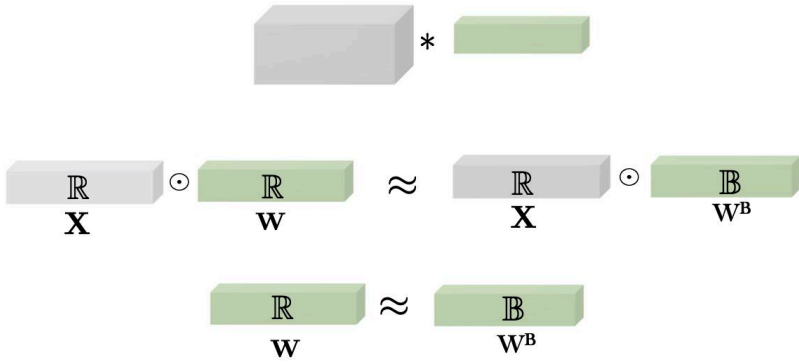| | Operations | Memory | Computation |
|---|---|---|---|
| $\mathbb{R}$ * $\mathbb{R}$ | + − × | 1x | 1x |

Binary Weight Networks

XNOR-Networks

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

| | | | Operations | Memory | Computation |
|---|---|---|---|---|---|
| $\mathbb{R}$ | * | $\mathbb{R}$ | + − × | 1x | 1x |
| $\mathbb{R}$ | * | $\mathbb{B}$ | + − | ~32x | ~2x |
| $\mathbb{B}$ | * | $\mathbb{B}$ | XNOR Bit-count | ~32x | ~58x |

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

$$\mathbf{W^B} = \text{sign}(\mathbf{W})$$

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

# Quantization Error

$$\mathbf{W^B} = \text{sign}(\mathbf{W})$$



$$\left\| \begin{array}{c} \mathbf{W} \\ \mathbb{R} \end{array} - \begin{array}{c} \mathbf{W^B} \\ \mathbb{B} \end{array} \right\| \approx 0.75$$

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

# Optimal Scaling Factor



$$\alpha^*, \mathbf{W^{B}}^* = \arg \min_{\mathbf{W^B}, \alpha} \{ ||\mathbf{W} - \alpha \mathbf{W^B}||^2 \}$$

$$\mathbf{W^{B}}^* = \text{sign}(\mathbf{W})$$
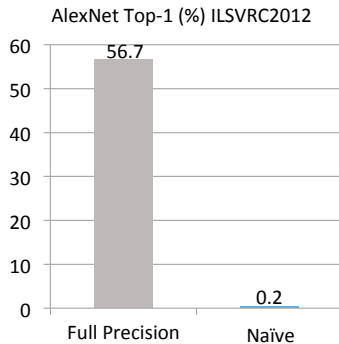$$\alpha^* = \frac{1}{n} ||\mathbf{W}||_{\ell_1}$$

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.
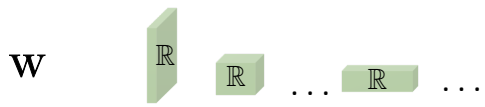
# How to train a CNN with binary filters?



$$\boxed{\mathbb{R}} * \boxed{\mathbb{R}} \approx \left( \boxed{\mathbb{R}} * \boxed{\mathbb{B}} \right) \alpha$$

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

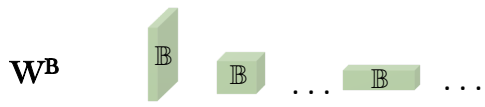# Training Binary Weight Networks

*Naive Solution:*

1. Train a network with real value parameters
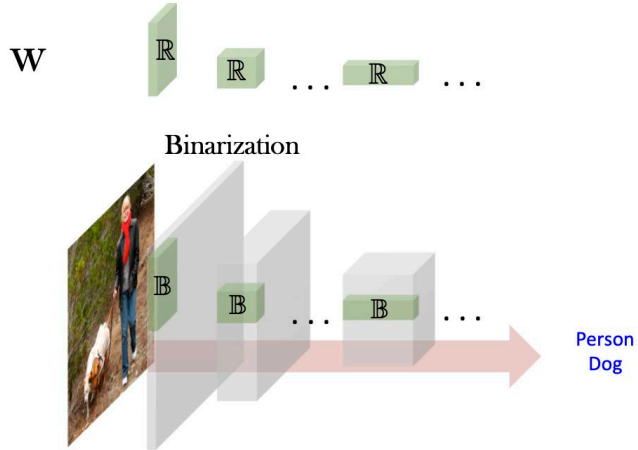2. Binarize the weight filters

---

[1]Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

AlexNet Top-1 (%) ILSVRC2012

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

$\mathbf{W}$

$\mathbb{R}$   $\mathbb{R}$   . . .   $\mathbb{R}$   . . .

Binarization

$\mathbb{B}$   $\mathbb{B}$   . . .   $\mathbb{B}$   . . .

Person
Dog

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.
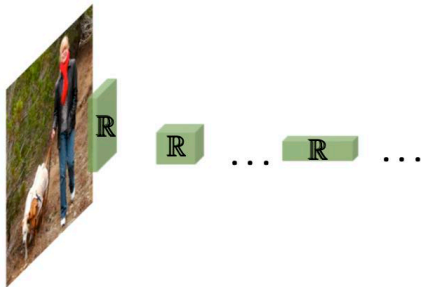
# Binary Weight Network

***Train for binary weights:***

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.    Load a random input image $\mathbf{X}$
4.    $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5.    $\alpha = \frac{\|W\|_{\ell_1}}{n}$
6.    Forward pass with $\alpha, \mathbf{W^B}$
7.    Compute loss function $\mathbf{C}$
8.    $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} = $ Backward pass with $\alpha, \mathbf{W^B}$
9.    Update $\mathbf{W}$ ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)
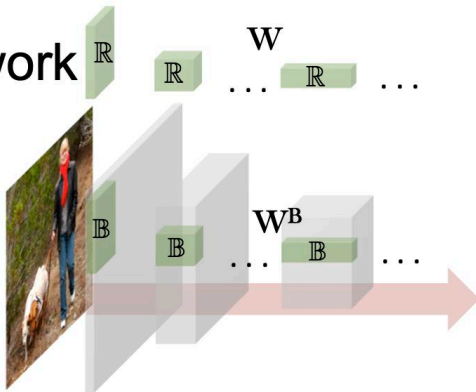
[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

# Binary Weight Network

$\mathbf{W}$

*Train for binary weights:*

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3. Load a random input image $\mathbf{X}$
4. $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5. $\alpha = \frac{\|W\|_{\ell_1}}{n}$
6. Forward pass with $\alpha, \mathbf{W^B}$
7. Compute loss function $\mathbf{C}$
8. $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} =$ Backward pass with $\alpha, \mathbf{W^B}$
9. Update $\mathbf{W}$ ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)



[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.
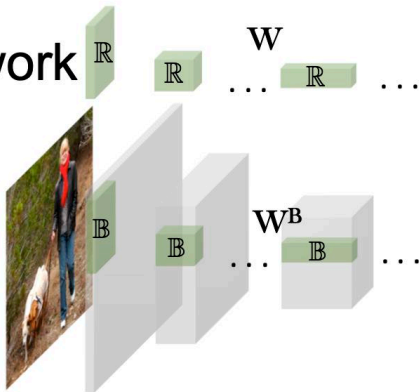
# Binary Weight Network

*Train for binary weights:*

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.    Load a random input image $\mathbf{X}$
4.    $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5.    $\alpha = \frac{\|W\|_{\ell 1}}{n}$
6.    Forward pass with $\alpha, \mathbf{W^B}$
7.    Compute loss function $\mathbf{C}$
8.    $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} =$ Backward pass with $\alpha, \mathbf{W^B}$
9.    Update $\mathbf{W}$ ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)



[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.
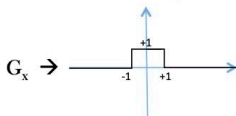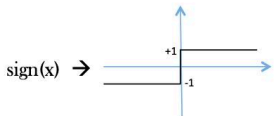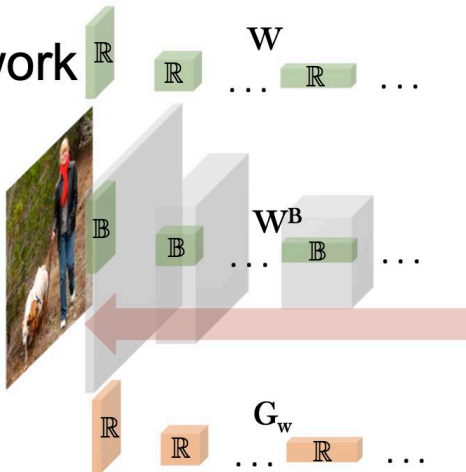
# Binary Weight Network



*Train for binary weights:*

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.     Load a random input image $\mathbf{X}$
4.     $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5.     $\alpha = \frac{\|W\|_{\ell_1}}{n}$
6.     Forward pass with $\alpha, \mathbf{W^B}$
7.     Compute loss function $\mathbf{C}$
8.     $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} = $ Backward pass with $\alpha, \mathbf{W^B}$
9.     Update $\mathbf{W}$ $\left(\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}\right)$
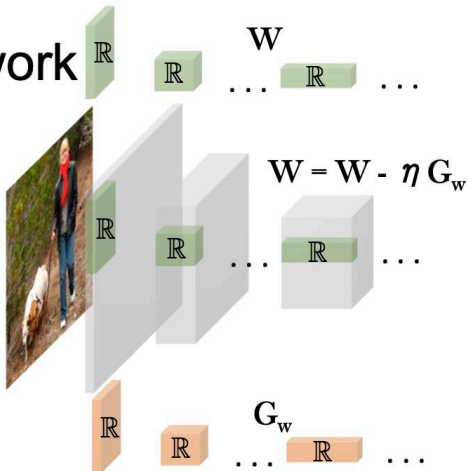
[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.
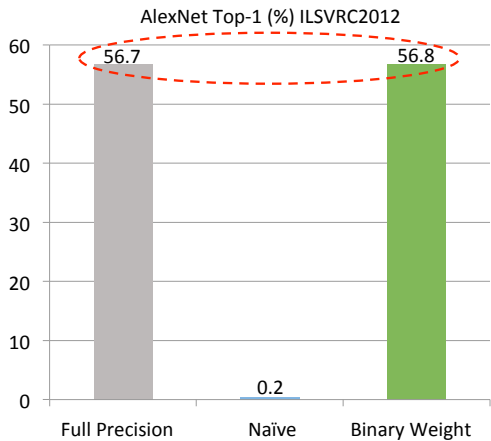
# Binary Weight Network



*Train for binary weights:*

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.     Load a random input image $\mathbf{X}$
4.     $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5.     $\alpha = \frac{\|W\|_{\ell 1}}{n}$
6.     Forward pass with $\alpha, \mathbf{W^B}$
7.     Compute loss function $\mathbf{C}$
8.     $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} =$ Backward pass with $\alpha, \mathbf{W^B}$
9.     Update $\mathbf{W}$ ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

# Binary Weight Network

*Train for binary weights:*

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.    Load a random input image $\mathbf{X}$
4.    $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5.    $\alpha = \frac{\|W\|_{\ell_1}}{n}$
6.    Forward pass with $\alpha, \mathbf{W^B}$
7.    Compute loss function $\mathbf{C}$
8.    $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} = $ Backward pass with $\alpha, \mathbf{W^B}$
9.    Update $\mathbf{W}$ ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)

$\text{sign}(x) \rightarrow$

$G_x \rightarrow$

**LOSS**

$\mathbf{W}$

$\mathbf{W^B}$

$\mathbf{G_w}$

[Hinton et al. 2012]

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

# Binary Weight Network

*Train for binary weights:*

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.    Load a random input image $\mathbf{X}$
4.    $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5.    $\alpha = \frac{\|W\|_{\ell_1}}{n}$
6.    Forward pass with $\alpha, \mathbf{W^B}$
7.    Compute loss function $\mathbf{C}$
8.    $\frac{\partial \mathbf{C}}{\partial \mathbf{W}}$ = Backward pass with $\alpha, \mathbf{W^B}$
9.    Update $\mathbf{W}$ ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)

$\mathbf{W}$

$\mathbf{W} = \mathbf{W} - \eta\, \mathbf{G_w}$

$\mathbf{G_w}$

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

AlexNet Top-1 (%) ILSVRC2012

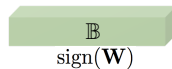(bar chart: Full Precision 56.7, Naïve 0.2, Binary Weight 56.8)

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: Proc. ECCV, pp. 525–542.

| | Operations | Memory | Computation |
|---|---|---|---|
| $\mathbb{R}$ * $\mathbb{R}$ | + − × | 1x | 1x |
| $\mathbb{R}$ * $\mathbb{B}$ | + − | ~32x | ~2x |
| $\mathbb{B}$ * $\mathbb{B}$ XNOR-Networks | XNOR Bit-count | ~32x | ~58x |

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

# Binary Input and Binary Weight (XNOR-Net)
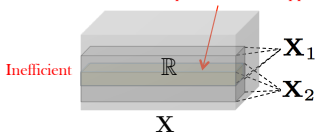


$$\underset{\mathbf{X}}{\mathbb{R}} \odot \underset{\mathbf{W}}{\mathbb{R}} \approx \beta \underset{\mathbf{X^B}}{\mathbb{B}} \odot \alpha \underset{\mathbf{W^B}}{\mathbb{B}}$$

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

# Binary Input and Binary Weight (XNOR-Net)



$$\mathbf{Y} \approx \gamma \, \mathbf{Y^B}$$

$$\mathbf{Y^{B^*}}, \gamma^* = \arg\min_{\mathbf{Y^B}, \gamma} \|\mathbf{Y} - \gamma \mathbf{Y^B}\|^2$$

$$\mathbf{Y^{B^*}} = \text{sign}(\mathbf{Y}) \quad \gamma^* = \frac{1}{n}\|\mathbf{Y}\|_{\ell 1}$$

$$\mathbf{X^{B^*}} = \text{sign}(\mathbf{X}) \quad \mathbf{W^{B^*}} = \text{sign}(\mathbf{W}) \qquad \alpha^* = \frac{1}{n}\|\mathbf{W}\|_{\ell 1} \quad \beta^* = \frac{1}{n}\|\mathbf{X}\|_{\ell 1}$$

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

**(1) Binarizing Weights**

$$\frac{1}{n}\|\mathbf{W}\|_{\ell 1} = \alpha$$

$\mathbb{R}$ $\mathbf{W}$     $\mathbb{B}$ $\text{sign}(\mathbf{W})$

**(2) Binarizing Input**

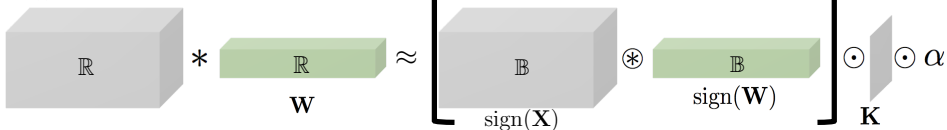Redundant computation in overlapping areas

Inefficient

$\mathbb{R}$    $\mathbf{X}_1$    $\mathbf{X}_2$
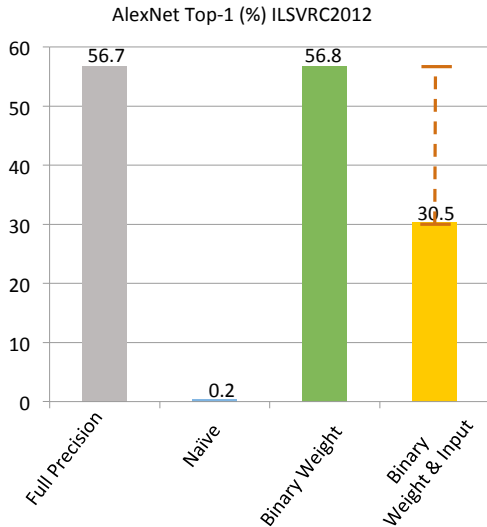
$$\frac{1}{n}\|\mathbf{X}_1\|_{\ell 1} = \beta_1$$
$$\frac{1}{n}\|\mathbf{X}_2\|_{\ell 1} = \beta_2$$

$\mathbf{K}$    $\mathbb{B}$ $\text{sign}(\mathbf{X})$

$\mathbf{X}$

**(2) Binarizing Input**

Efficient

$$\frac{\sum |\mathbf{X}_{:,:,\mathbf{i}}|}{c} =$$

$\xleftarrow{\hspace{2cm}} c \xrightarrow{\hspace{2cm}}$

$* = \begin{cases} \beta_1 \\ \beta_2 \end{cases}$

Average Filter   $\mathbf{K}$

$\mathbb{B}$ $\text{sign}(\mathbf{X})$

**(3) Convolution with XNOR-Bitcount**

$$\mathbb{R} * \underset{\mathbf{W}}{\mathbb{R}} \approx \left[ \underset{\text{sign}(\mathbf{X})}{\mathbb{B}} \circledast \underset{\text{sign}(\mathbf{W})}{\mathbb{B}} \right] \odot \underset{\mathbf{K}}{\phantom{x}} \odot \alpha$$
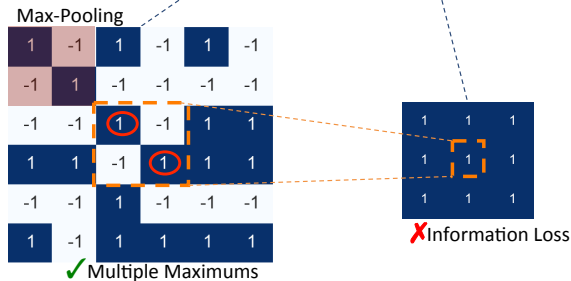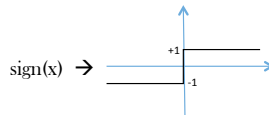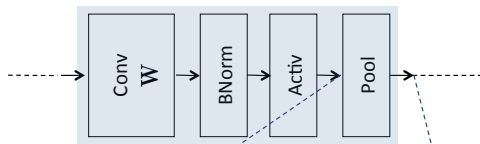
[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.   Load a random input image $\mathbf{X}$
4.   $\mathbf{W^B} = \text{sign}(\mathbf{W})$
5.   $\alpha = \frac{\|W\|_{\ell_1}}{n}$
6.   Forward pass with $\alpha, \mathbf{W^B}$
7.   Compute loss function $\mathbf{C}$
8.   $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} = $ Backward pass with $\alpha, \mathbf{W^B}$
9.   Update $\mathbf{W}$ ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: Proc. ECCV, pp. 525–542.

AlexNet Top-1 (%) ILSVRC2012

[1]Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.
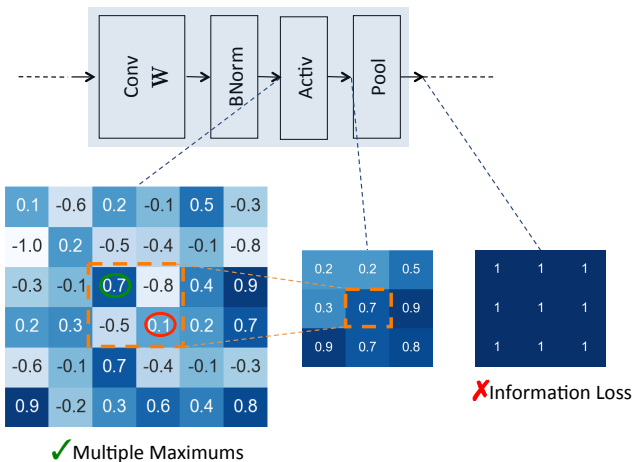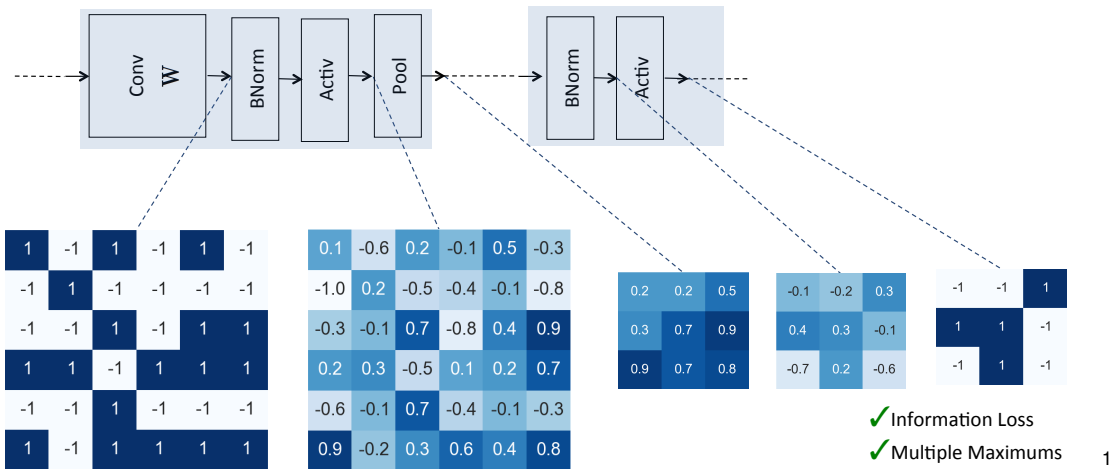
# Network Structure in XNOR-Networks



A typical block in CNN

Max-Pooling

sign(x) →
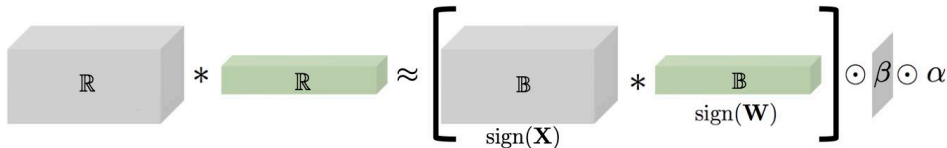
✓ Multiple Maximums

✗ Information Loss

---

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

# Network Structure in XNOR-Networks



✓ Multiple Maximums

✗ Information Loss

[1]Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.
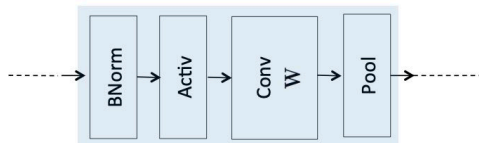
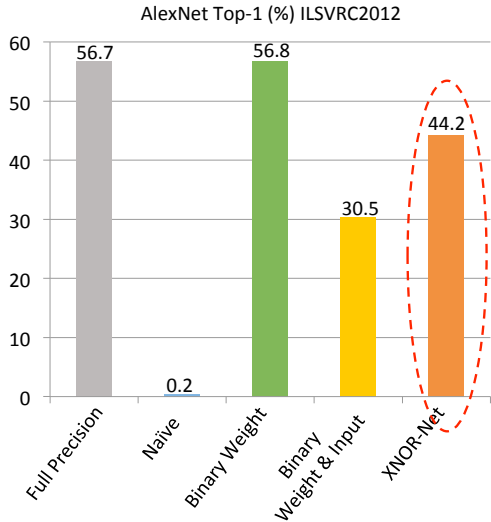# Network Structure in XNOR-Networks



✔ Information Loss
✔ Multiple Maximums

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.
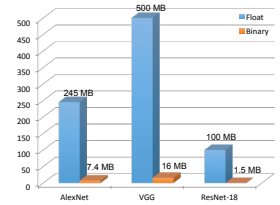
1. Randomly initialize $\mathbf{W}$
2. For $iter = 1$ to $N$
3.     Load a random input image $\mathbf{X}$
4.     $\mathbf{W}^{\mathbf{B}} = \text{sign}(\mathbf{W})$
5.     $\alpha = \frac{\|W\|_{\ell 1}}{n}$
6.     Forward pass with $\alpha, \mathbf{W}^{\mathbf{B}}$
7.     Compute loss function $\mathbf{C}$
8.     $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} = $ Backward pass with $\alpha, \mathbf{W}^{\mathbf{B}}$
9.     Update $\mathbf{W}$ $(\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}})$

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

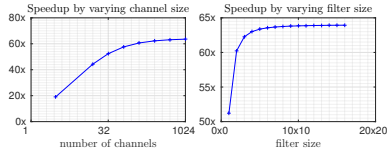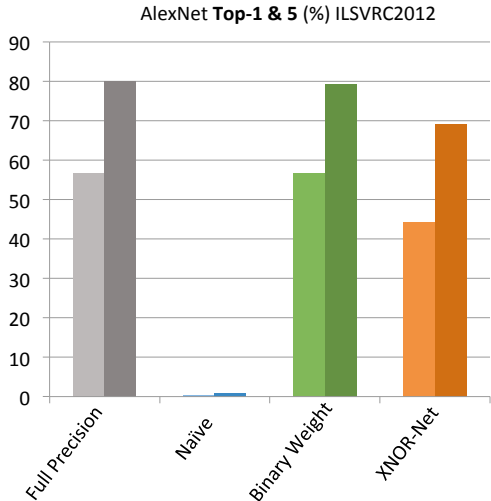AlexNet Top-1 (%) ILSVRC2012

✓ 32x Smaller Model

✓ 58x Less Computation

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

AlexNet **Top-1 & 5** (%) ILSVRC2012

[1] Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

# Motivation and Intuition

## Motivation

▶ Naive methods (Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David (2015). "Binaryconnect: Training deep neural networks with binary weights during propagations". In: *Advances in neural information processing systems*, pp. 3123–3131, Matthieu Courbariaux, Itay Hubara, et al. (2016). "Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1". In: *arXiv preprint arXiv:1602.02830*) suffer the accuracy loss

## Intuition

▶ Quantized parameter should approximate the full precision parameter as closely as possible

Towards Accurate Binary Convolutional Neural Network

# ABC-Net

**Contribution**

▶ Approximate full-precision weights with the linear combination of multiple binary weight bases

▶ Introduce multiple binary activations

# ABC-Net

# ABC-Net

## Forward and Backward

▶ Forward

$$B_1, B_2, \cdots, B_M = F_{u_1}(W), F_{w_2}(W), \cdots, F_{u,u}(W)$$

$$solve \min_a J(\alpha) = \|W - B\alpha\|^2 \ for \ \alpha$$

$$O = \sum_{m=1}^{M} \alpha_m \, \mathrm{Conv}\,(B_m, A)$$

▶ Backward

$$\frac{\partial c}{\partial W} = \frac{\partial c}{\partial O} \left( \sum_{m=1}^{M} \alpha_m \frac{\partial O}{\partial B_m} \frac{\partial B_m}{\partial W} \right) \overset{STE}{=} \frac{\partial c}{\partial O} \left( \sum_{m=1}^{M} \alpha_m \frac{\partial O}{\partial B_m} \right) = \sum_{m=1}^{M} \alpha_m \frac{\partial c}{\partial B_m}$$

# ABC-Net

## Multiple Binary Activations

▶ Bounded Activation Function

$$h(x) \in [0, 1]$$
$$h_r(x) = \text{clip}(x + v, 0, 1)$$
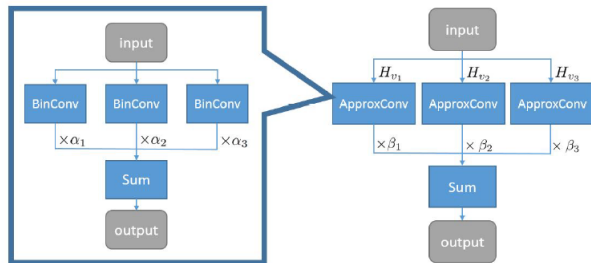
where $v$ is a shift parameter

▶ Binarization Function

$$H_v(\boldsymbol{R}) := 2\mathbb{I}_{h_v(\boldsymbol{R}) \geq 0.5} - 1$$
$$A_1, A_2, \ldots, A_N = H_{v_1}(R), H_{v_2}(R), \ldots, H_{v_N}(R)$$
$$R \approx \beta_1 A_1 + \beta_2 A_2 + \ldots + \beta_N A_N$$

where $R$ is the real-value activation

▶ $A_1, A_2, \ldots, A_N$ is the base to represent the real-valued activations

# ABC-Net



- ▶ ApproxConv is expected to approximate the conventional full-precision convolution with linear combination of binary convolutions
- ▶ The right part is the overall block structure of the convolution in ABC-Net. The input is binarized using different functions $H_v1, H_v2, H_v3$

$$\text{Conv}(\boldsymbol{W}, \boldsymbol{R}) \approx \text{Conv}\left(\sum_{m=1}^{M} \alpha_m \boldsymbol{B}_m, \sum_{n=1}^{N} \beta_n \boldsymbol{A}_n\right) =$$
$$\sum_{m=1}^{M} \sum_{n=1}^{N} \alpha_m \beta_n \text{Conv}(\boldsymbol{B}_m, \boldsymbol{A}_n)$$

Read the paper[2] if you want to learn the specific details of the algorithm

---

**Towards Accurate Binary Convolutional Neural Network**

Xiaofan Lin          Cong Zhao          Wei Pan*

DJI Innovations Inc, Shenzhen, China

{xiaofan.lin, cong.zhao, wei.pan}@dji.com

---

[2] Xiaofan Lin, Cong Zhao, and Wei Pan (2017). "Towards accurate binary convolutional neural network". In: *Advances in Neural Information Processing Systems*, pp. 345–353.

# Overview

# Motivation and Intuition

## Motivation

▶ Although STE is often adopted to estimate the gradients in BP, there exists obvious gradient mismatch between the gradient of the binarization function

▶ With the restriction of STE, the parameters outside the range of $[-1 : +1]$ will not be updated.

Bi-real net: Enhancing the performance of 1-bit CNNs with improved representational capability and advanced training algorithm
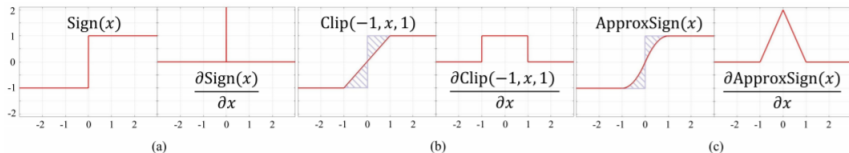
# Bi-Real

▶ Recall the partial derivative calculation in back propagation

$$\frac{\partial \mathcal{L}}{\partial \mathbf{A}_r^{l,t}} = \frac{\partial \mathcal{L}}{\partial \mathbf{A}_b^{l,t}} \frac{\partial \mathbf{A}_b^{l,t}}{\partial \mathbf{A}_r^{l,t}} = \frac{\partial \mathcal{L}}{\partial \mathbf{A}_b^{l,t}} \frac{\partial \operatorname{Sign}\left(\mathbf{A}_r^{l,t}\right)}{\partial \mathbf{A}_r^{l,t}} \approx \frac{\partial \mathcal{L}}{\partial \mathbf{A}_b^{l,t}} \frac{\partial F\left(\mathbf{A}_r^{l,t}\right)}{\partial \mathbf{A}_r^{l,t}}$$

▶ *Sign* function is a non-differentiable function, so $F$ is an approximation differentiable function of *Sign* function

# Bi-Real

$$\frac{\partial \mathcal{L}}{\partial \mathbf{A}_r^{l,t}} = \frac{\partial \mathcal{L}}{\partial \mathbf{A}_b^{l,t}} \frac{\partial \mathbf{A}_b^{l,t}}{\partial \mathbf{A}_r^{l,t}} = \frac{\partial \mathcal{L}}{\partial \mathbf{A}_b^{l,t}} \frac{\partial \operatorname{Sign}\left(\mathbf{A}_r^{l,t}\right)}{\partial \mathbf{A}_r^{l,t}} \approx \frac{\partial \mathcal{L}}{\partial \mathbf{A}_b^{l,t}} \frac{\partial F\left(\mathbf{A}_r^{l,t}\right)}{\partial \mathbf{A}_r^{l,t}}$$



## Approximation of *Sign* function

▶ Naive Approximation $F(x) = clip(x, 0, 1)$, see fig(b)

▶ More Precious Approximation in Bi-Real, see fig(c)

$$Approxsign(x) = \begin{cases} -1, & \text{if } x < -1 \\ 2x + x^2, & \text{if } -1 \le x < 0 \\ 2x - x^2, & \text{if } 0 \le x < 1 \\ 1, & \text{otherwise} \end{cases} \qquad \frac{\partial Approxsign(x)}{\partial x} = \begin{cases} 2 + 2x, & \text{if } -1 \le x < 0 \\ 2 - 2x, & \text{if } 0 \le x < 1 \\ 0, & \text{otherwise} \end{cases}$$

Read the paper[3] if you want to learn the specific details of the algorithm

## Bi-Real Net: Enhancing the Performance of 1-bit CNNs With Improved Representational Capability and Advanced Training Algorithm

Zechun Liu[1], Baoyuan Wu[2], Wenhan Luo[2], Xin Yang[3*], Wei Liu[2], and Kwang-Ting Cheng[1]
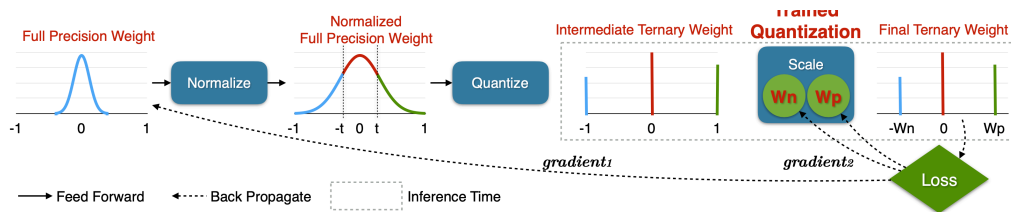
[1] Hong Kong University of Science and Technology
[2] Tencent AI lab
[3] Huazhong University of Science and Technology

---

[3] Zechun Liu et al. (2018). "Bi-real net: Enhancing the performance of 1-bit cnns with improved representational capability and advanced training algorithm". In: *Proceedings of the European conference on computer vision (ECCV)*, pp. 722–737.
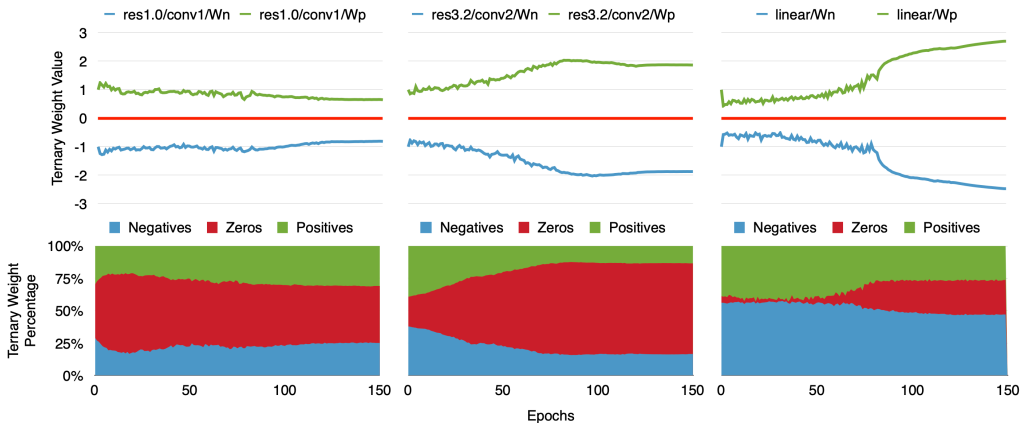
# Trained Ternary Quantization[4]



Overview of the trained ternary quantization procedure.

[4]Chenzhuo Zhu et al. (2017). "Trained ternary quantization". In: *Proc. ICLR*.

# Trained Ternary Quantization[4]



Ternary weights value (above) and distribution (below) with iterations for different layers of ResNet-20 on CIFAR-10.

[4]Chenzhuo Zhu et al. (2017). "Trained ternary quantization". In: *Proc. ICLR*.

# Reading List

- Hyeonuk Kim et al. (2017). "A Kernel Decomposition Architecture for Binary-weight Convolutional Neural Networks". In: *Proc. DAC*, 60:1–60:6

- Jungwook Choi et al. (2018). "Pact: Parameterized clipping activation for quantized neural networks". In: *arXiv preprint arXiv:1805.06085*

- Dongqing Zhang et al. (2018). "Lq-nets: Learned quantization for highly accurate and compact deep neural networks". In: *Proceedings of the European conference on computer vision (ECCV)*, pp. 365–382

- Aojun Zhou et al. (2017). "Incremental network quantization: Towards lossless cnns with low-precision weights". In: *arXiv preprint arXiv:1702.03044*

- Zhaowei Cai et al. (2017). "Deep learning with low precision by half-wave gaussian quantization". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5918–5926