



香港中文大學  
The Chinese University of Hong Kong

CENG5030

# Part 2-4: CNN Inaccurate Speedup-2 — Quantization

**Bei Yu**

(Latest update: March 25, 2019)

Spring 2019

## These slides contain/adapt materials developed by

- ▶ Suyog Gupta et al. (2015). “Deep learning with limited numerical precision”. In: *Proc. ICML*, pp. 1737–1746
- ▶ 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





I want to research  
on a topic with DEAP  
LEARNING in it?



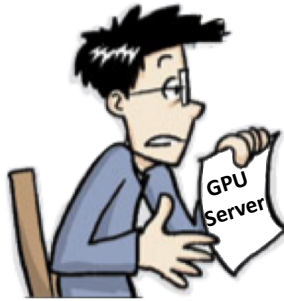
What should I learn  
to do well in  
computer vision  
research?





DEEP LEARNING



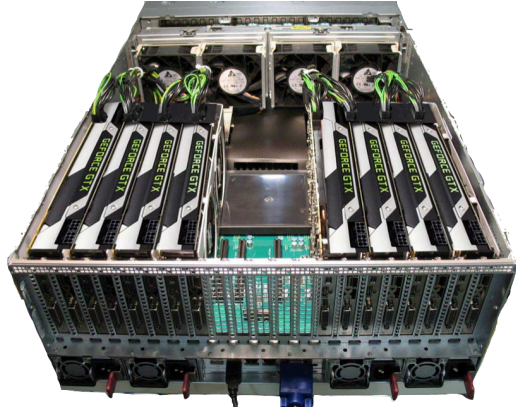


Ohhh No!!!



# State of the art recognition methods

- Very Expensive
  - Memory
  - Computation
  - Power



# Overview

Fixed-Point Representation

Binary/Ternary Network

Reading List



# Overview

Fixed-Point Representation

Binary/Ternary Network

Reading List



# Fixed-Point v.s. Floating-Point

# Fixed-Point v.s. Floating-Point

Fixed Point

6.5

10.75

$\begin{array}{|c|c|c|c|} \hline 8 & 4 & 2 & 1 \\ \hline 0 & 1 & 1 & 0 \\ \hline \end{array} \cdot \begin{array}{|c|c|c|c|} \hline 5 & 25 & 125 & 625 \\ \hline 1 & 0 & 0 & 0 \\ \hline \end{array}$   
6 . 5

$\begin{array}{|c|c|c|c|} \hline 8 & 4 & 2 & 1 \\ \hline 1 & 0 & 1 & 0 \\ \hline \end{array} \cdot \begin{array}{|c|c|c|c|} \hline 5 & 25 & 125 & 625 \\ \hline 1 & 1 & 0 & 0 \\ \hline \end{array}$

range →  $1010 \cdot 1100$  ← accuracy



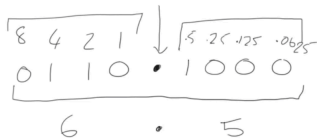


# Fixed-Point v.s. Floating-Point

Fixed Point

6.5

10.75

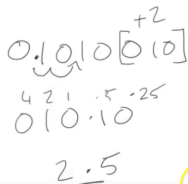
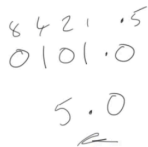


Floating Point

Mantissa 5 Exponent 3

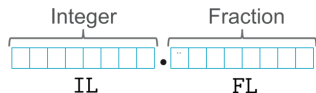


2's comp.



# Fixed-Point Arithmetic

*Number representation*  $\langle \text{IL}, \text{FL} \rangle$



Word Length  $\text{WL} = \text{IL} + \text{FL}$

Granularity  $2^{-\text{FL}}$

Range  $[-2^{\text{IL}-1}, 2^{\text{IL}-1} - 2^{-\text{FL}}]$

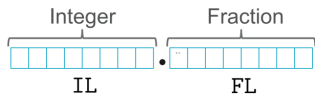
*Convert*  $(x, \langle \text{IL}, \text{FL} \rangle) =$

$$\begin{cases} -2^{\text{IL}-1} & \text{if } x \leq -2^{\text{IL}-1} \\ 2^{\text{IL}-1} - 2^{-\text{FL}} & \text{if } x \geq 2^{\text{IL}-1} - 2^{-\text{FL}} \\ \text{Round}(x, \langle \text{IL}, \text{FL} \rangle) & \text{otherwise} \end{cases}$$



# Fixed-Point Arithmetic

## Number representation $\langle \text{IL}, \text{FL} \rangle$



Word Length  $WL = IL + FL$

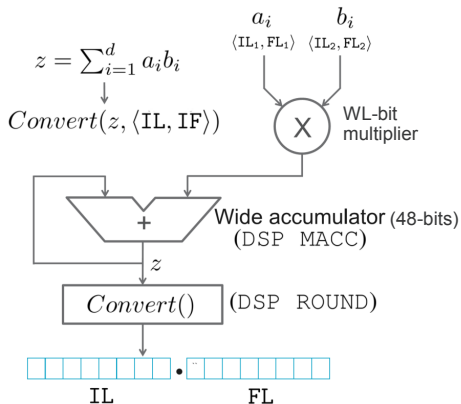
Granularity  $2^{-FL}$

Range  $[-2^{IL-1}, 2^{IL-1} - 2^{-FL}]$

$Convert(x, \langle \text{IL}, \text{FL} \rangle) =$

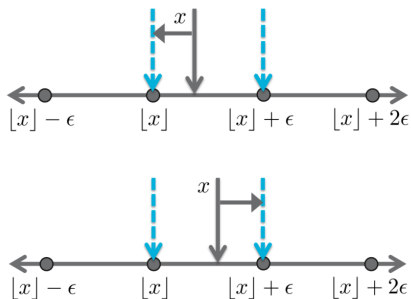
$$\begin{cases} -2^{IL-1} & \text{if } x \leq -2^{IL-1} \\ 2^{IL-1} - 2^{-FL} & \text{if } x \geq 2^{IL-1} - 2^{-FL} \\ Round(x, \langle \text{IL}, \text{FL} \rangle) & \text{otherwise} \end{cases}$$

## Multiply-and-ACCumulate



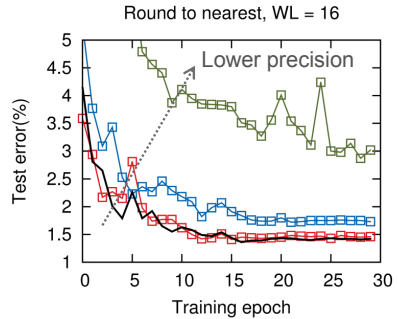
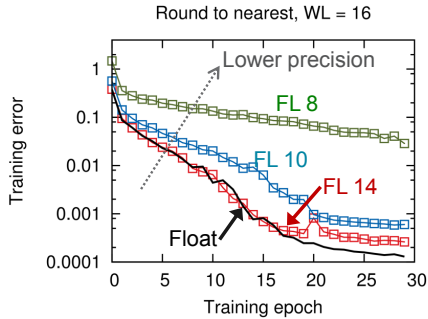
# Fixed-Point Arithmetic: Rounding Modes

*Round-to-nearest*

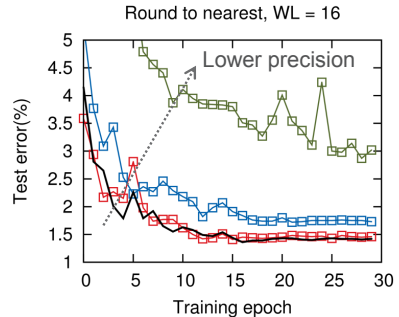
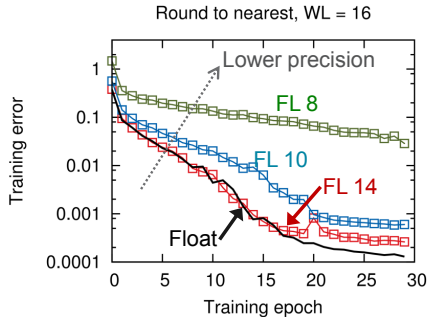




# MNIST: *Fully-connected DNNs*



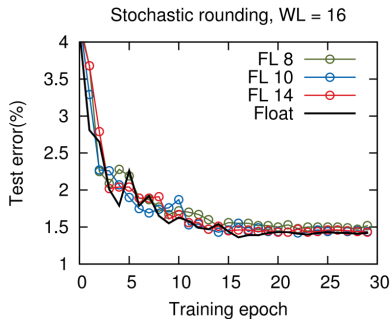
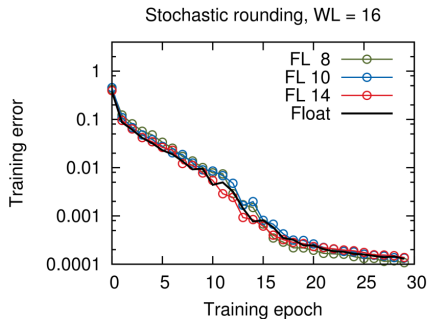
# MNIST: *Fully-connected DNNs*



- For small fractional lengths (FL < 12), a large majority of weight updates are rounded to zero when using the round-to-nearest scheme.
  - Convergence slows down
- For FL < 12, there is a noticeable degradation in the classification accuracy



# MNIST: *Fully-connected DNNs*

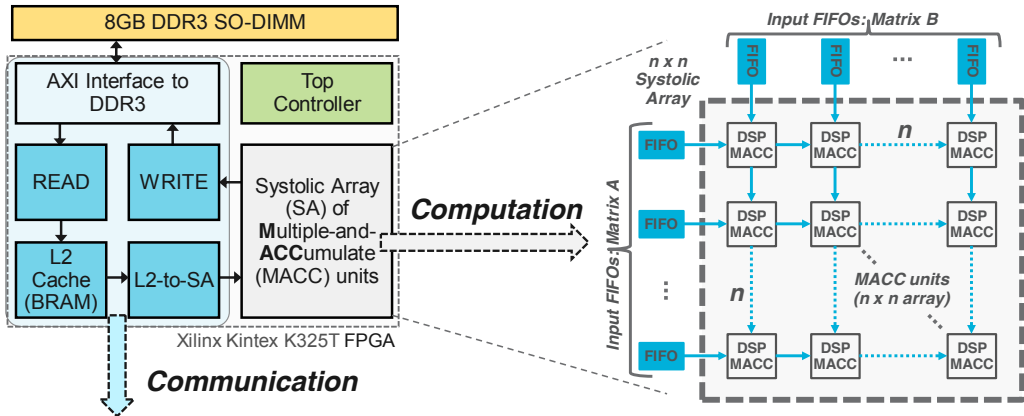


- Stochastic rounding preserves gradient information (statistically)
  - No degradation in convergence properties
- Test error nearly equal to that obtained using 32-bit floats





# FPGA prototyping: GEMM with stochastic rounding

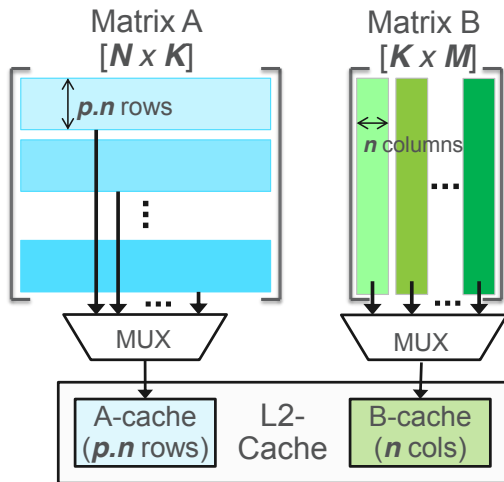


Top-level controller and memory hierarchy designed to maximize data reuse

Wavefront systolic array for computing matrix product  $AB$ . Arrows indicate dataflow



# Maximizing data reuse



## Inner Loop:

Cycle through columns of Matrix B  
( $M/n$  iterations)

## Outer Loop:

Cycle through rows of Matrix A  
( $K/p \cdot n$  iterations)

Re-use factor for Matrix A:  $M$  times

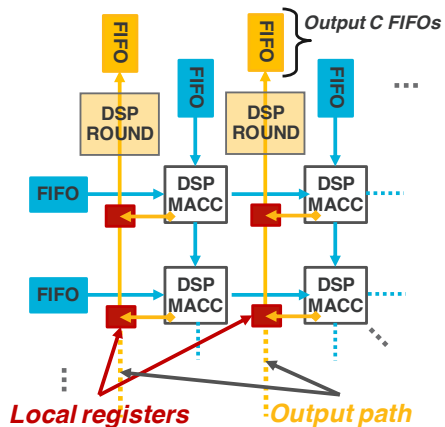
Re-use factor for Matrix B:  $p \cdot n$  times

$n$  : dimension of the systolic array

$p$  : parameter chosen based on available BRAM resources



# Stochastic rounding

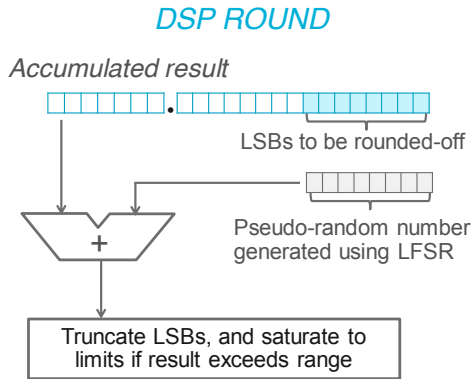
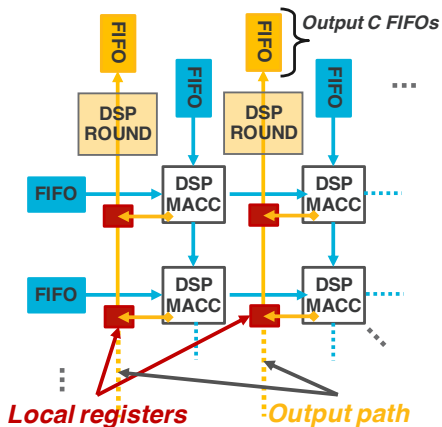


23

<sup>1</sup>Suyog Gupta et al. (2015). "Deep learning with limited numerical precision". In: *Proc. ICML*, pp. 1737-1746. > < ≡ ↺ ↻



# Stochastic rounding



These operations can be implemented efficiently using a single DSP unit



# Overview

Fixed-Point Representation

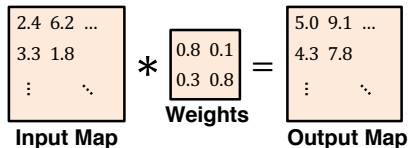
Binary/Ternary Network

Reading List



# Binarized Neural Networks (BNN)

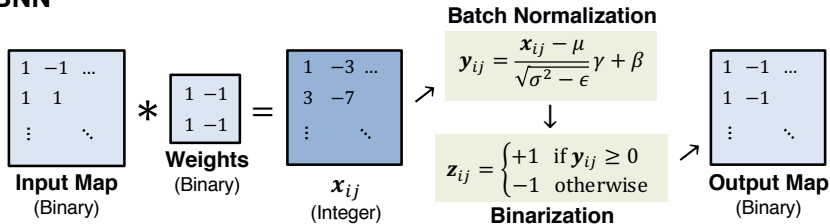
## CNN



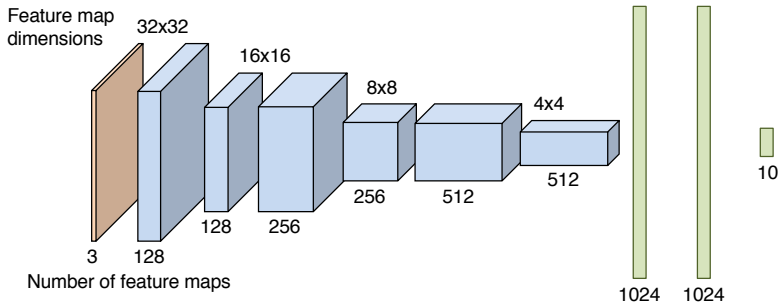
## Key Differences

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

## BNN



## BNN CIFAR-10 Architecture [2]



- ▶ 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)

# 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

$b_1$	$b_2$	$b_1 \text{ XOR } b_2$
0	0	0
0	1	1
1	0	1
1	1	0

- Encode  $\{+1, -1\}$  as  $\{0, 1\}$   $\rightarrow$  multiplies become XORs
- Conv/dense layers do dot products  $\rightarrow$  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





# 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

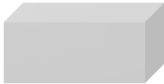

## ► Comparison:

- 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*.



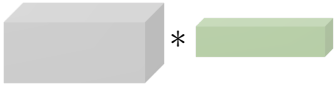
 * 	Operations	Memory	Computation
$\mathbb{R}$ * $\mathbb{R}$	+ - ×	1x	1x

Binary Weight Networks

XNOR-Networks

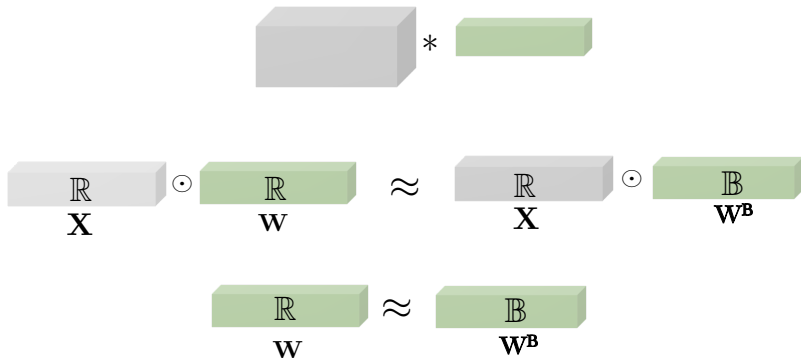
<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.



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

<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.





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



# Quantization Error

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

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

---

<sup>2</sup>[Mohammad Rastegari et al. \(2016\)](#). “XNOR-NET: Imagenet classification using binary convolutional neural networks”. In: *Proc. ECCV*, pp. 525–542.



# Optimal Scaling Factor

$$\begin{array}{c} \mathbb{R} \\ \mathbf{W} \end{array} \approx \alpha \begin{array}{c} \mathbb{B} \\ \mathbf{W}^{\mathbb{B}} \end{array}$$

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

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



## How to train a CNN with binary filters?

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

<sup>2</sup>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:*

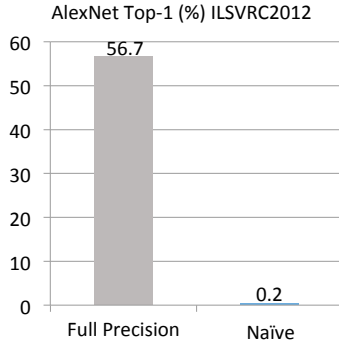
1. Train a network with real value parameters
2. Binarize the weight filters

---

<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.







---

<sup>2</sup>[Mohammad Rastegari et al. \(2016\)](#). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.



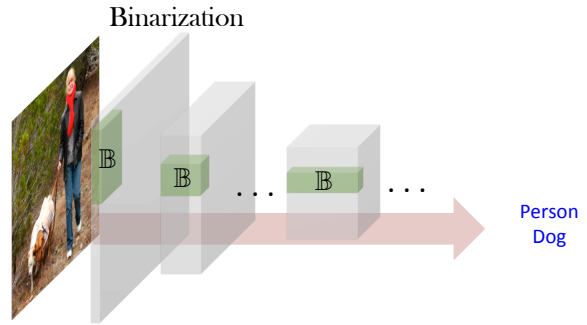


Binarization



---

<sup>2</sup>[Mohammad Rastegari et al. \(2016\)](#). “XNOR-NET: Imagenet classification using binary convolutional neural networks”. In: *Proc. ECCV*, pp. 525–542.



<sup>2</sup>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:*

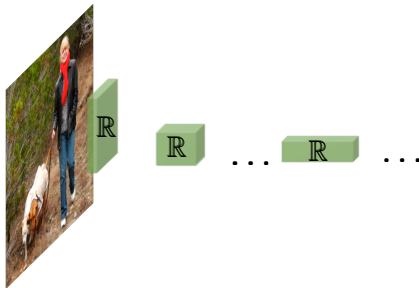
1. Randomly initialize  $\mathbf{W}$
2. For  $iter = 1$  to  $N$
3. Load a random input image  $\mathbf{X}$
4.  $\mathbf{W}^B = \text{sign}(\mathbf{W})$
5.  $\alpha = \frac{\|\mathbf{W}\|_{\ell_1}}{n}$
6. Forward pass with  $\alpha, \mathbf{W}^B$
7. Compute loss function  $\mathbf{C}$
8.  $\frac{\partial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



# Binary Weight Network

*Train for binary weights:*

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



# Binary Weight Network

*Train for binary weights:*

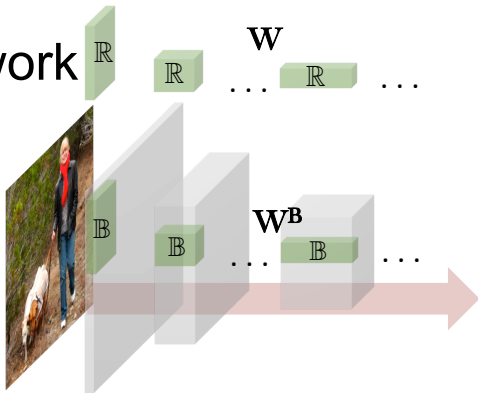
1. Randomly initialize  $W$
2. For  $iter = 1$  to  $N$
3. Load a random input image  $X$
4.  $W^B = \text{sign}(W)$
5.  $\alpha = \frac{\|W\|_{\ell_1}}{n}$
6. Forward pass with  $\alpha, W^B$
7. Compute loss function  $C$
8.  $\frac{\partial C}{\partial W} =$  Backward pass with  $\alpha, W^B$
9. Update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



# Binary Weight Network

*Train for binary weights:*

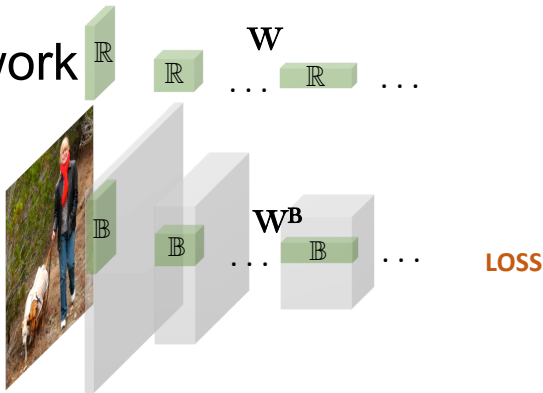
1. Randomly initialize  $W$
2. For  $iter = 1$  to  $N$
3. Load a random input image  $X$
4.  $W^B = \text{sign}(W)$
5.  $\alpha = \frac{\|W\|_{l1}}{n}$
6. Forward pass with  $\alpha, W^B$
7. Compute loss function  $C$
8.  $\frac{\partial C}{\partial W} =$  Backward pass with  $\alpha, W^B$
9. Update  $W$  ( $W = W - \frac{\partial C}{\partial W}$ )



# Binary Weight Network

*Train for binary weights:*

1. Randomly initialize  $W$
2. For  $iter = 1$  to  $N$
3. Load a random input image  $X$
4.  $W^B = \text{sign}(W)$
5.  $\alpha = \frac{\|W\|_{l1}}{n}$
6. Forward pass with  $\alpha, W^B$
7. Compute loss function  $C$
8.  $\frac{\partial C}{\partial W} =$  Backward pass with  $\alpha, W^B$
9. Update  $W$  ( $W = W - \frac{\partial C}{\partial W}$ )

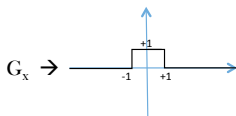
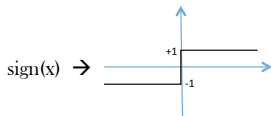
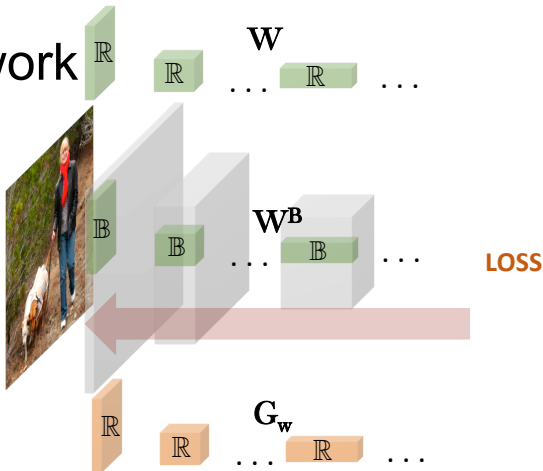




# Binary Weight Network

*Train for binary weights:*

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



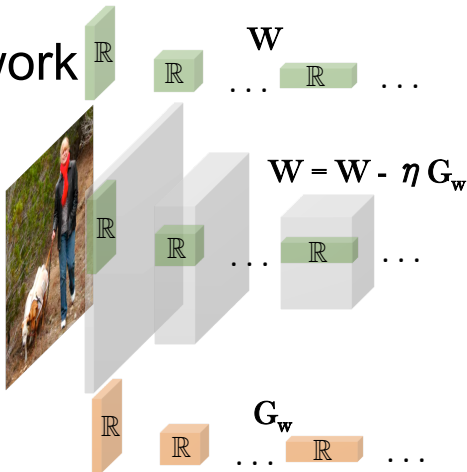
[Hinton et al. 2012]

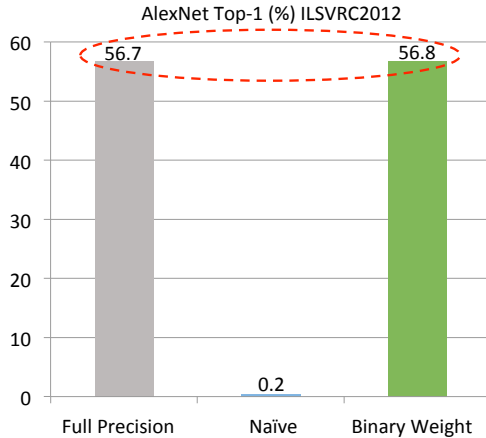


# Binary Weight Network

*Train for binary weights:*



1. Randomly initialize  $W$
2. For  $iter = 1$  to  $N$
3. Load a random input image  $X$
4.  $W^B = \text{sign}(W)$
5.  $\alpha = \frac{\|W\|_{l1}}{n}$
6. Forward pass with  $\alpha, W^B$
7. Compute loss function  $C$
8.  $\frac{\partial C}{\partial W} =$  Backward pass with  $\alpha, W^B$
9. Update  $W$  ( $W = W - \frac{\partial C}{\partial W}$ )





<sup>2</sup>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

<sup>2</sup>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)

$$\begin{array}{c} \mathbb{R} \\ \mathbf{X} \end{array} \odot \begin{array}{c} \mathbb{R} \\ \mathbf{W} \end{array} \approx \beta \begin{array}{c} \mathbb{B} \\ \mathbf{X}^{\mathbb{B}} \end{array} \odot \alpha \begin{array}{c} \mathbb{B} \\ \mathbf{W}^{\mathbb{B}} \end{array}$$

<sup>2</sup>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)

$$\underbrace{\begin{matrix} \mathbb{R} \\ \mathbf{X} \end{matrix}} \odot \underbrace{\begin{matrix} \mathbb{R} \\ \mathbf{W} \end{matrix}} \approx \underbrace{\beta \alpha}_{\gamma} \underbrace{\begin{matrix} \mathbb{B} \\ \mathbf{X}^{\mathbb{B}} \end{matrix}} \odot \underbrace{\begin{matrix} \mathbb{B} \\ \mathbf{W}^{\mathbb{B}} \end{matrix}}_{\mathbf{Y}^{\mathbb{B}}}$$

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

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

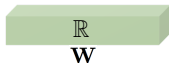
$$\mathbf{Y}^{\mathbb{B}*} = \text{sign}(\mathbf{Y}) \quad \gamma^* = \frac{1}{n} \|\mathbf{Y}\|_{\ell_1}$$

$$\mathbf{X}^{\mathbb{B}*} = \text{sign}(\mathbf{X}) \quad \mathbf{W}^{\mathbb{B}*} = \text{sign}(\mathbf{W})$$

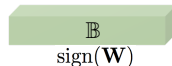
$$\alpha^* = \frac{1}{n} \|\mathbf{W}\|_{\ell_1} \quad \beta^* = \frac{1}{n} \|\mathbf{X}\|_{\ell_1}$$



(1) Binarizing Weights

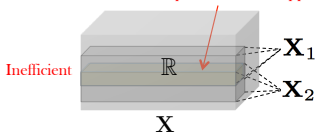


$$\frac{1}{n} \|\mathbf{W}\|_{\ell_1} = \alpha$$



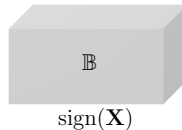
(2) Binarizing Input

Redundant computation in overlapping areas

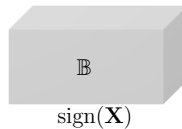
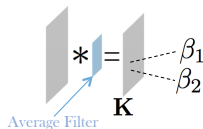
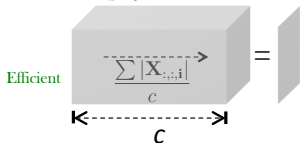


$$\frac{1}{n} \|\mathbf{X}_1\|_{\ell_1} = \beta_1$$

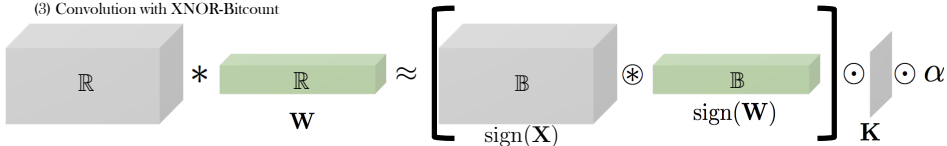
$$\frac{1}{n} \|\mathbf{X}_2\|_{\ell_1} = \beta_2$$

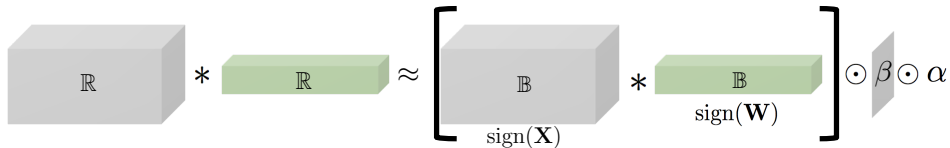


(2) Binarizing Input



(3) Convolution with XNOR-Bitcount



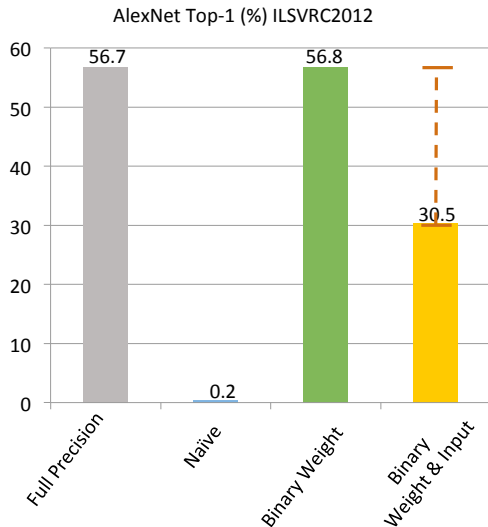


1. Randomly initialize  $\mathbf{W}$
2. For  $iter = 1$  to  $N$
3. Load a random input image  $\mathbf{X}$
4.  $\mathbf{W}^B = \text{sign}(\mathbf{W})$
5.  $\alpha = \frac{\|\mathbf{W}\|_{\ell_1}}{n}$
6. Forward pass with  $\alpha, \mathbf{W}^B$
7. Compute loss function  $\mathbf{C}$
8.  $\frac{\partial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

<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.



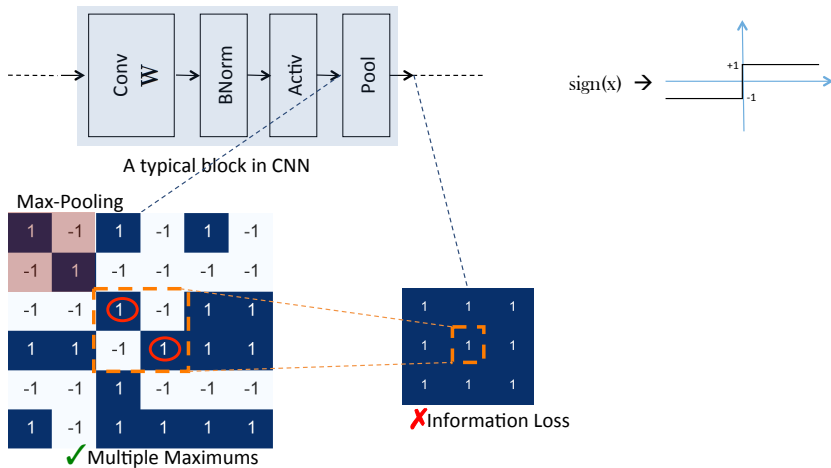




<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.



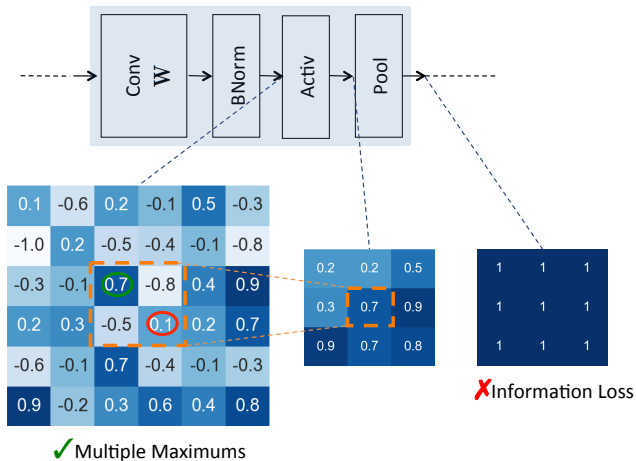
# Network Structure in XNOR-Networks



<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.



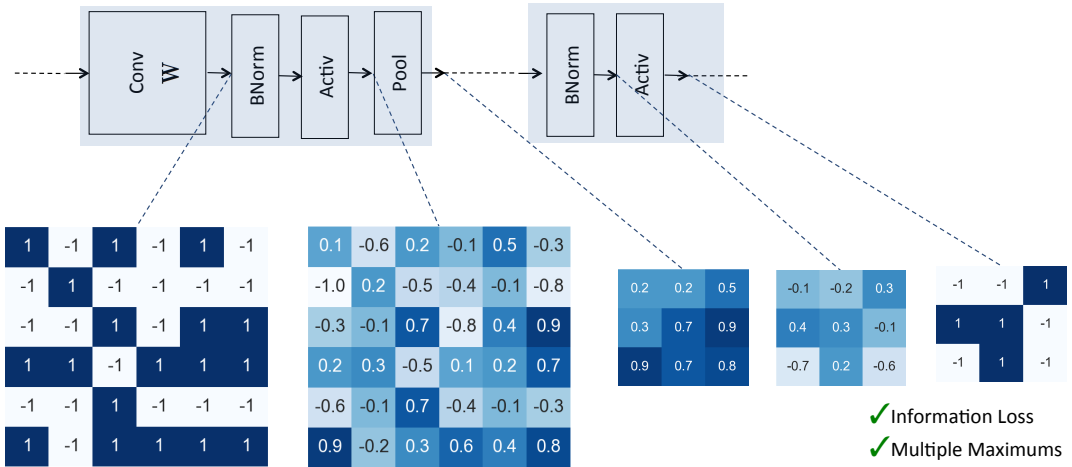
# Network Structure in XNOR-Networks



<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.

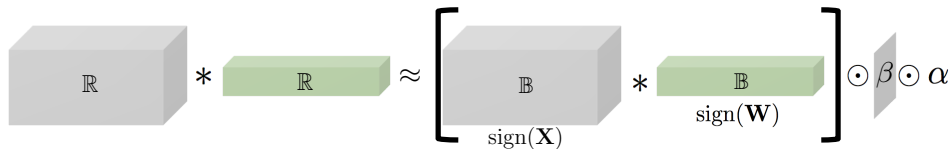


# Network Structure in XNOR-Networks

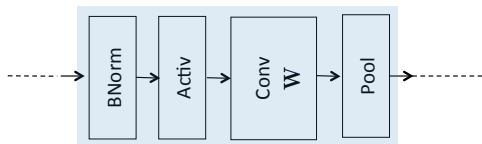


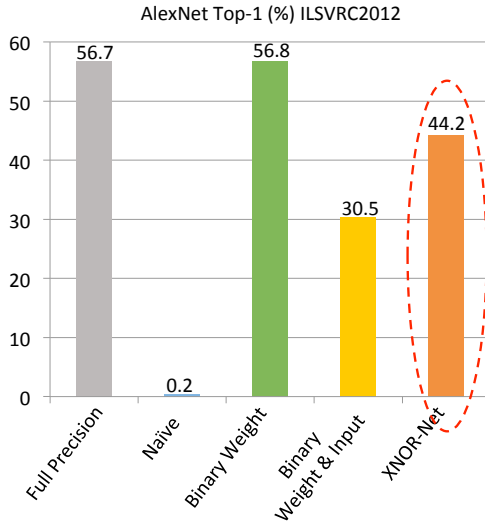
<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.



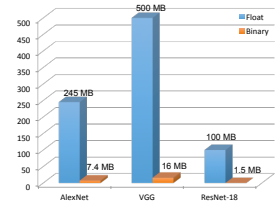


1. Randomly initialize  $\mathbf{W}$
2. For  $iter = 1$  to  $N$
3. Load a random input image  $\mathbf{X}$
4.  $\mathbf{W}^B = \text{sign}(\mathbf{W})$
5.  $\alpha = \frac{\|\mathbf{W}\|_{\ell_1}}{n}$
6. Forward pass with  $\alpha, \mathbf{W}^B$
7. Compute loss function  $\mathbf{C}$
8.  $\frac{\partial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

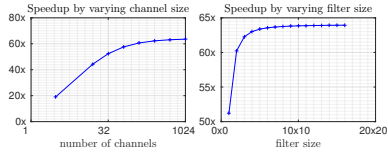




✓ 32x Smaller Model

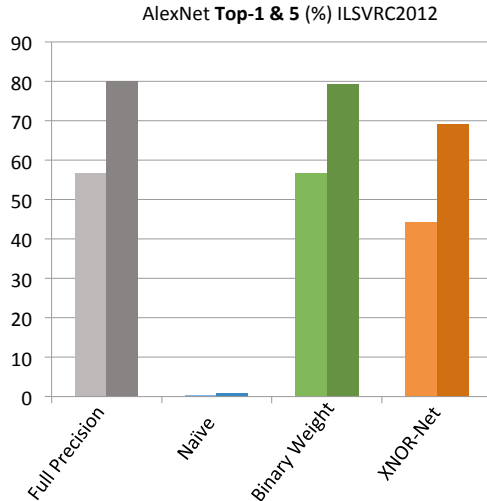


✓ 58x Less Computation



<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: Proc. ECCV, pp. 525–542.





<sup>2</sup>Mohammad Rastegari et al. (2016). "XNOR-NET: Imagenet classification using binary convolutional neural networks". In: *Proc. ECCV*, pp. 525–542.



# Overview

Fixed-Point Representation

Binary/Ternary Network

Reading List





# Further Reading List

## Fixed-Point Representation:

- ▶ Darryl Lin, Sachin Talathi, and Sreekanth Annapureddy (2016). “Fixed point quantization of deep convolutional networks”. In: *Proc. ICML*, pp. 2849–2858
- ▶ Soroosh Khoram and Jing Li (2018). “Adaptive quantization of neural networks”. In: *Proc. ICLR*

## Binary/Ternary Network:

- ▶ Hyeonuk Kim et al. (2017). “A Kernel Decomposition Architecture for Binary-weight Convolutional Neural Networks”. In: *Proc. DAC*, 60:1–60:6
- ▶ Chenzhuo Zhu et al. (2017). “Trained ternary quantization”. In: *Proc. ICLR*

