

XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks

Mohammad Rastegari

Vicente Ordonez

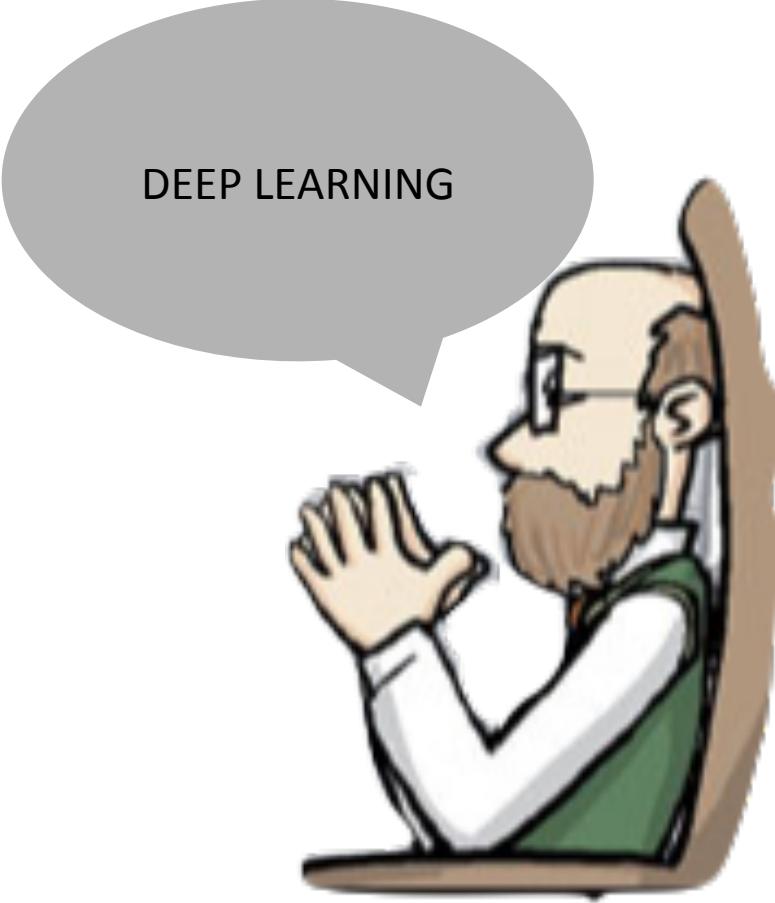
Joseph Redmon

Ali Farhadi

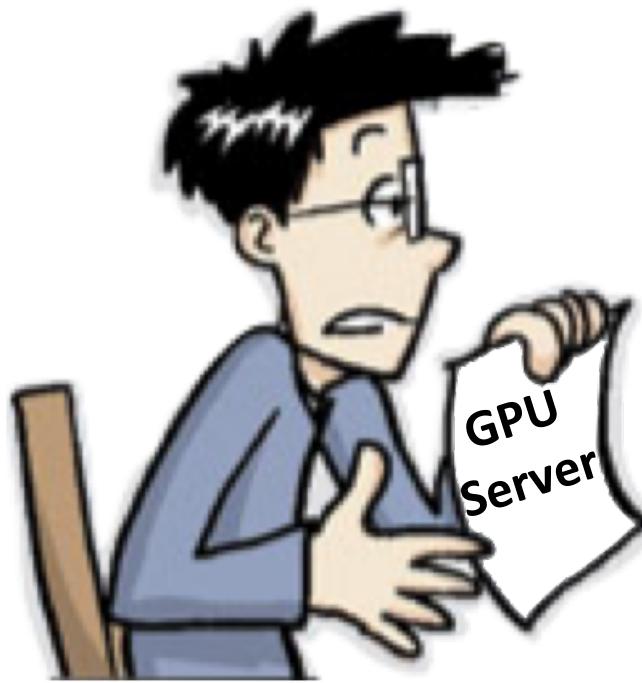








DEEP LEARNING



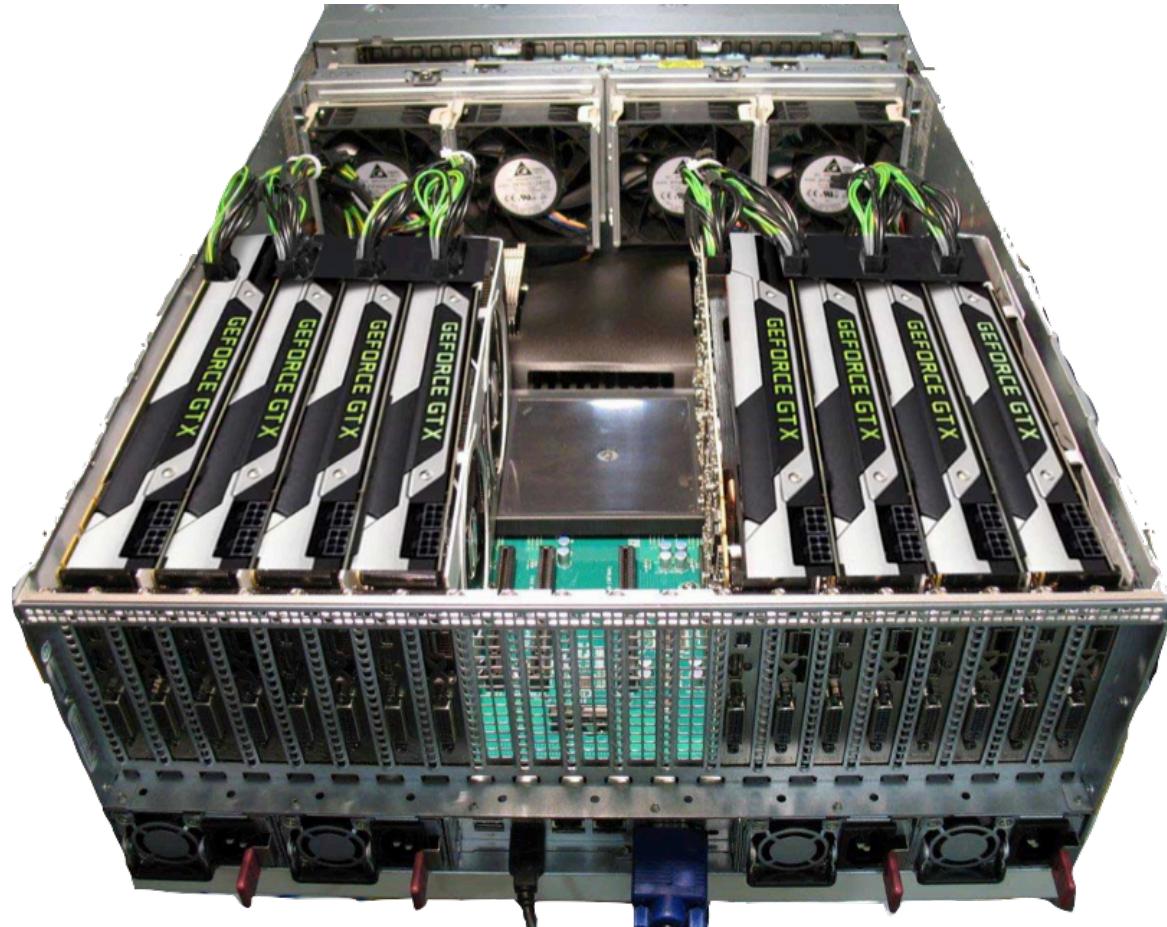


Ohhh No!!!



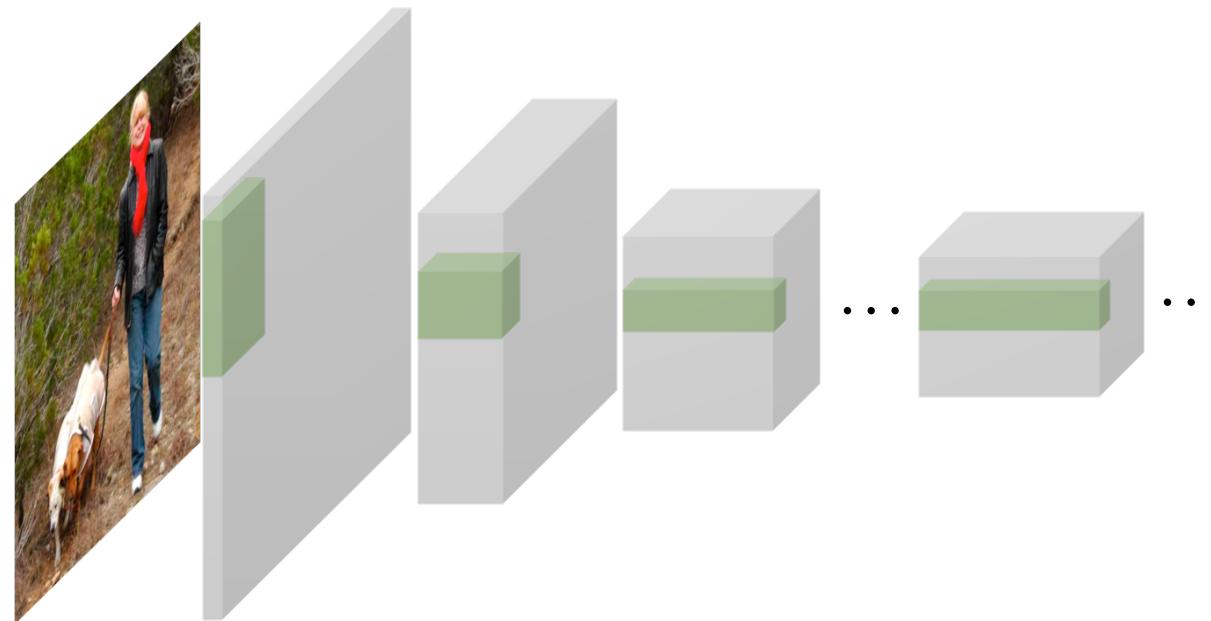
State of the art recognition methods

- Very Expensive
 - Memory
 - Computation
 - Power

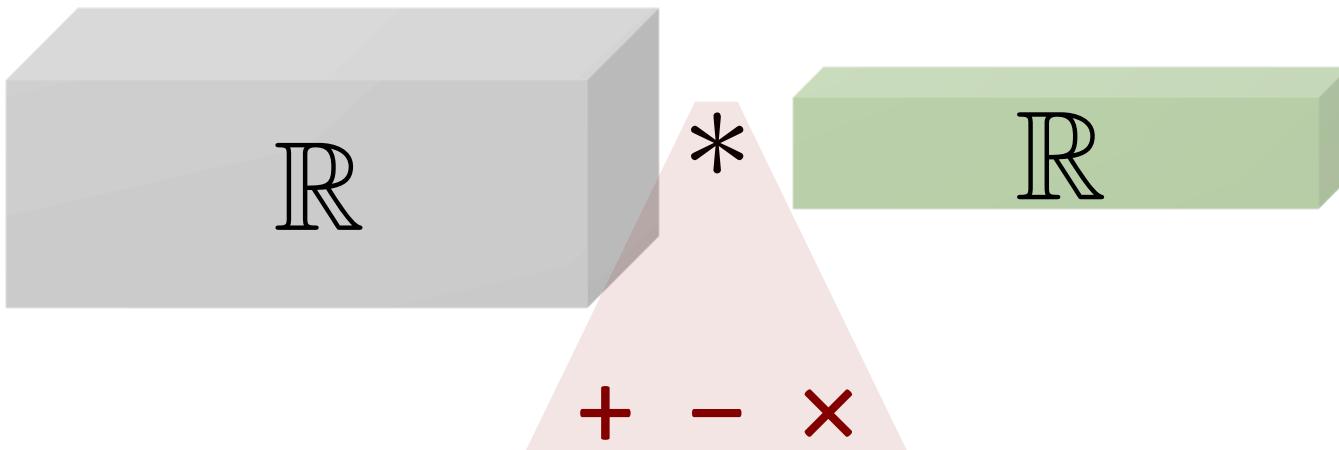




Convolutional Neural Networks



GPU !



Number of Operations :

- AlexNet → 1.5B FLOPs
- VGG → 19.6B FLOPs

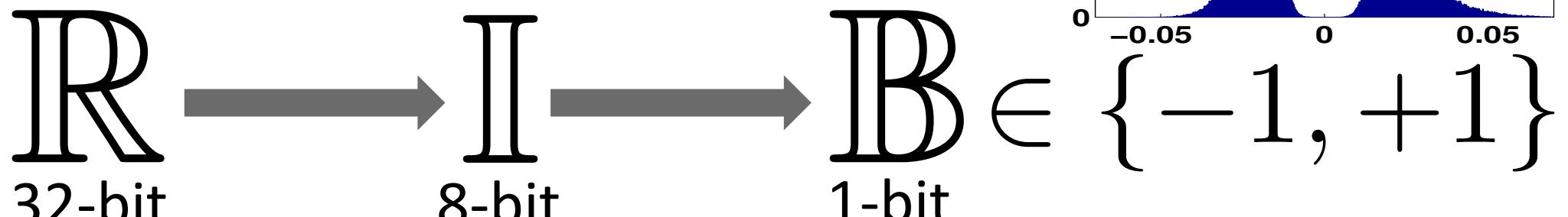
Inference time on CPU :

- AlexNet → ~3 fps
- VGG → ~0.25 fps

Lower Precision

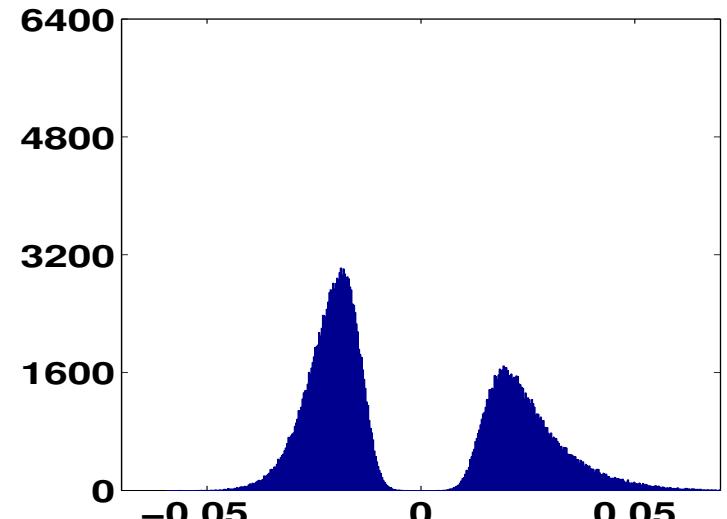
Reducing Precision

- Saving Memory
- Saving Computation



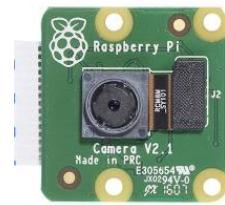
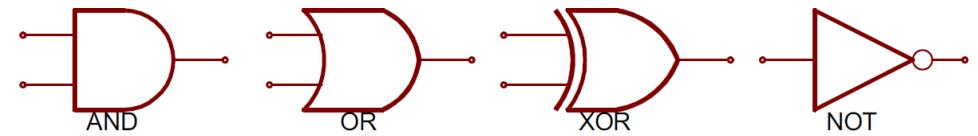
$\{-1,+1\}$	$\{0,1\}$
MUL	XNOR
ADD, SUB	Bit-Count (popcount)

[Han et al. 2016]



Why Binary?

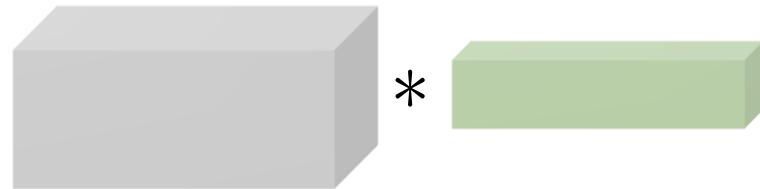
- Binary Instructions
 - AND, OR, XOR, XNOR, PoPCount (Bit-Count)
- Low Power Device



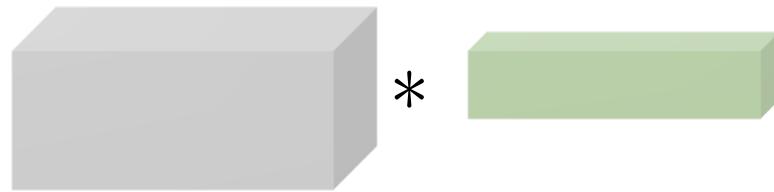
	*		Operations	Memory	Computation
\mathbb{R}	*	\mathbb{R}	+ - \times	1x	1x

Binary Weight Networks

XNOR-Networks



		Operations	Memory	Computation
R	* R	+ - ×	1x	1x
R	* B	+ -	~32x	~2x
B	* B	XNOR Bit-count	~32x	~58x



$$\begin{matrix} \text{R} \\ \text{X} \end{matrix} \odot \begin{matrix} \text{R} \\ \text{W} \end{matrix} \approx \begin{matrix} \text{R} \\ \text{X} \end{matrix} \odot \begin{matrix} \mathbb{B} \\ \text{W}^{\mathbb{B}} \end{matrix}$$

$$\begin{matrix} \text{R} \\ \text{W} \end{matrix} \approx \begin{matrix} \mathbb{B} \\ \text{W}^{\mathbb{B}} \end{matrix}$$

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

Quantization Error

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

$$\left\| \begin{matrix} W \\ R \end{matrix} - \begin{matrix} W^B \\ B \end{matrix} \right\| \approx 0.75$$

Optimal Scaling Factor

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

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

$$\boxed{\begin{aligned} \mathbf{W}^{\mathbf{B}^*} &= \text{sign}(\mathbf{W}) \\ \alpha^* &= \frac{1}{n} \|\mathbf{W}\|_{\ell 1} \end{aligned}}$$

How to train a CNN with binary filters?

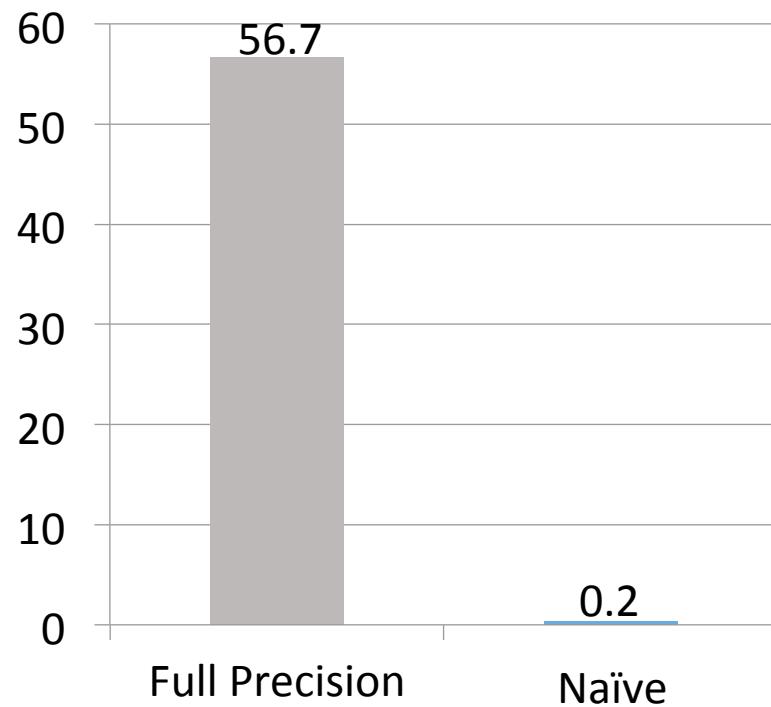
$$\mathbb{R} \star \mathbb{R} \approx (\mathbb{R} \star \mathbb{B})^\alpha$$

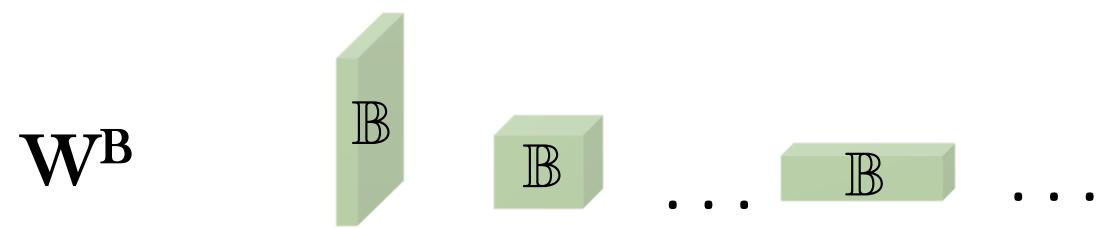
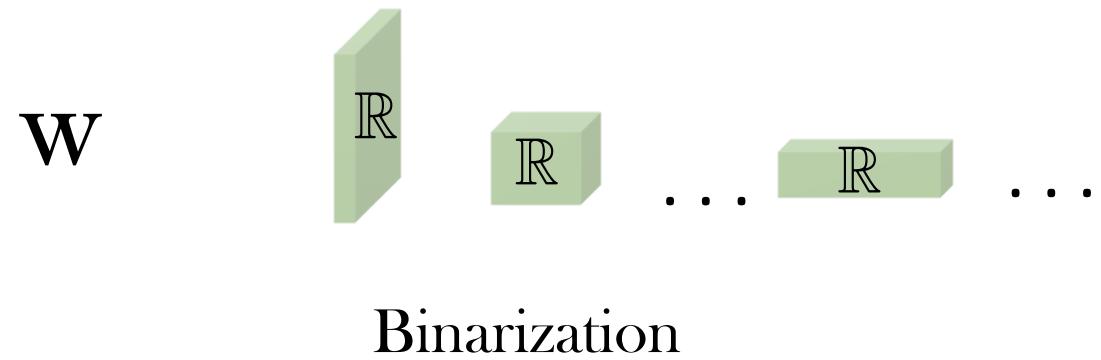
Training Binary Weight Networks

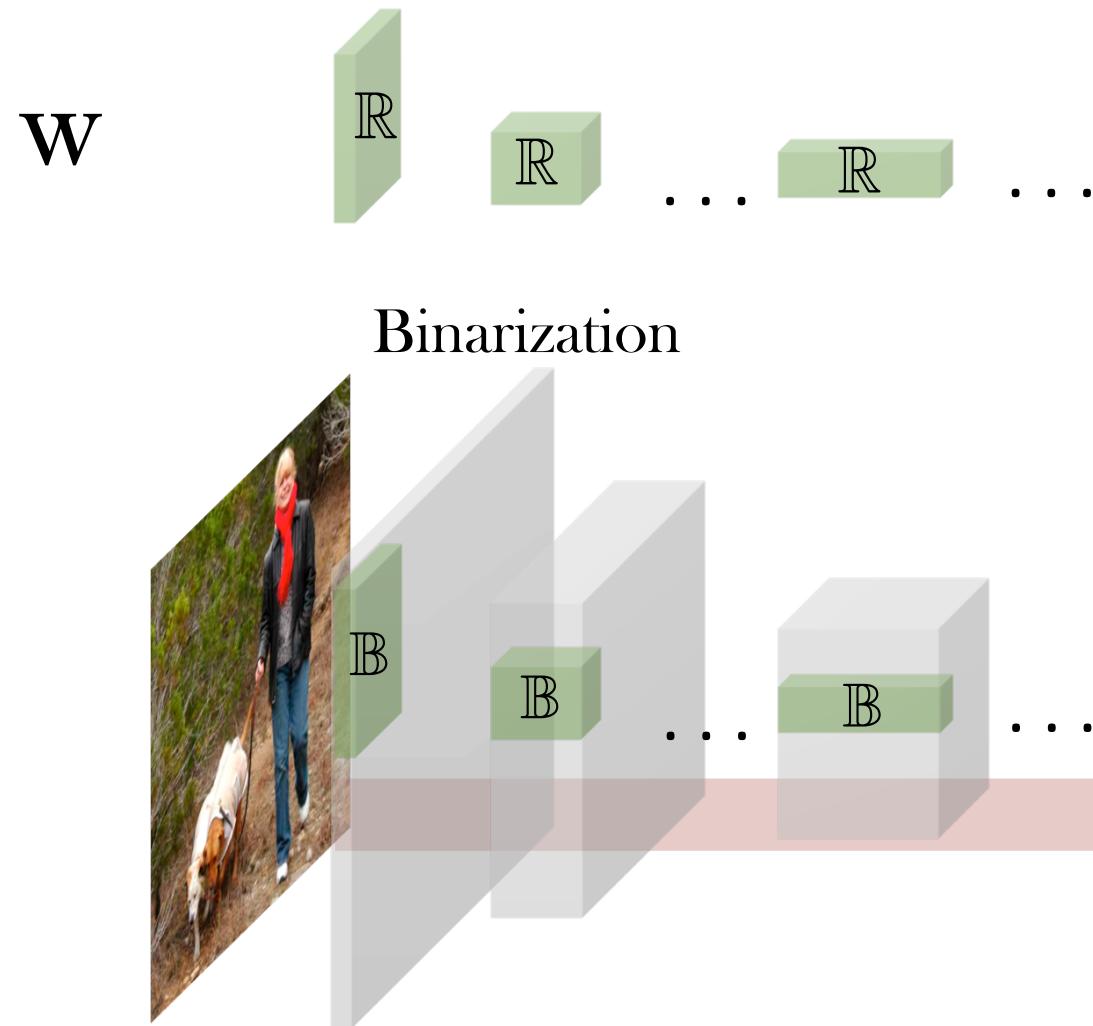
Naive Solution:

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

AlexNet Top-1 (%) ILSVRC2012



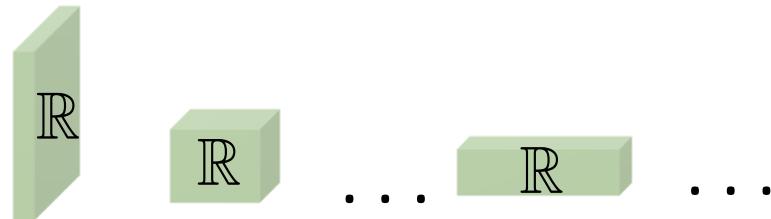




Binary Weight Network

Train for binary weights:

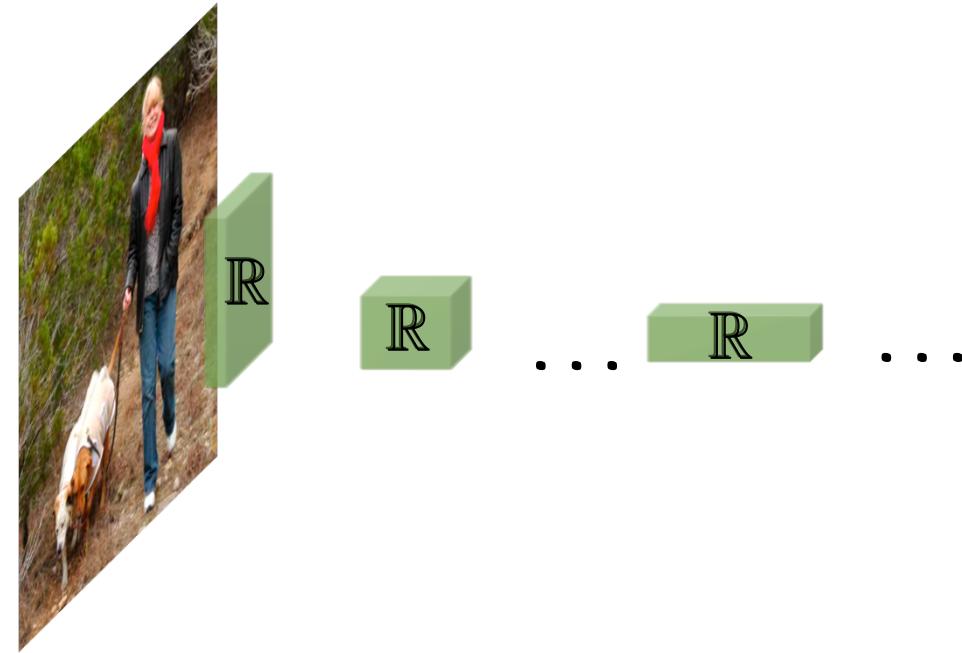
1. Randomly initialize \mathbf{W}
2. For $iter = 1$ to N
3. Load a random input image \mathbf{X}
4. $\mathbf{W}^B = \text{sign}(\mathbf{W})$
5. $\alpha = \frac{\|\mathbf{W}\|_{\ell_1}}{n}$
6. Forward pass with α, \mathbf{W}^B
7. Compute loss function \mathbf{C}
8. $\frac{\partial \mathbf{C}}{\partial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



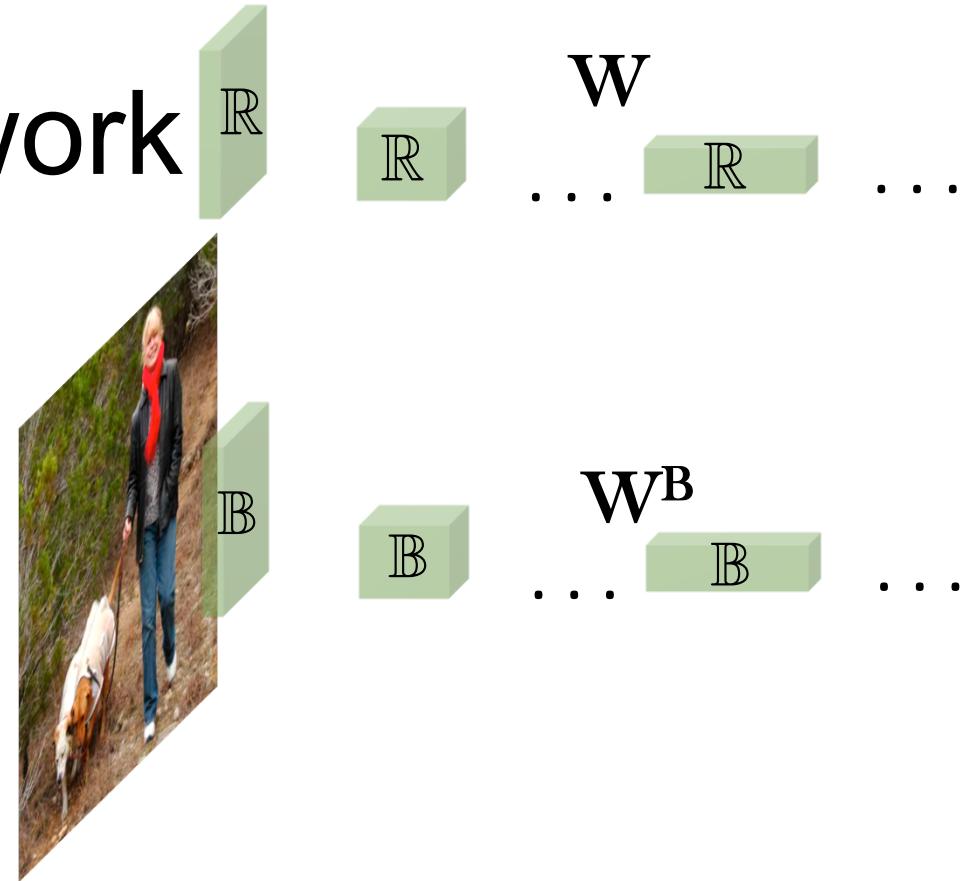
Binary Weight Network

Train for binary weights:

1. Randomly initialize \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 α, \mathbf{W}^B
7. Compute loss function \mathbf{C}
8. $\frac{\partial \mathbf{C}}{\partial \mathbf{W}} =$ Backward pass with α, \mathbf{W}^B
9. Update \mathbf{W} ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)



Binary Weight Network



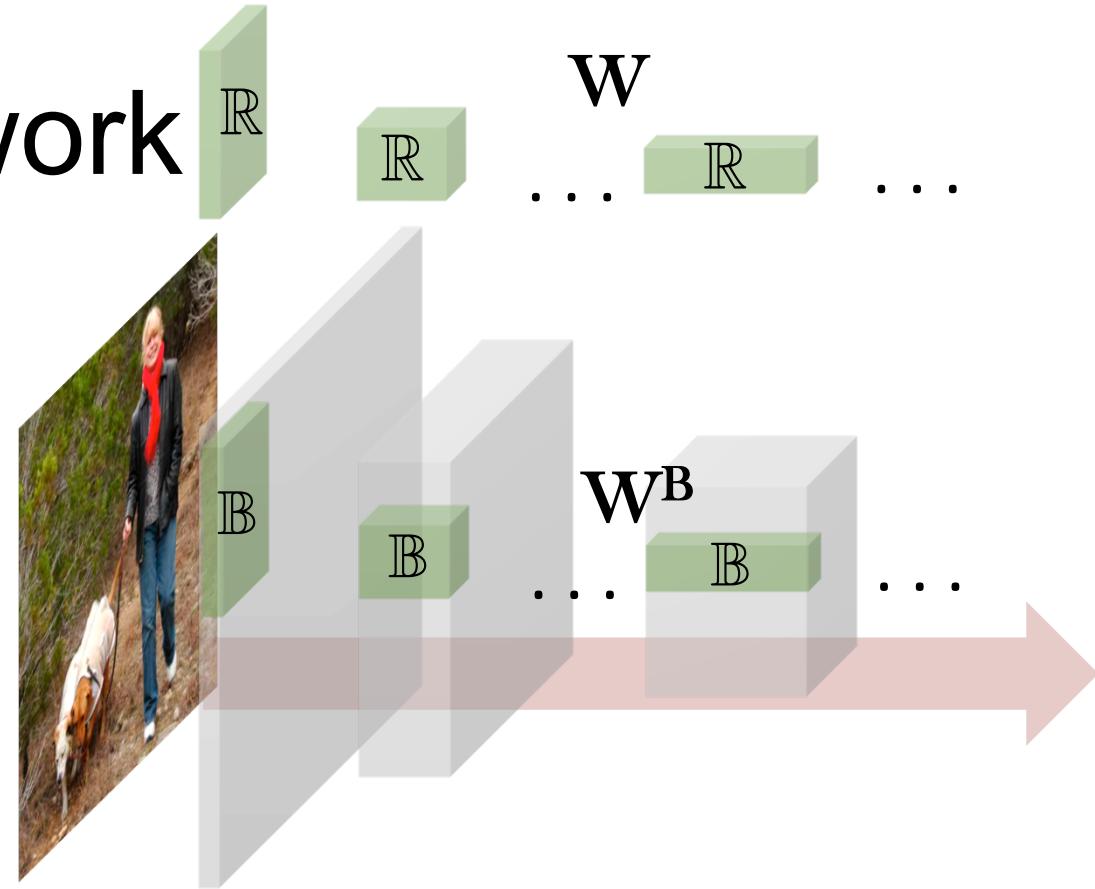
Train for binary weights:

1. Randomly initialize \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 α , \mathbf{W}^B
7. Compute loss function C
8. $\frac{\partial C}{\partial \mathbf{W}}$ = Backward pass with α , \mathbf{W}^B
9. Update \mathbf{W} ($\mathbf{W} = \mathbf{W} - \frac{\partial C}{\partial \mathbf{W}}$)

Binary Weight Network

Train for binary weights:

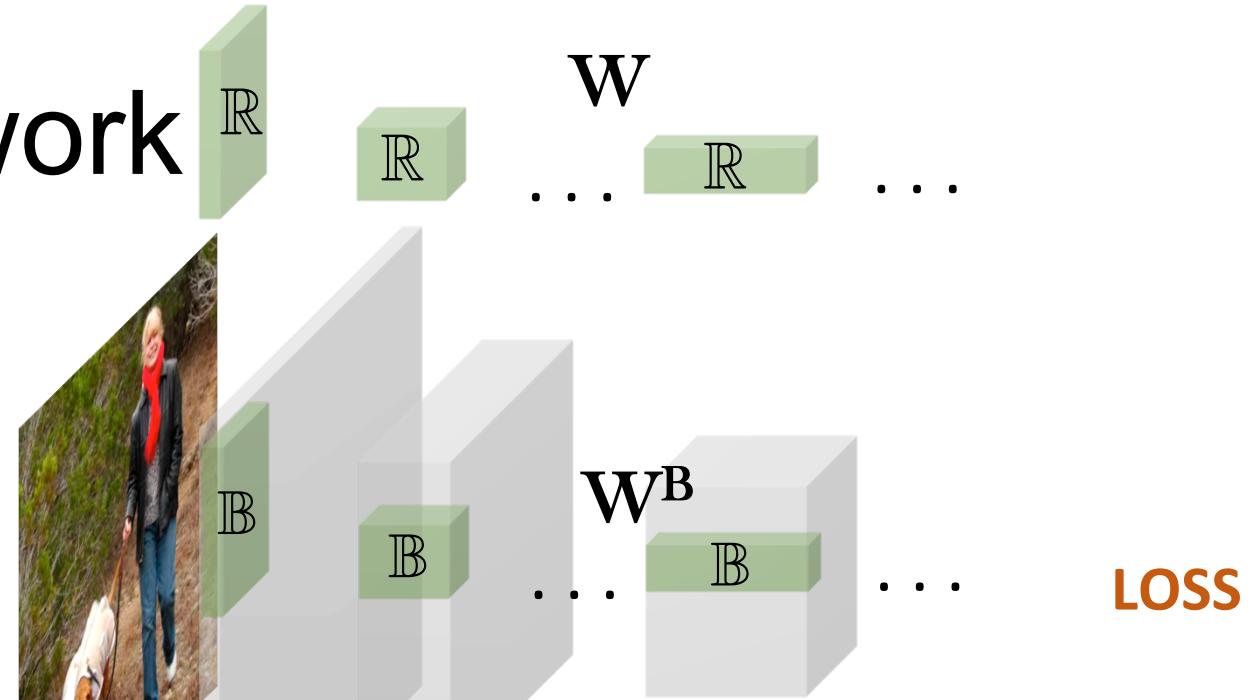
1. Randomly initialize \mathbf{W}
2. For $iter = 1$ to N
3. Load a random input image \mathbf{X}
4. $\mathbf{W}^B = \text{sign}(\mathbf{W})$
5. $\alpha = \frac{\|\mathbf{W}\|_{\ell_1}}{n}$
6. **Forward pass with α, \mathbf{W}^B**
7. Compute loss function \mathbf{C}
8. $\frac{\partial \mathbf{C}}{\partial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



Binary Weight Network

Train for binary weights:

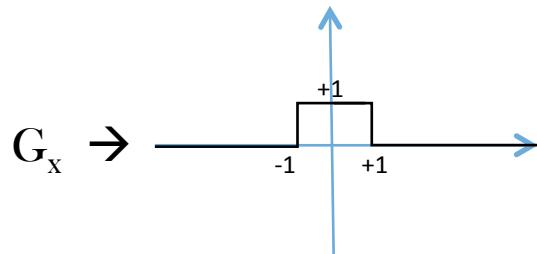
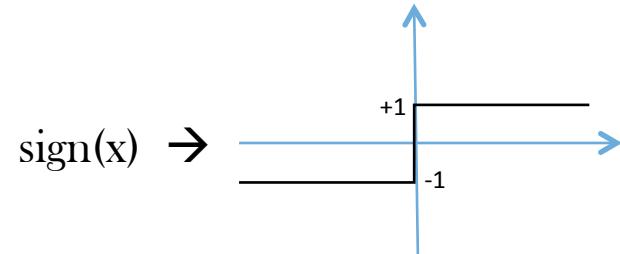
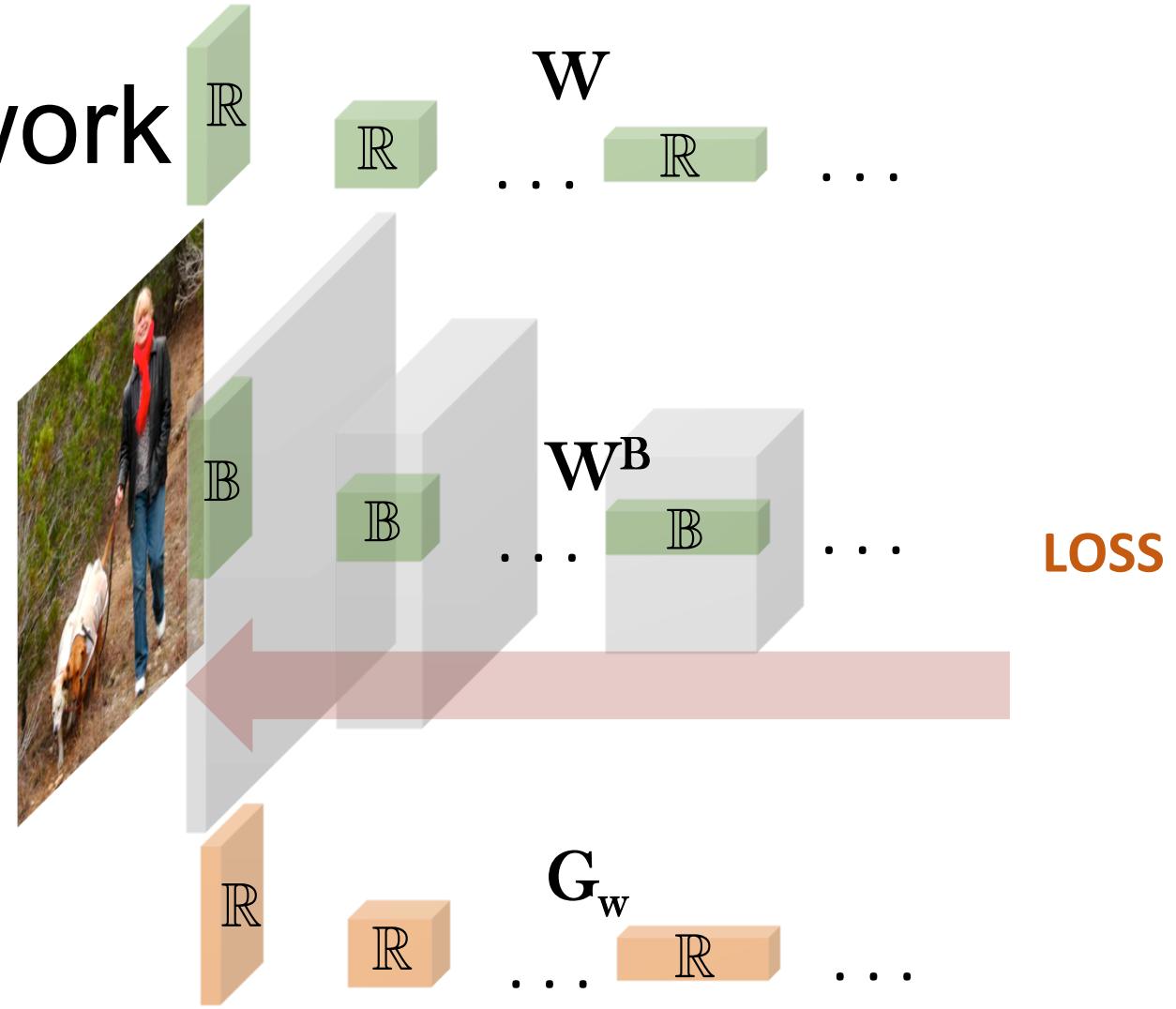
1. Randomly initialize \mathbf{W}
2. For $iter = 1$ to N
3. Load a random input image \mathbf{X}
4. $\mathbf{W}^B = \text{sign}(\mathbf{W})$
5. $\alpha = \frac{\|\mathbf{W}\|_{\ell_1}}{n}$
6. Forward pass with α, \mathbf{W}^B
7. Compute loss function \mathbf{C}
8. $\frac{\partial \mathbf{C}}{\partial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



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 α, \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} ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)

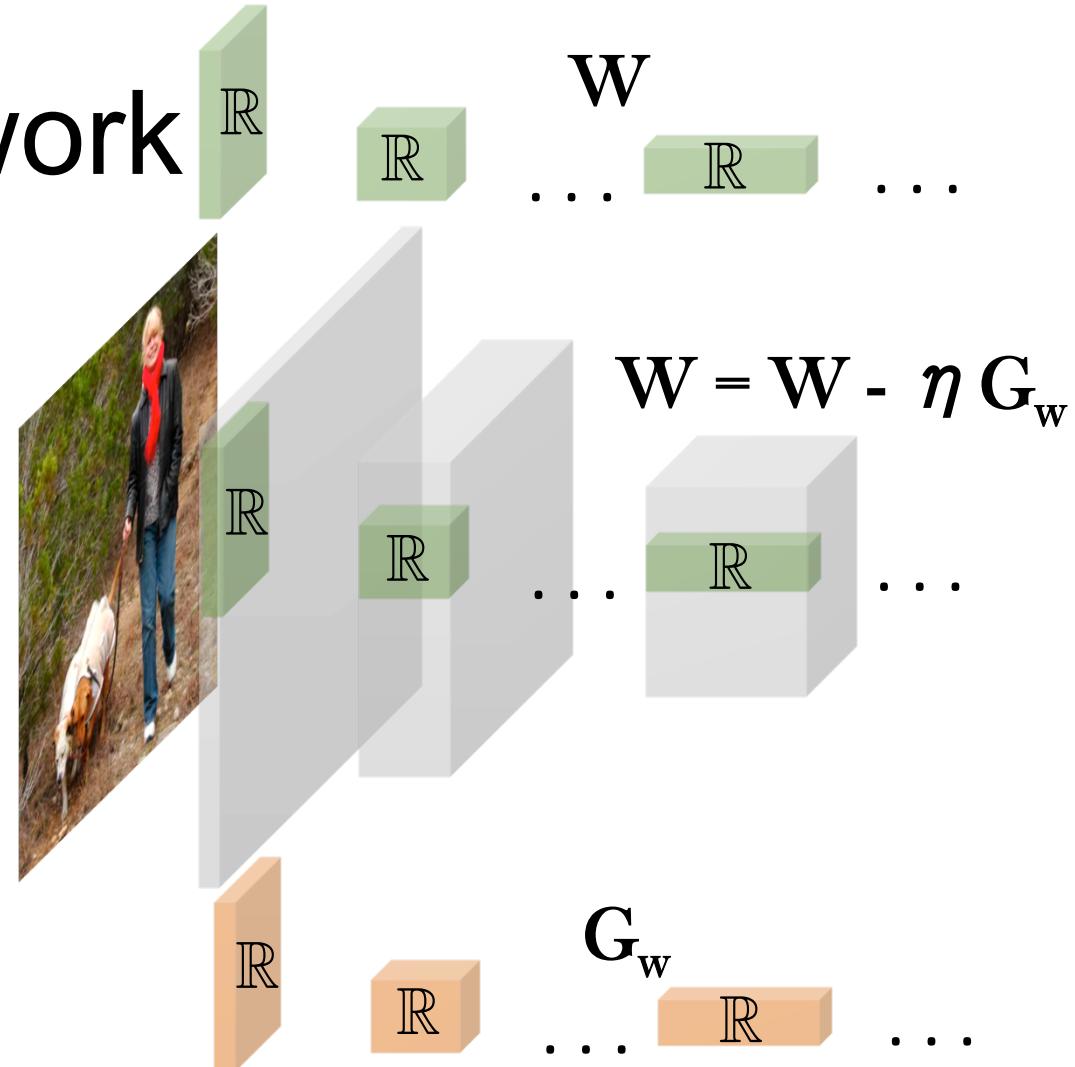


[Hinton et al. 2012]

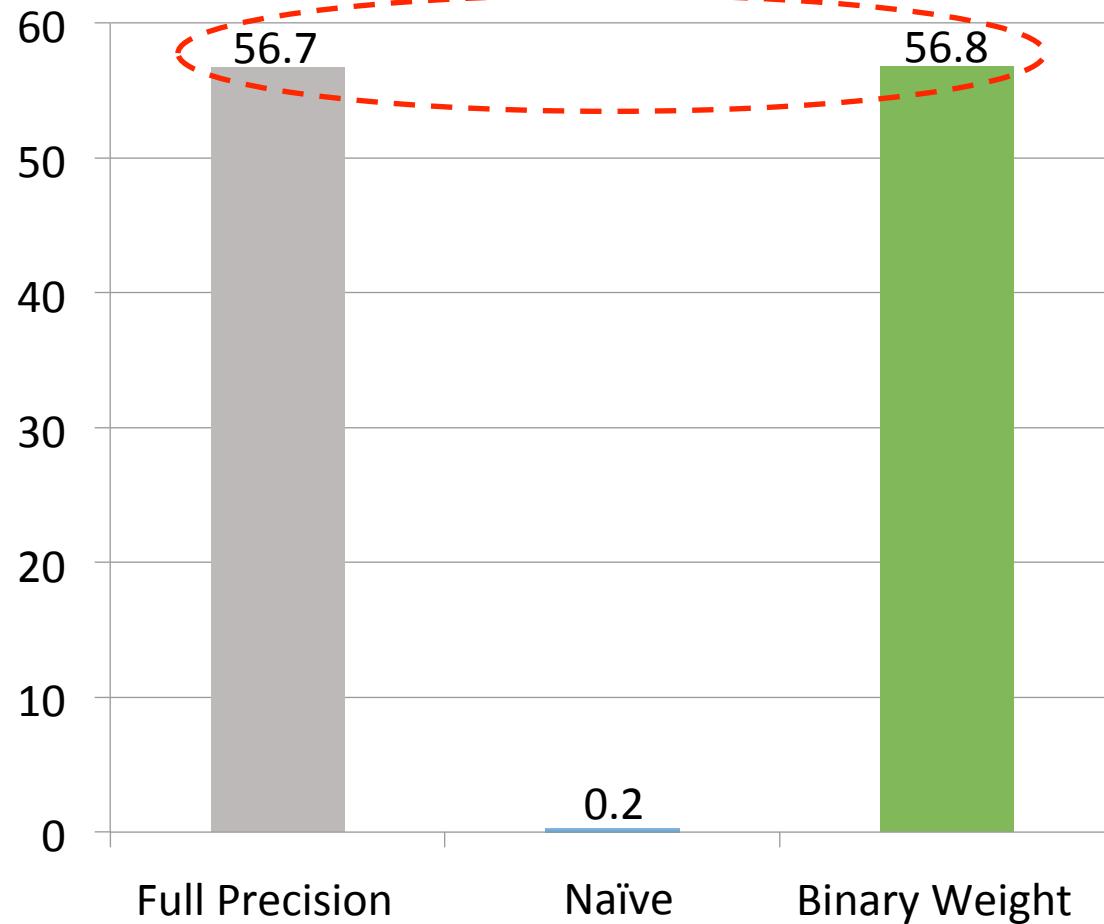
Binary Weight Network

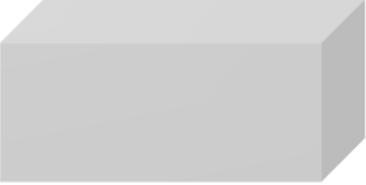
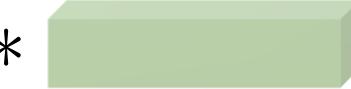
Train for binary weights:

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 α, \mathbf{W}^B
7. Compute loss function \mathbf{C}
8. $\frac{\partial \mathbf{C}}{\partial \mathbf{W}}$ = Backward pass with α, \mathbf{W}^B
9. Update \mathbf{W} ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)



AlexNet Top-1 (%) ILSVRC2012

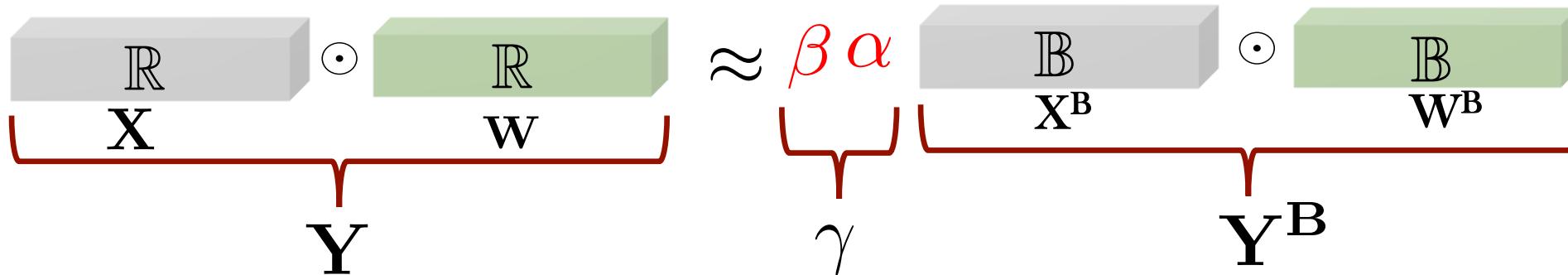


 * 	Operations	Memory	Computation
$R \quad * \quad R$	+ - \times	1x	1x
$R \quad * \quad B$	+ -	~32x	~2x
$B \quad * \quad B$ XNOR-Networks	XNOR Bit-count	~32x	~58x

Binary Input and Binary Weight (XNOR-Net)

$$\begin{matrix} \text{R} \\ \text{X} \end{matrix} \odot \begin{matrix} \text{R} \\ \text{W} \end{matrix} \approx \beta \quad \begin{matrix} \text{B} \\ \text{X}^B \end{matrix} \odot \alpha \begin{matrix} \text{B} \\ \text{W}^B \end{matrix}$$

Binary Input and Binary Weight (XNOR-Net)



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

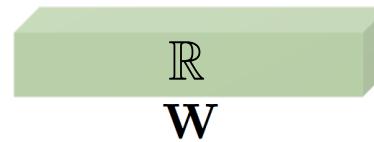
$$\mathbf{Y}^{B*}, \gamma^* = \arg \min_{\mathbf{Y}^B, \gamma} \|\mathbf{Y} - \gamma \mathbf{Y}^B\|^2$$

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

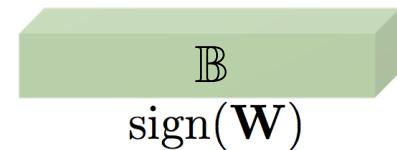
$$\boxed{\mathbf{X}^{B*} = \text{sign}(\mathbf{X}) \quad \mathbf{W}^{B*} = \text{sign}(\mathbf{W})}$$

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

(1) Binarizing Weights

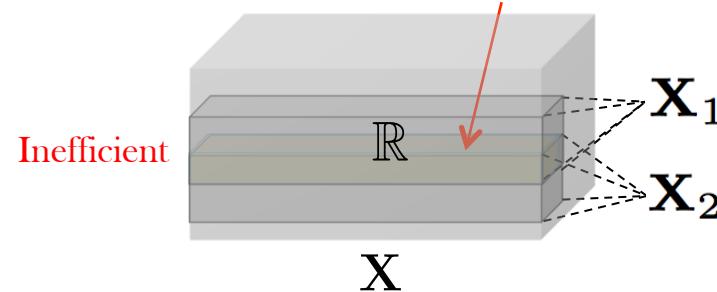


$$\frac{1}{n} \|\mathbf{W}\|_{\ell_1} = \alpha$$

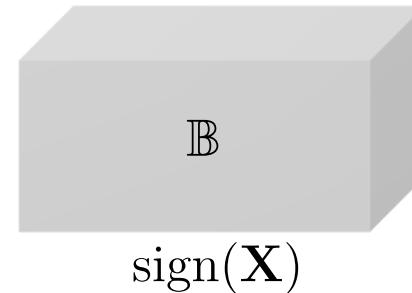


(2) Binarizing Input

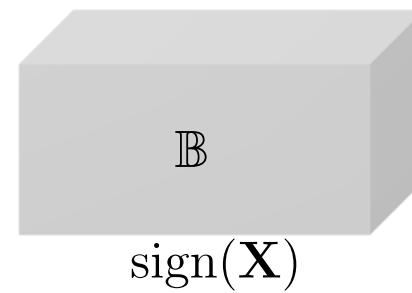
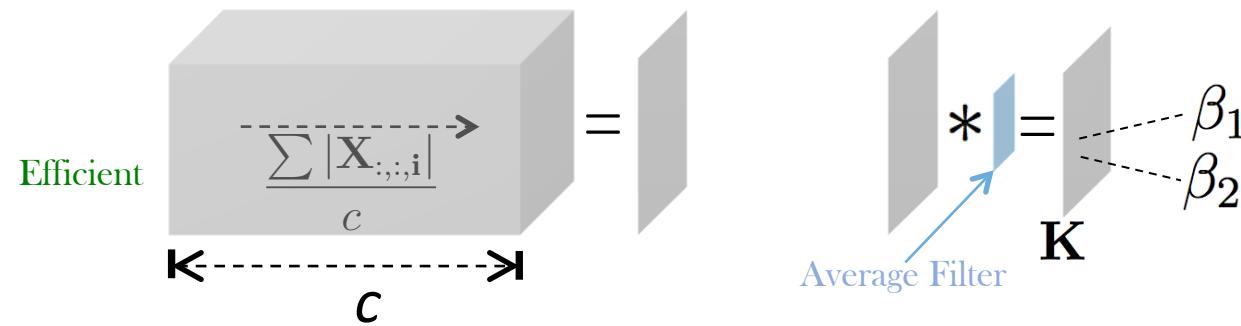
Redundant computation in overlapping areas



$$\begin{aligned} \frac{1}{n} \|\mathbf{X}_1\|_{\ell_1} &= \beta_1 \\ \frac{1}{n} \|\mathbf{X}_2\|_{\ell_1} &= \beta_2 \end{aligned}$$



(2) Binarizing Input



(3) Convolution with XNOR-Bitcount

$$\mathbb{R} * \mathbb{R} \approx$$

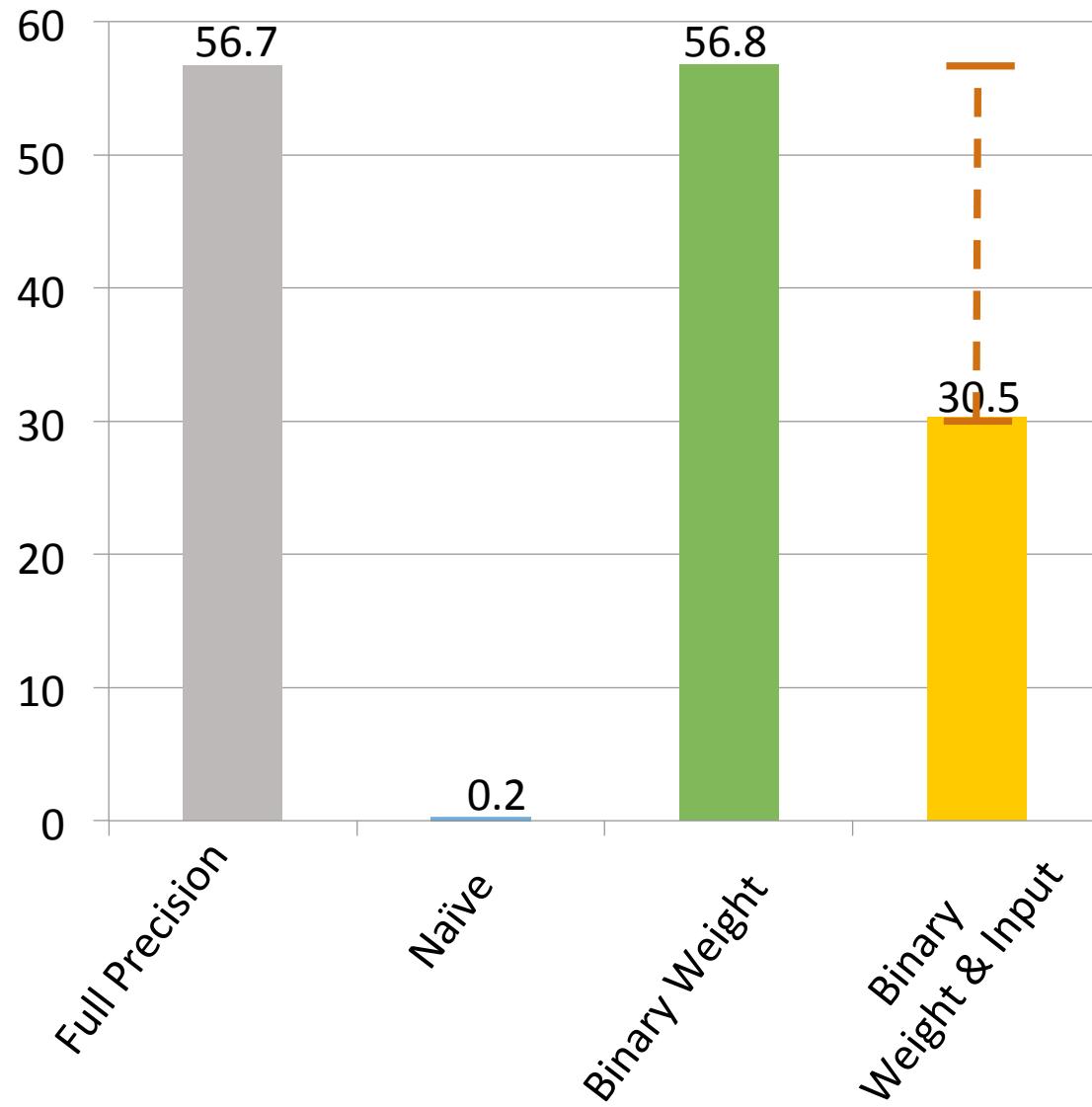
$$\left[\mathbb{B} * \mathbb{B} \right] \odot \circledcirc \alpha$$

sign(X) sign(W)

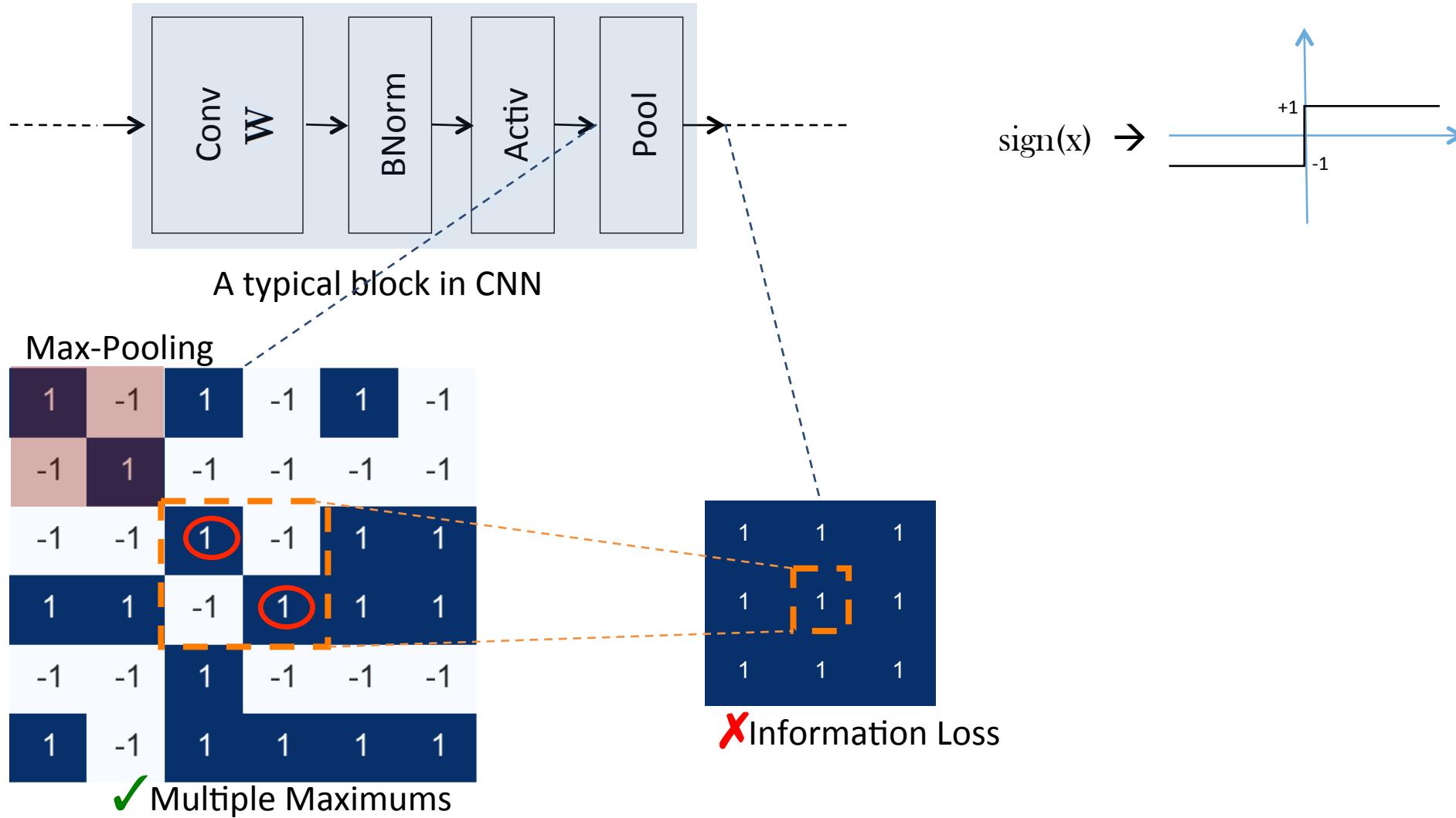
$$\mathbb{R} \times \mathbb{R} \approx \left[\begin{matrix} \mathbb{B} \\ \text{sign}(X) \end{matrix} * \begin{matrix} \mathbb{B} \\ \text{sign}(W) \end{matrix} \right] \odot \beta \odot \alpha$$

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 α, \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} ($\mathbf{W} = \mathbf{W} - \frac{\partial \mathbf{C}}{\partial \mathbf{W}}$)

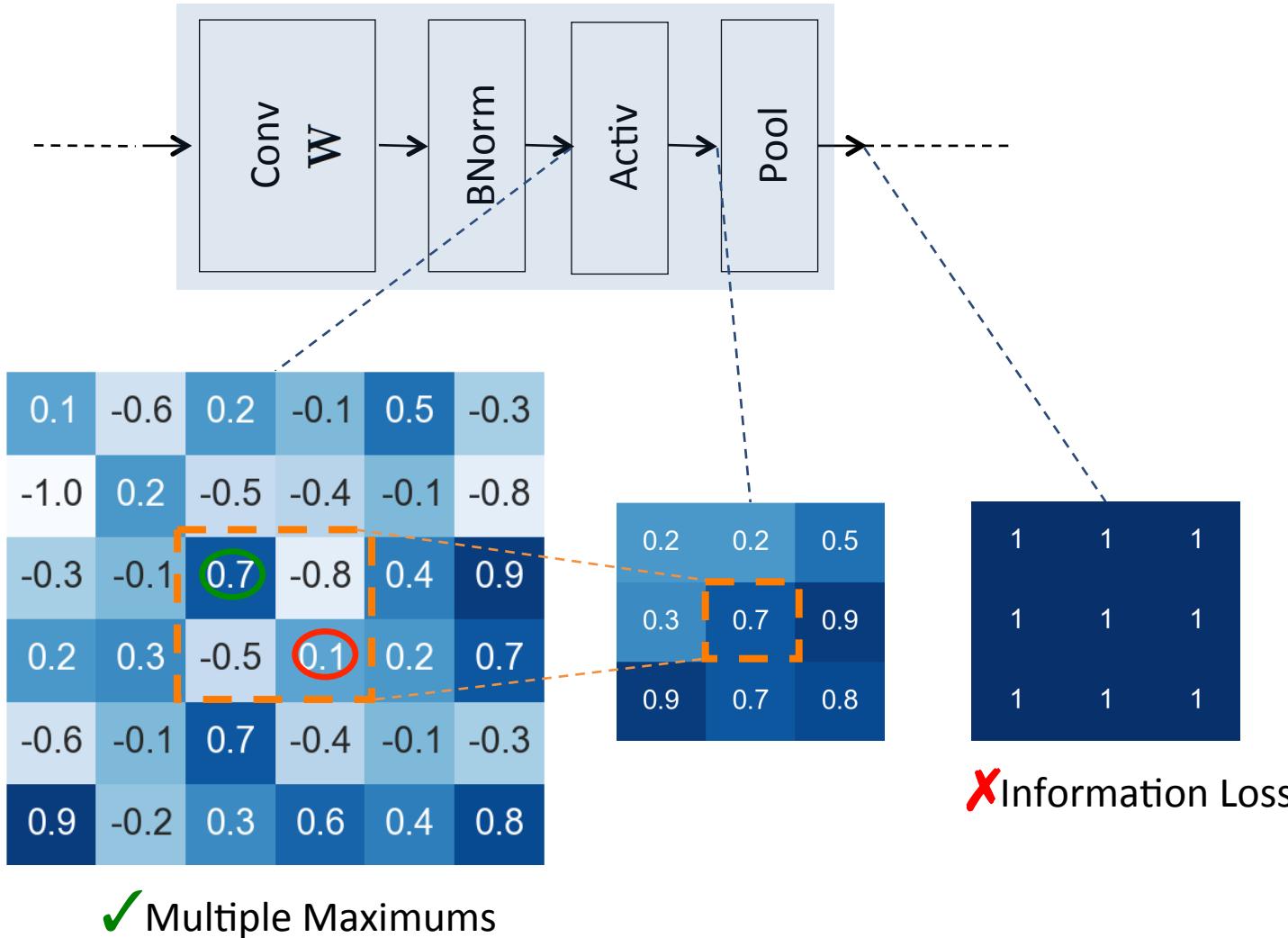
AlexNet Top-1 (%) ILSVRC2012



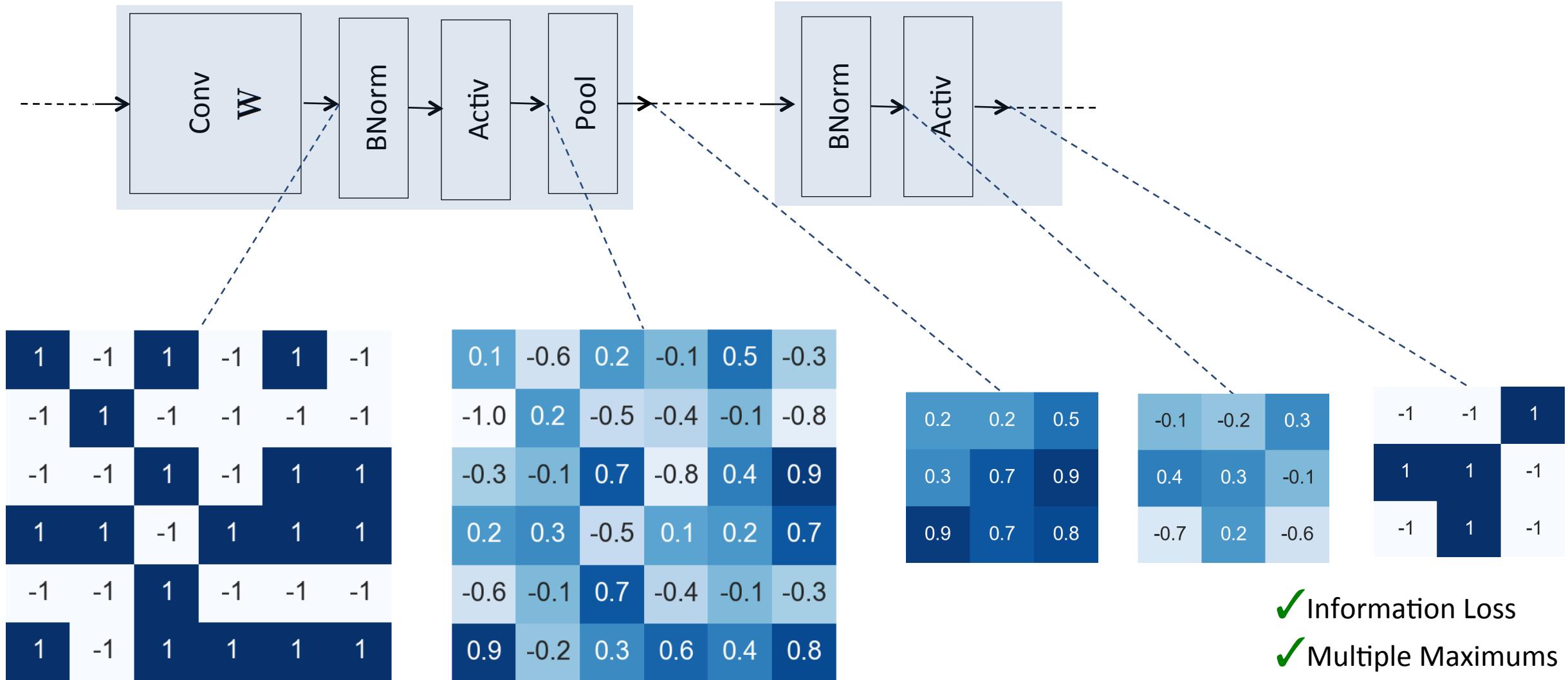
Network Structure in XNOR-Networks



Network Structure in XNOR-Networks

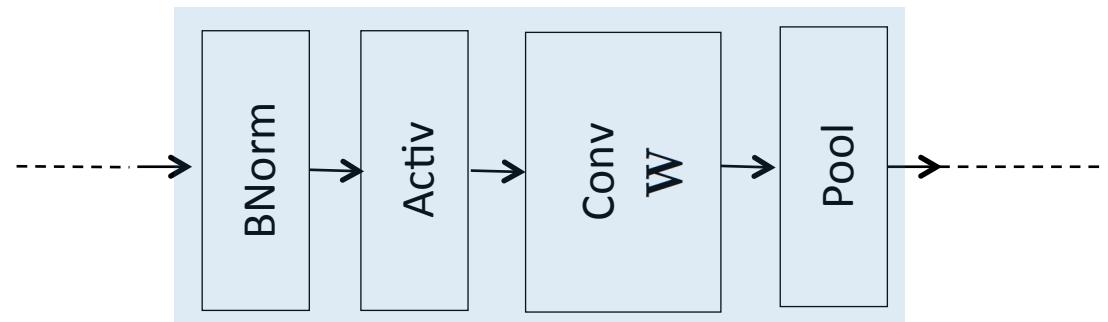


Network Structure in XNOR-Networks

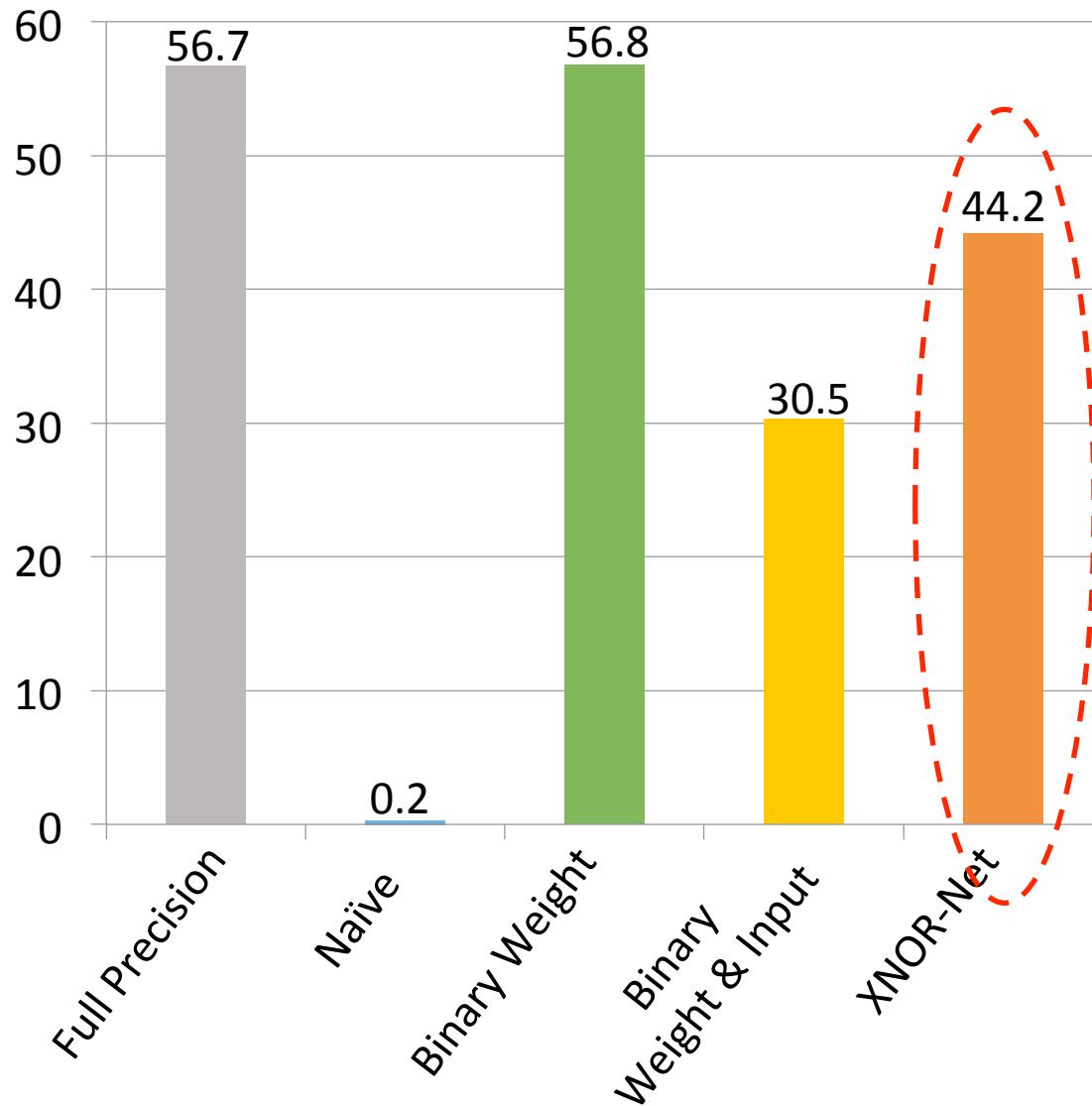


$$\mathbb{R} * \mathbb{R} \approx \left[\begin{matrix} \mathbb{B} \\ \text{sign}(\mathbf{X}) \end{matrix} \right] * \left[\begin{matrix} \mathbb{B} \\ \text{sign}(\mathbf{W}) \end{matrix} \right] \odot \beta \odot \alpha$$

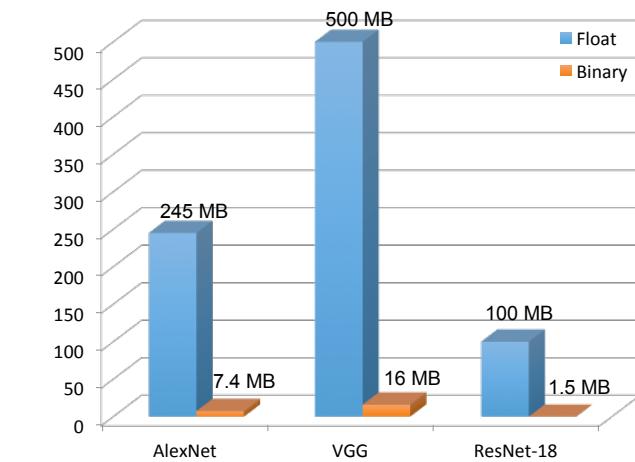
1. Randomly initialize \mathbf{W}
2. For $iter = 1$ to N
3. Load a random input image \mathbf{X}
4. $\mathbf{W}^B = \text{sign}(\mathbf{W})$
5. $\alpha = \frac{\|\mathbf{W}\|_{\ell_1}}{n}$
6. Forward pass with α, \mathbf{W}^B
7. Compute loss function C
8. $\frac{\partial C}{\partial \mathbf{W}}$ = Backward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



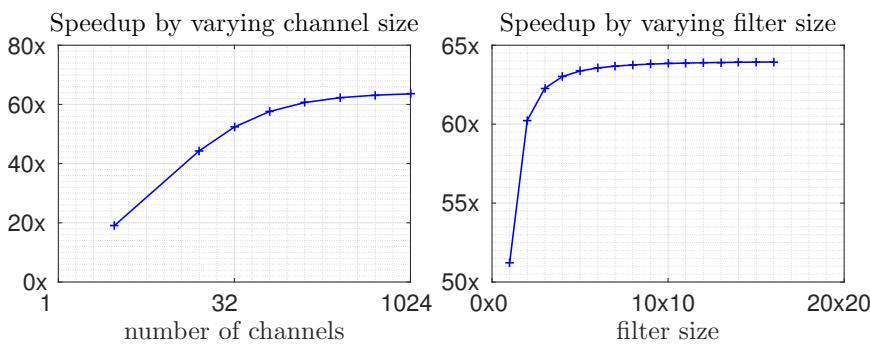
AlexNet Top-1 (%) ILSVRC2012



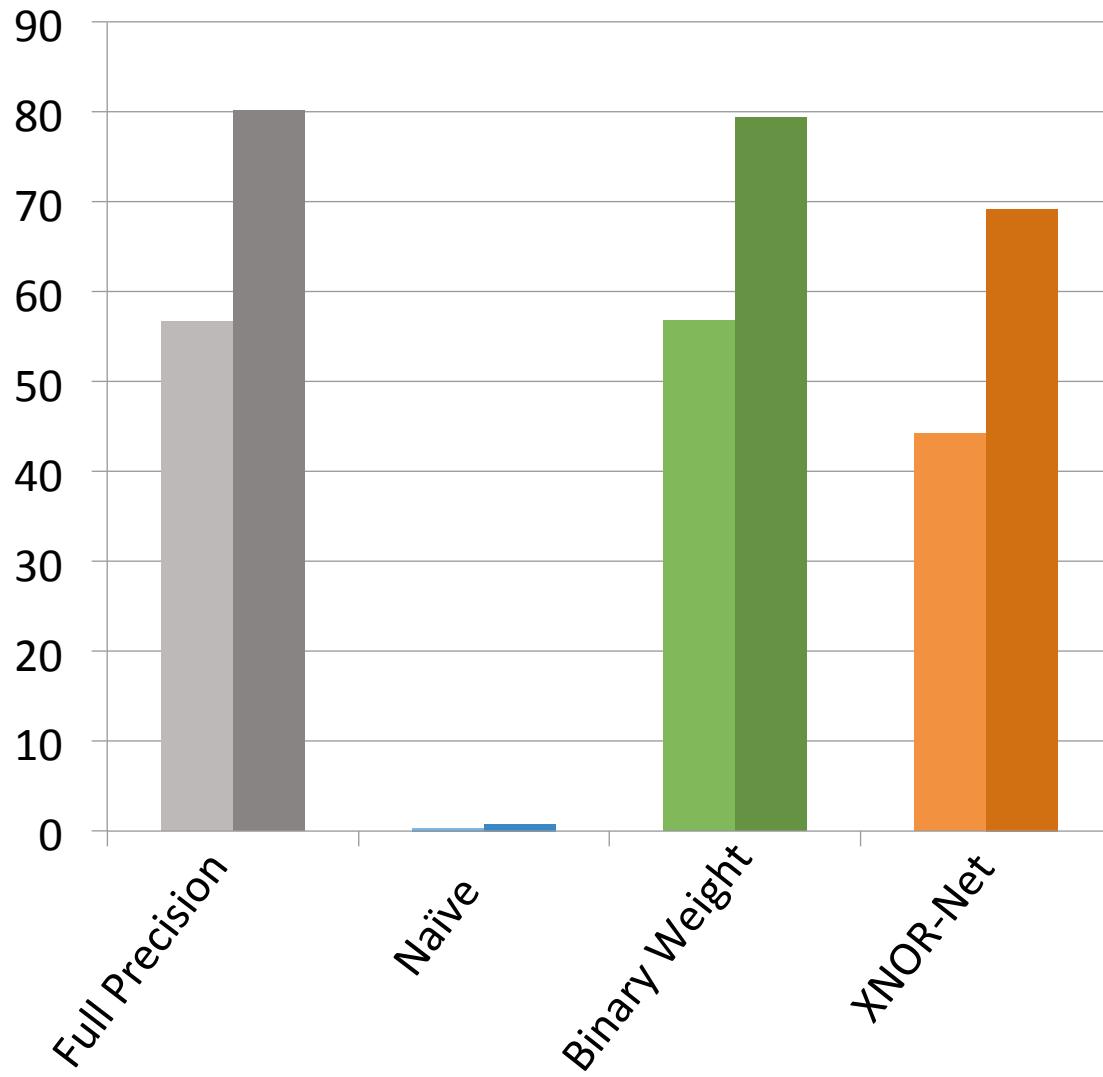
✓ 32x Smaller Model



✓ 58x Less Computation

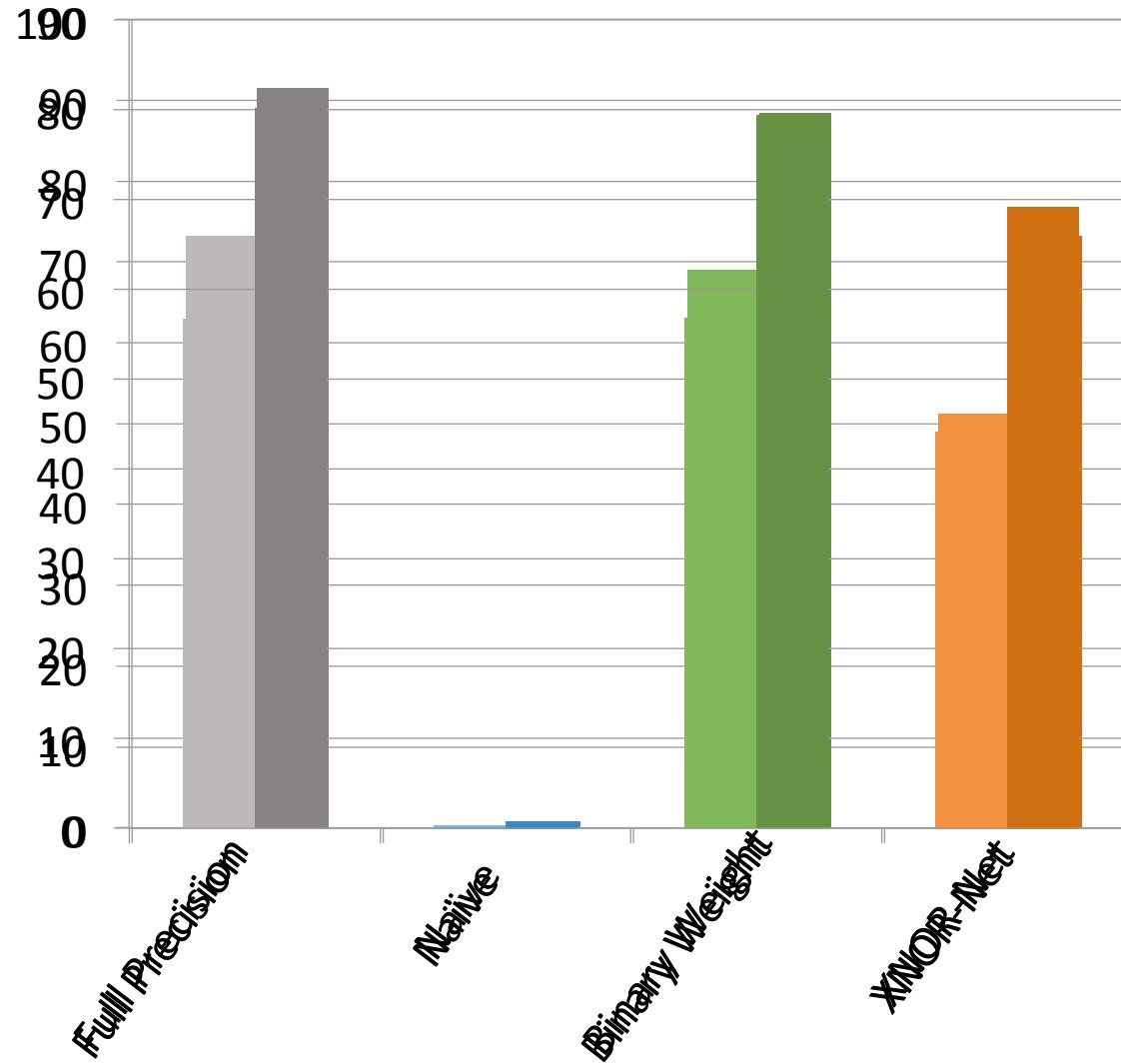


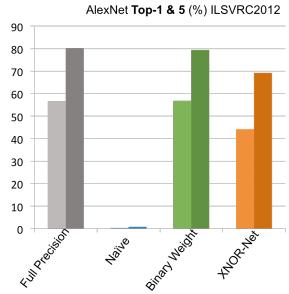
AlexNet **Top-1 & 5 (%)** ILSVRC2012



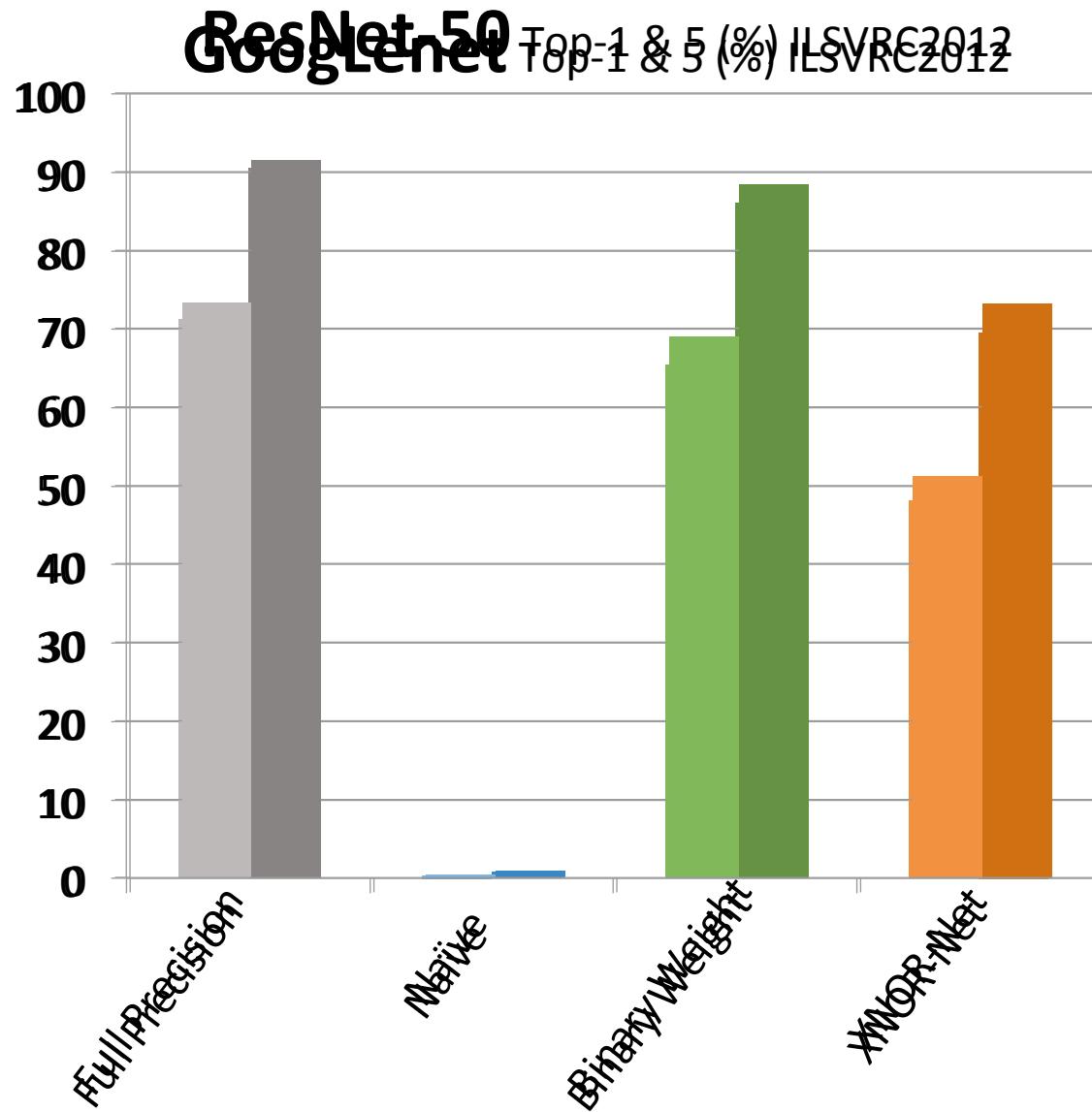
ResNet-50 Top-5 (%) & LSVRC2012

AlexNet

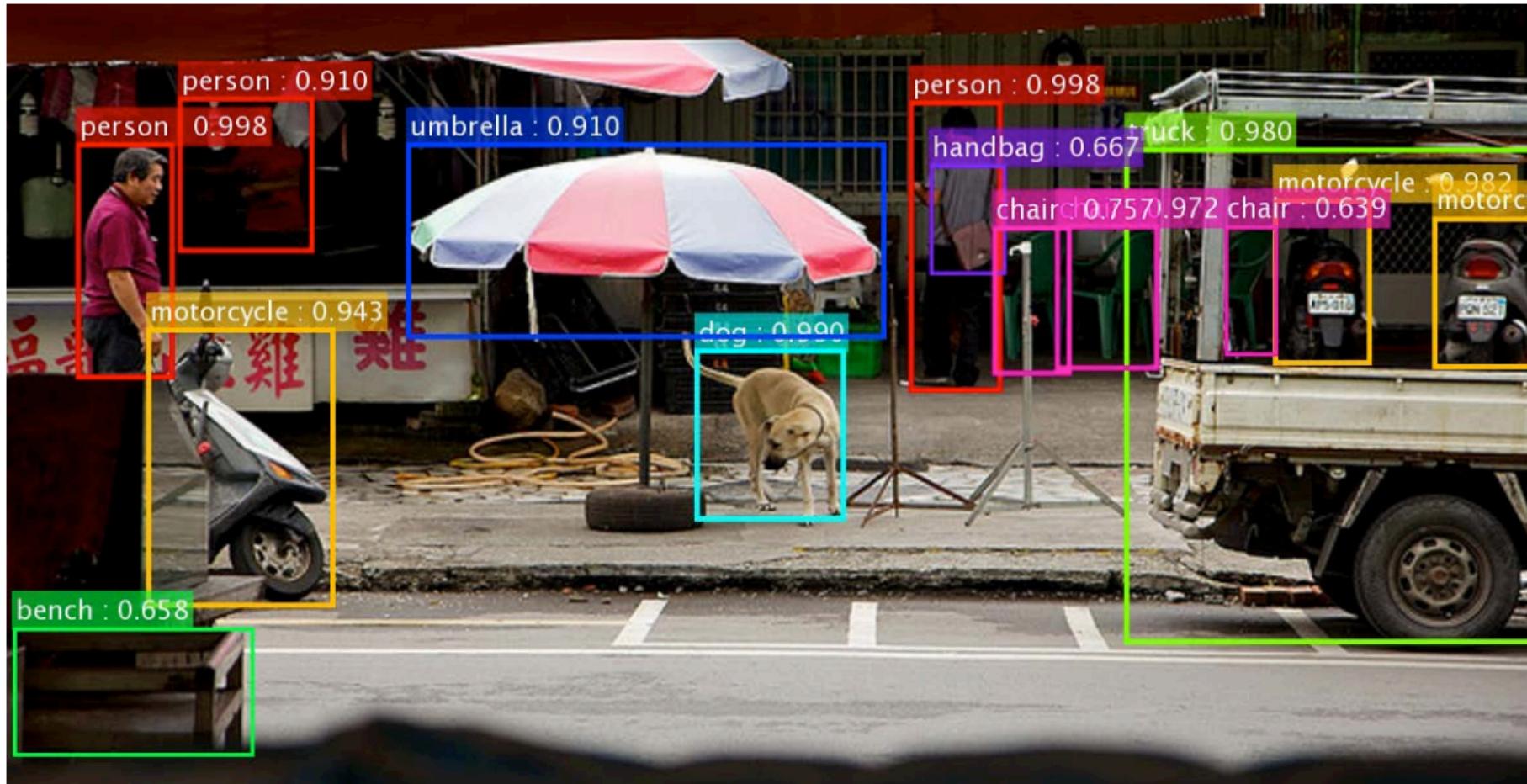




AlexNet



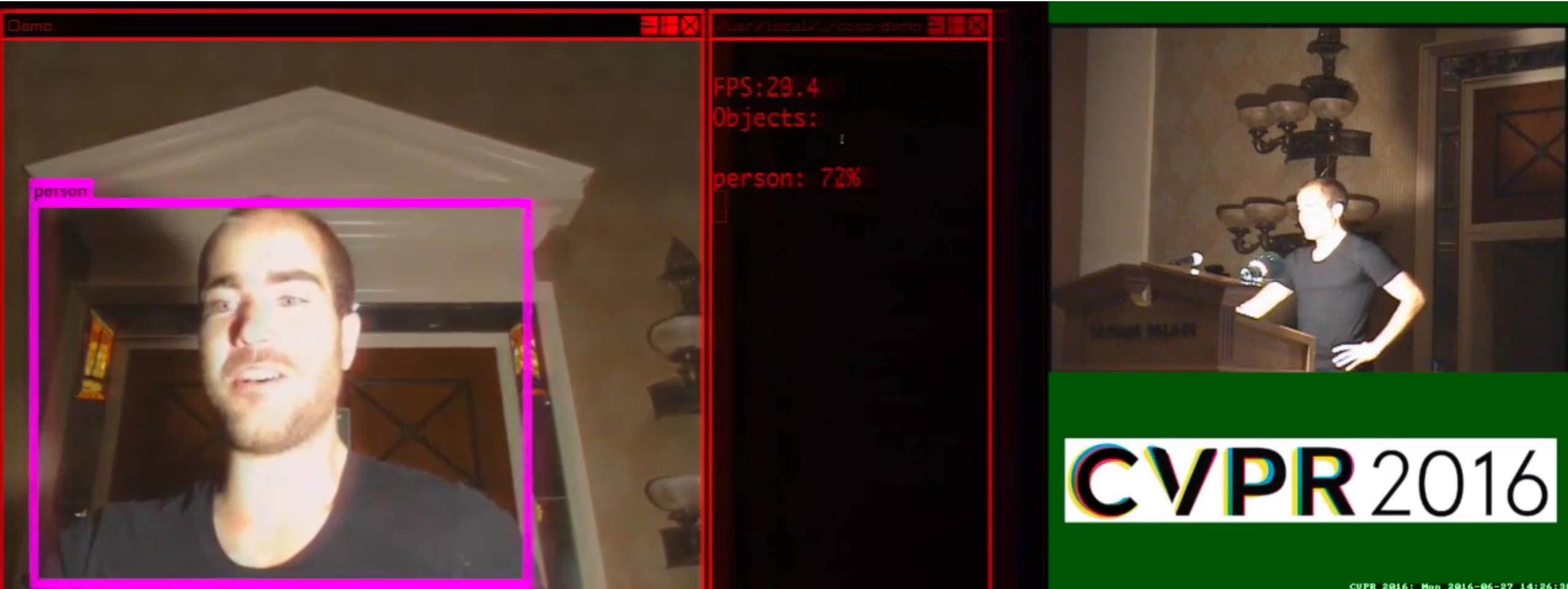
Object Detection



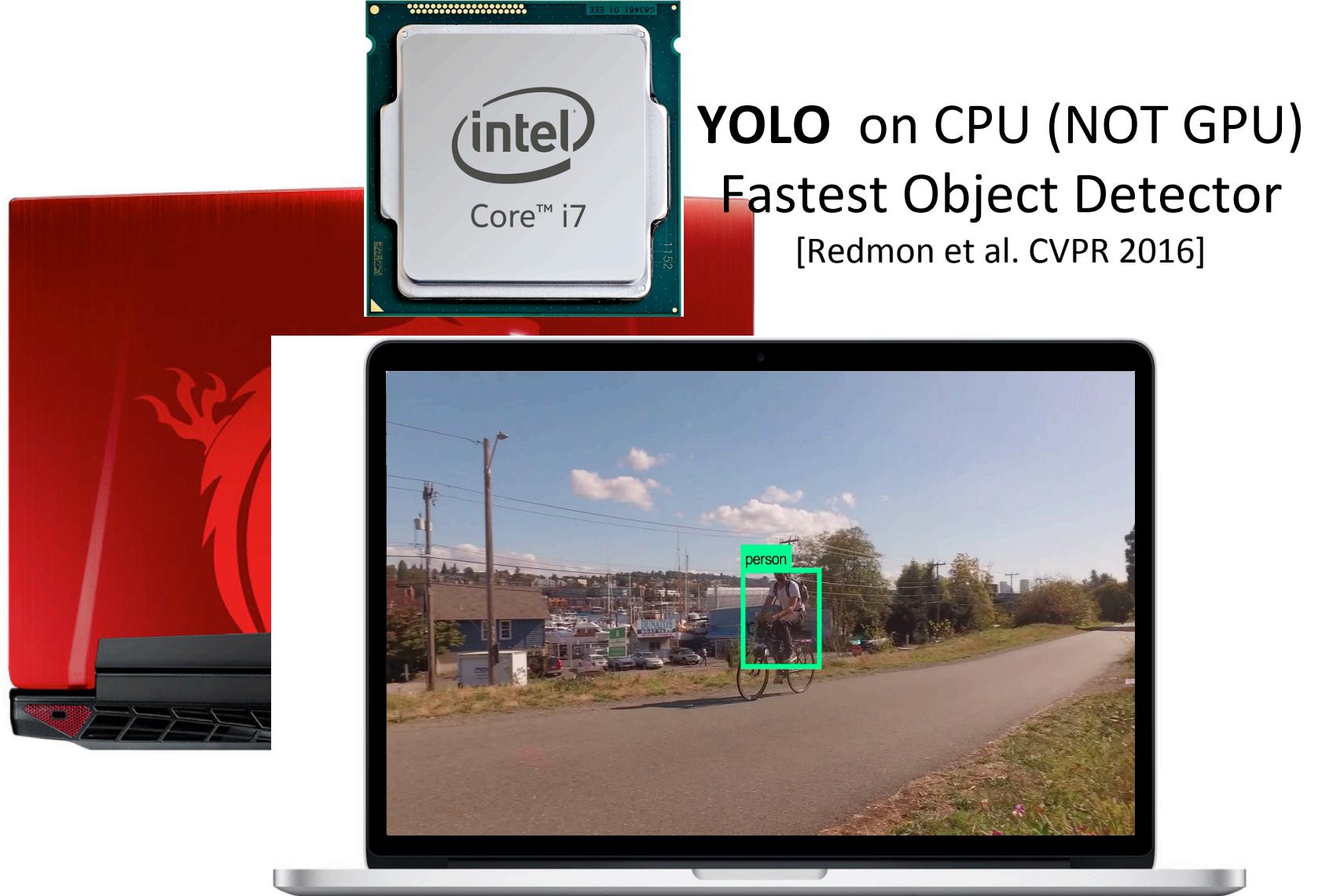
[He et al, 2015]

YOLO: Fastest Object Detector

[Redmon et al. CVPR 2016]



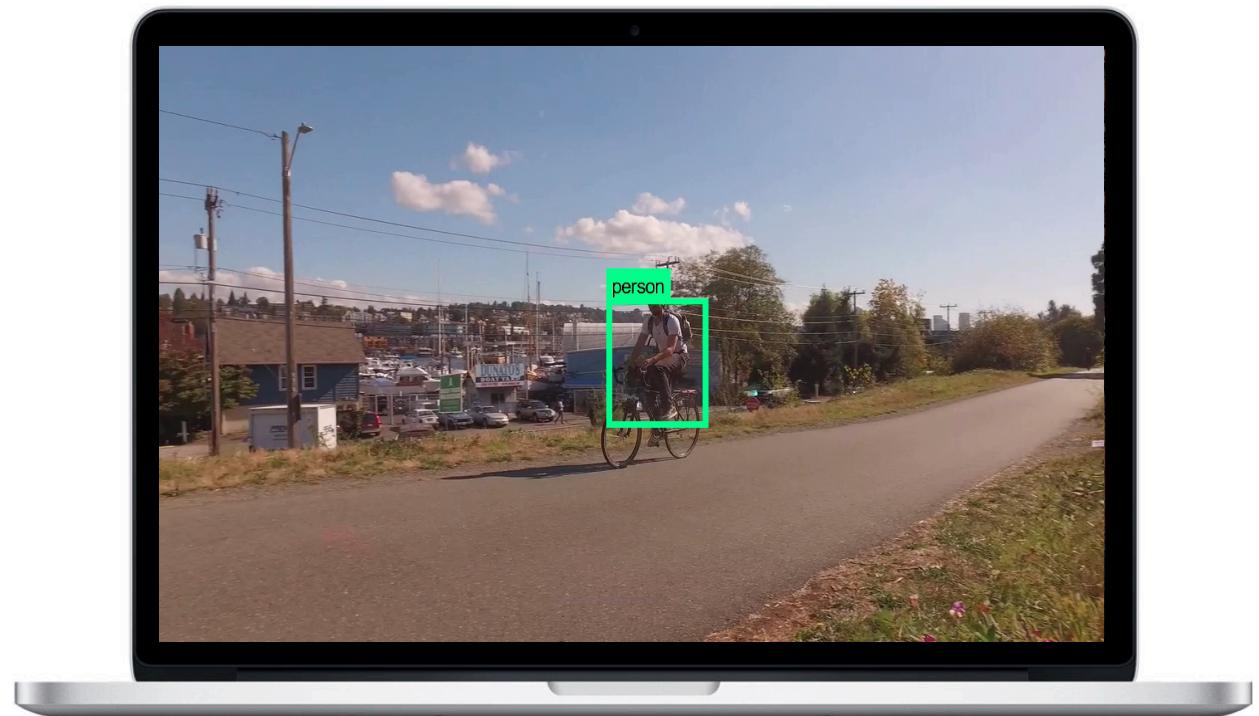






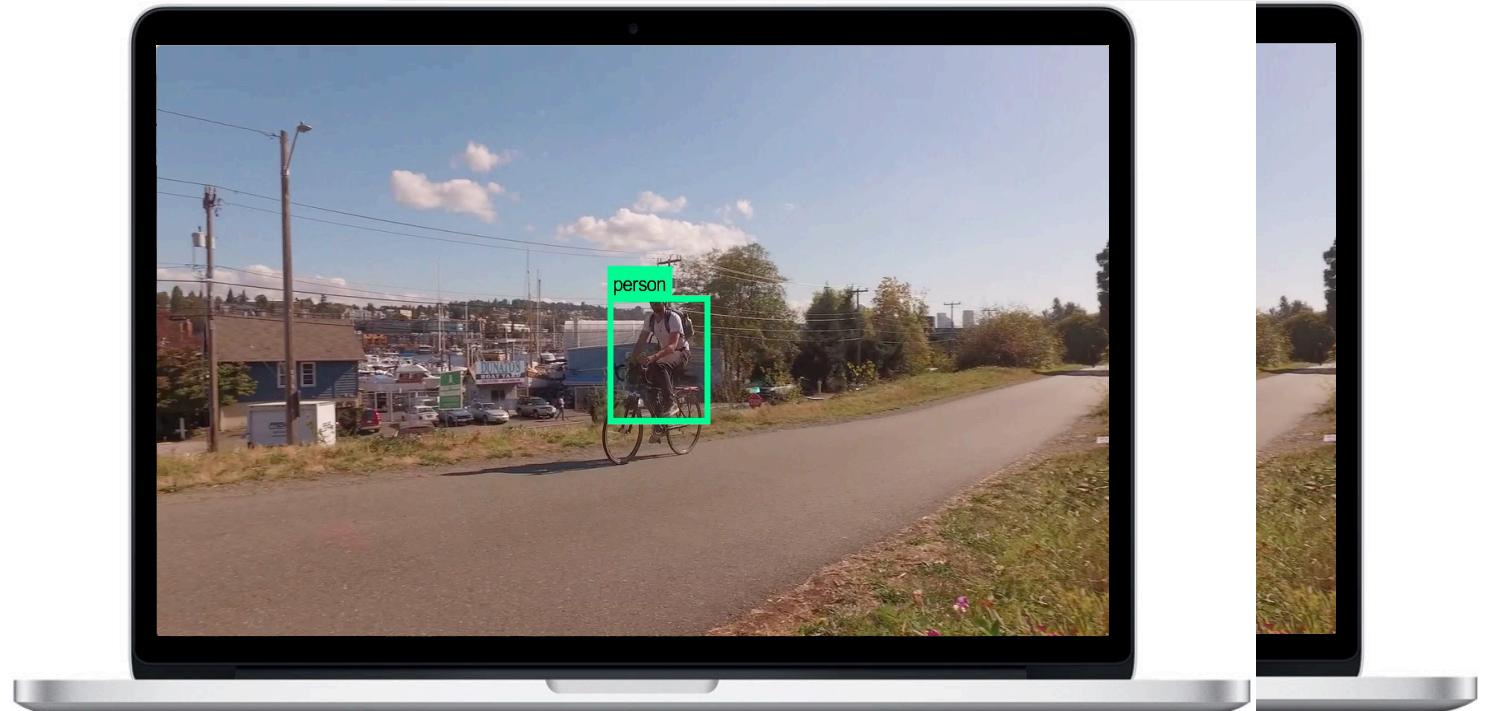
YOLO on CPU (NOT GPU)
Fastest Object Detector

[Redmon et al. CVPR 2016]



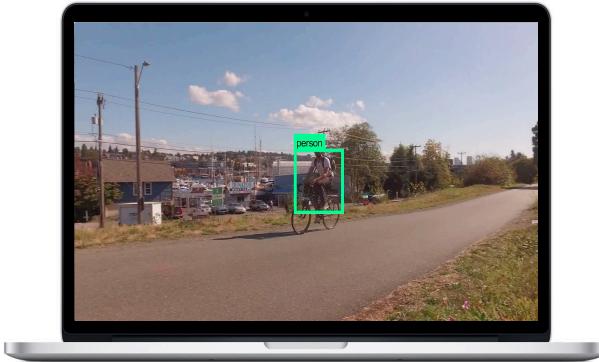


YOLO Method
CPU (NOT GPU)
Fastest Object Detector
[Redmon et al. CVPR 2016]

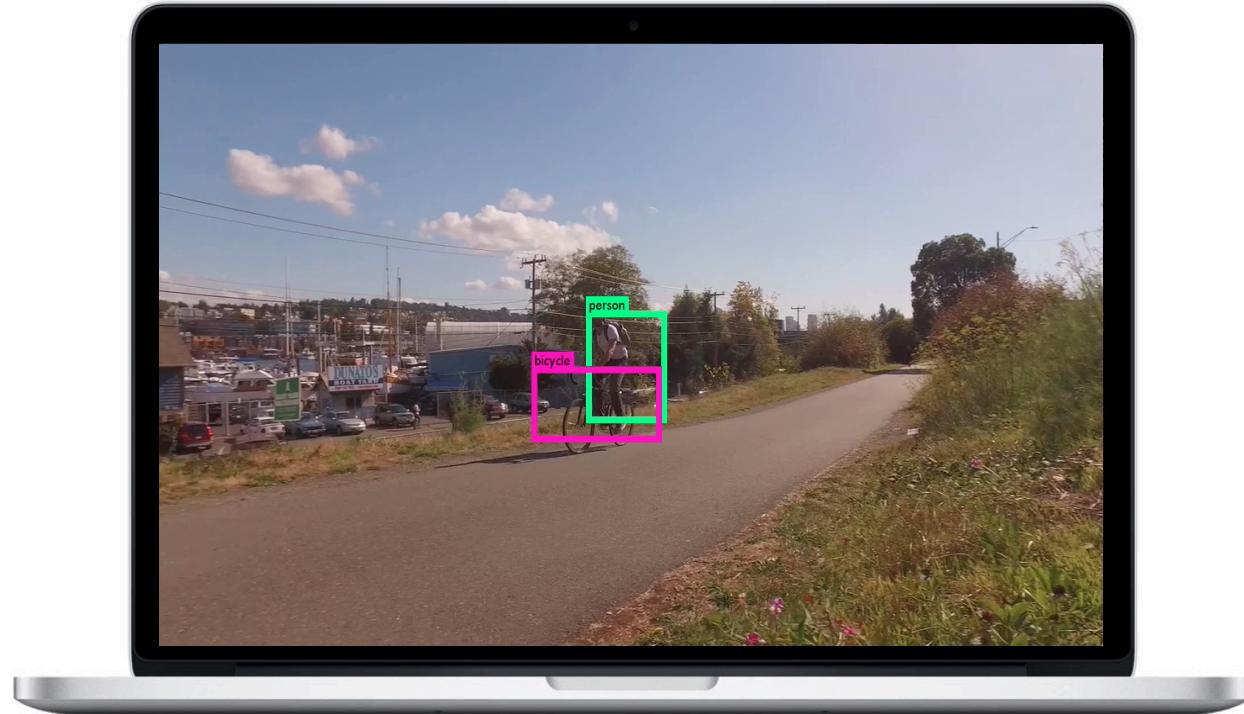




YOLO
Fastest Object Detector
[Redmon et al. CVPR 2016]

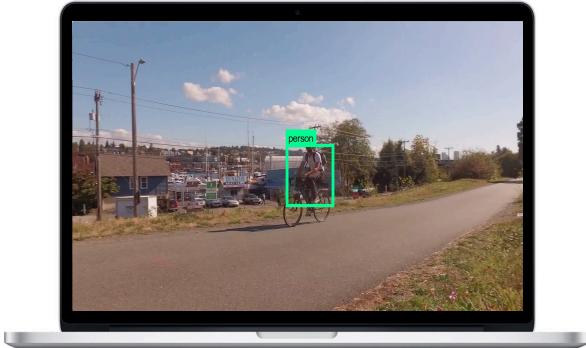


Our Method On CPU (NOT GPU)

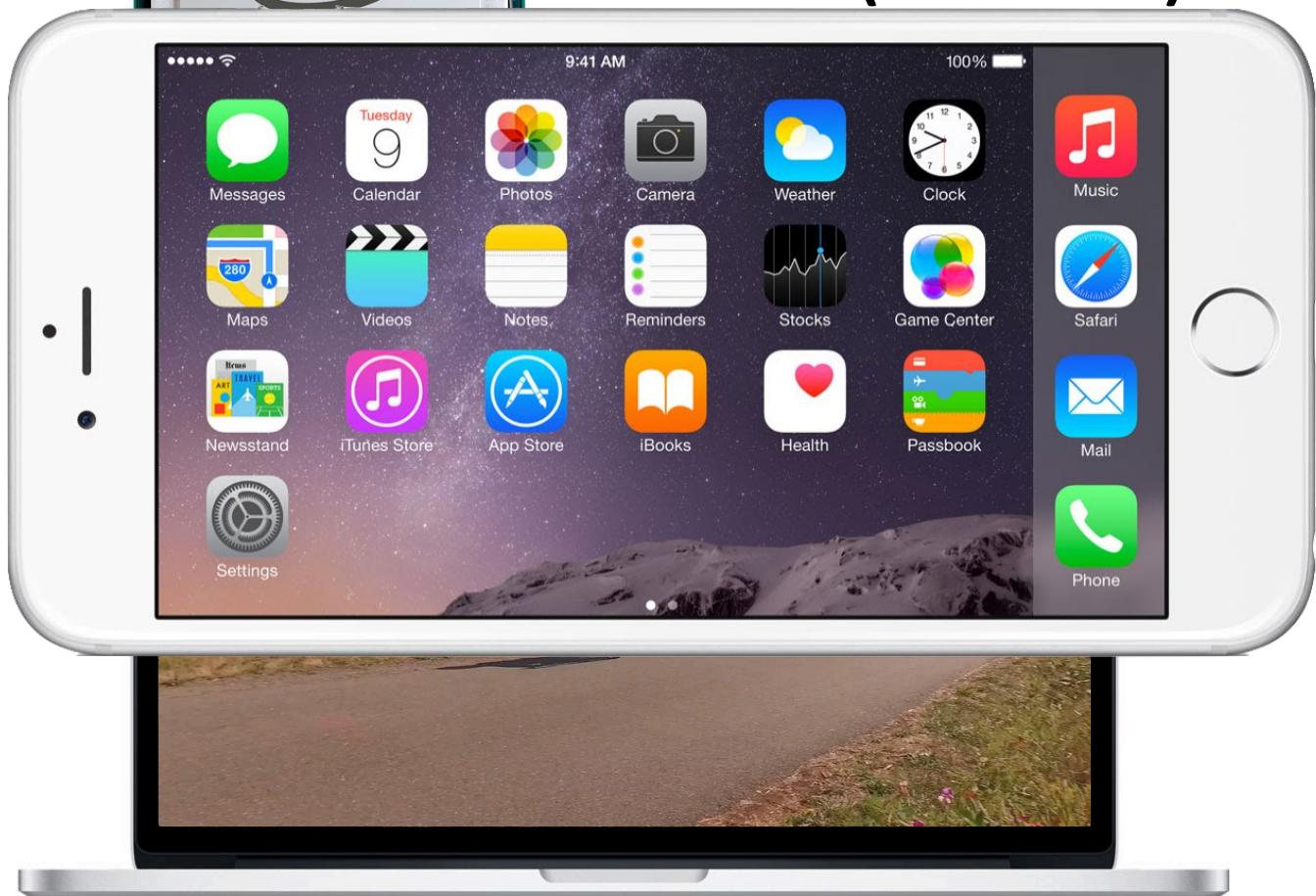




YOLO
Fastest Object Detector
[Redmon et al. CVPR 2016]

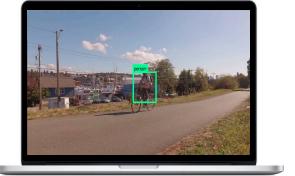


Our Method On CPU (NOT GPU)

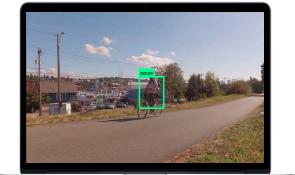




 YOLO
Fastest Object Detector
[Redmon et al. CVPR 2016]



 Our Method

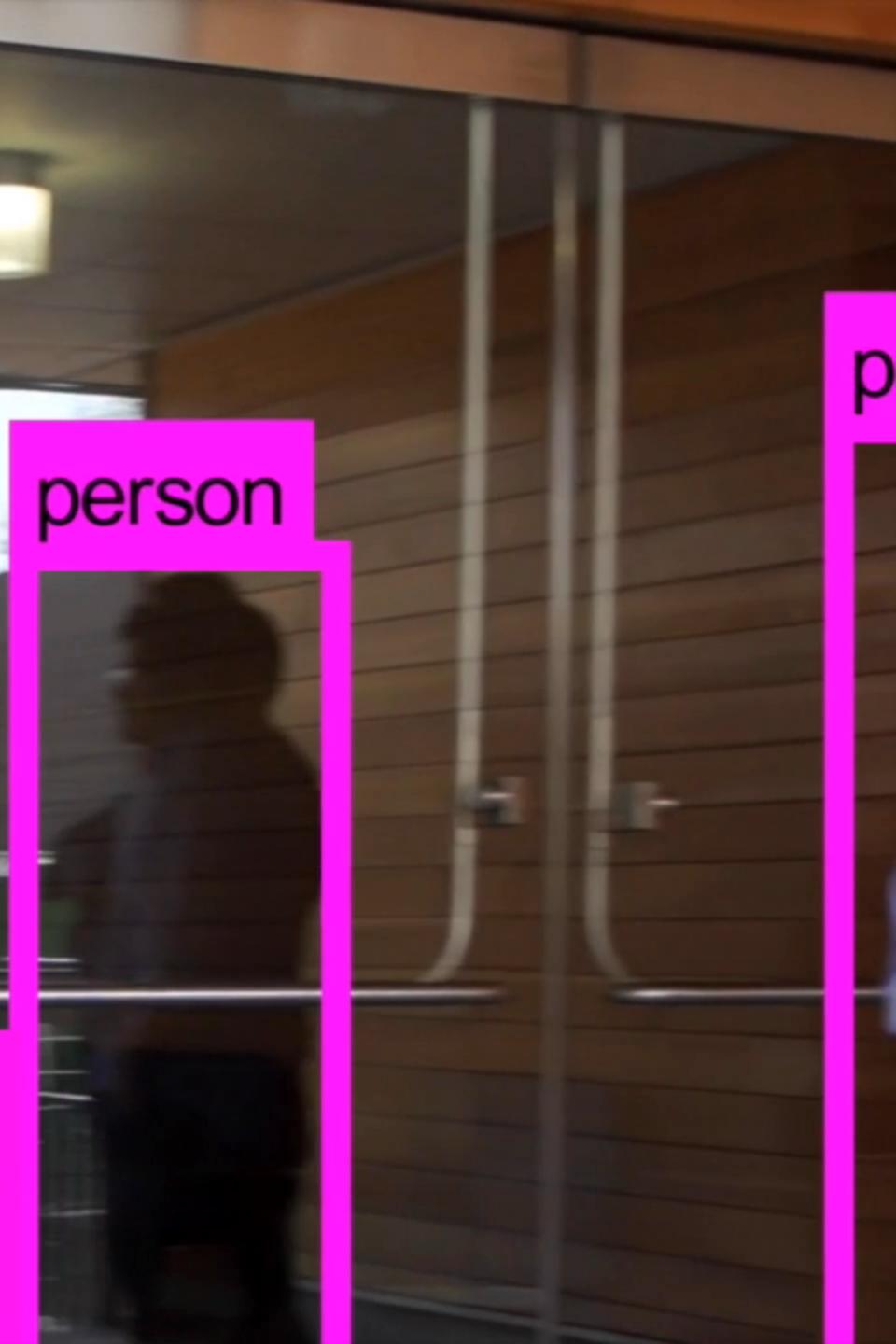


Thank You !

person



person



p

Thank You !



Carlo C. del Mundo

Dmitry Belenko



<https://github.com/allenai/XNOR-Net>