

Binarized Neural Networks

outlines

- **Binarized Neural Network**
- **XNOR Neural Network**
- Brief introduction
 - FINN
 - Residual Binarized Neural Network (ReBennet)
 - BI-Real NETWORK
 - DOREFA-NET

Binarized Neural Networks: Training Neural Networks with Weights and Activations Constrained to +1 or -1

Matthieu Courbariaux*¹

Itay Hubara*²

Daniel Soudry³

Ran El-Yaniv²

Yoshua Bengio^{1,4}

¹Université de Montréal

²Technion Israel Institute of Technology

³Columbia University

⁴CIFAR Senior Fellow

*Indicates equal contribution. Ordering determined by coin flip.

MATTHIEU.COURBARIAUX@GMAIL.COM

ITAYHUBARA@GMAIL.COM

DANIEL.SOUDRY@GMAIL.COM

RANI@CS.TECHNION.AC.IL

YOSHUA.UMONTREAL@GMAIL.COM

Binarized : deterministic & stochastic

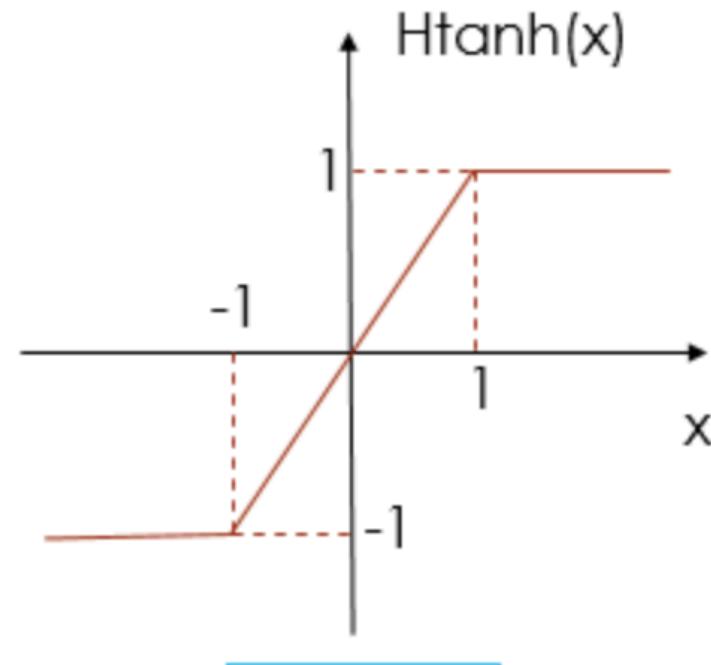
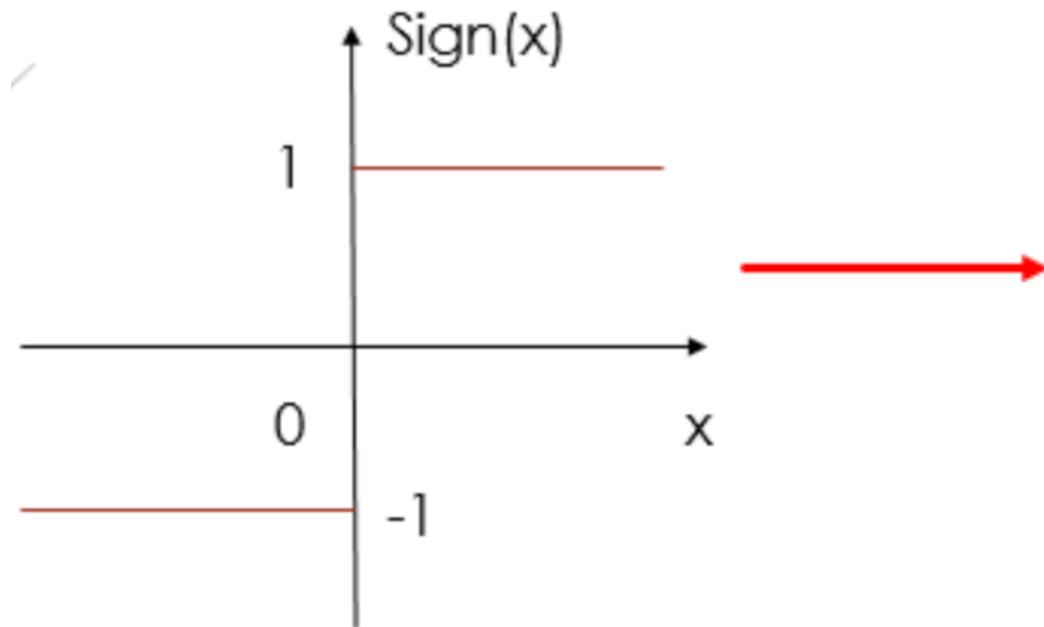
$$x^b = \text{Sign}(x) = \begin{cases} +1 & \text{if } x \geq 0, \\ -1 & \text{otherwise.} \end{cases} \quad \text{Used}$$

$$x^b = \begin{cases} +1 & \text{with probability } p = \sigma(x), \\ -1 & \text{with probability } 1 - p. \end{cases}$$

$$\sigma(x) = \text{clip}\left(\frac{x+1}{2}, 0, 1\right) = \max(0, \min(1, \frac{x+1}{2}))$$

hard tanh

$$q = \text{Sign}(r)$$



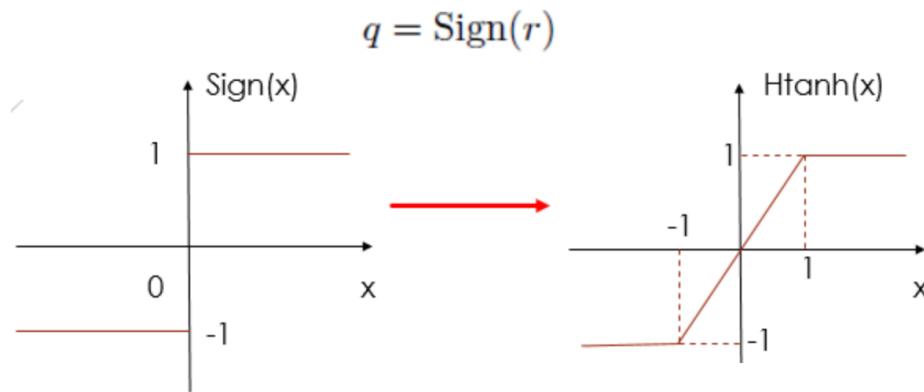
Weight : Binarized

Gradient : ?

- The **magnitude** of the gradient is very **small**
- Additive effect : the noise is generally considered to be a normal distribution, so the cumulative gradient can be used to **average the noise**.

Gradient : real value (high precision)

How to do the propagation (forward & backward)



$$\text{Htanh}(x) = \text{Clip}(x, -1, 1) = \max(-1, \min(1, x))$$

$$q = \text{Sign}(r)$$

$$\frac{\partial C}{\partial r} \quad \frac{\partial C}{\partial q}$$

$$g_r = g_q \mathbf{1}_{|r| \leq 1}.$$

Tips :

- For weight :
- clipping the weights during training
- When using a weight w^r , quantize it using $w^b = \text{Sign}(w^r)$.

1.Computing the parameters' gradient

1.1 Forward propagation

```
for  $k = 1$  to  $L$  do
     $W_k^b \leftarrow \text{Binarize}(W_k)$ 
     $s_k \leftarrow a_{k-1}^b W_k^b$ 
     $a_k \leftarrow \text{BatchNorm}(s_k, \theta_k)$ 
    if  $k < L$  then
         $a_k^b \leftarrow \text{Binarize}(a_k)$ 
    end if
end for
```

1.2. Backward propagation

{Please note that the gradients are not binary.}

Compute $g_{a_L} = \frac{\partial C}{\partial a_L}$ knowing a_L and a^*

for $k = L$ to 1 **do**

if $k < L$ **then**

$$g_{a_k} \leftarrow g_{a_k^b} \circ 1_{|a_k| \leq 1}$$

$$\text{Htanh}(x) = \text{Clip}(x, -1, 1) = \max(-1, \min(1, x))$$

end if

$$(g_{s_k}, g_{\theta_k}) \leftarrow \text{BackBatchNorm}(g_{a_k}, s_k, \theta_k)$$

$$g_{a_{k-1}^b} \leftarrow g_{s_k} W_k^b$$

$$g_{W_k^b} \leftarrow g_{s_k}^\top a_{k-1}^b$$

end for

{2. Accumulating the parameters' gradient:}

for $k = 1$ to L **do**

$$\theta_k^{t+1} \leftarrow \text{Update}(\theta_k, \eta, g_{\theta_k})$$

$$W_k^{t+1} \leftarrow \text{Clip}(\text{Update}(W_k, \gamma_k \eta, g_{W_k^b}), -1, 1)$$

$$\eta^{t+1} \leftarrow \lambda \eta$$

talk later

Weight is binary

end for

Shift BN

trick

Algorithm 2 Shift based Batch Normalizing Transform, applied to activation x over a mini-batch. $AP2(x) = \text{sign}(x) \times 2^{\text{round}(\log_2|x|)}$ is the approximate power-of-2, and $\ll\gg$ stands for **both** left and right binary shift.

Estimate the value using the base 2 exponential function

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{Var[x^{(k)}]}}$$

Require: Values of x over a mini-batch: $B = \{x_1 \dots m\}$;
Parameters to be learned: γ, β

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

Ensure: $\{y_i = \text{BN}(x_i, \gamma, \beta)\}$

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad \{\text{mini-batch mean}\}$$

$$C(x_i) \leftarrow (x_i - \mu_B) \quad \{\text{centered input}\}$$

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (C(x_i) \ll\gg AP2(C(x_i))) \quad \{\text{apx variance}\}$$

$$\hat{x}_i \leftarrow C(x_i) \ll\gg AP2((\sqrt{\sigma_B^2 + \epsilon})^{-1}) \quad \{\text{normalize}\}$$

denominator>0

$$y_i \leftarrow AP2(\gamma) \ll\gg \hat{x}_i \quad \{\text{scale and shift}\}$$

1 SGD

2 Momentum

3 Nesterov

4 Adagrad

5 Adadelta

6 RMSprop

7 Adam

8 AdaMax

9 Nadam

Shift AdaMax

Algorithm 3 Shift based AdaMax learning rule (Kingma

& Ba, 2014). g_t^2 indicates the elementwise square $g_t \circ g_t$.

Good default settings are $\alpha = 2^{-10}$, $1 - \beta_1 = 2^{-3}$, $1 - \beta_2 = 2^{-10}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

Require: Previous parameters θ_{t-1} and their gradient g_t , and learning rate α .

Ensure: Updated parameters θ_t

{Biased 1st and 2nd raw moment estimates:}

$$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$$

$$v_t \leftarrow \max(\beta_2 \cdot v_{t-1}, |g_t|)$$

{Updated parameters:}

$$\theta_t \leftarrow \theta_{t-1} - (\alpha \ll\gg (1 - \beta_1)) \cdot \hat{m} \ll\gg v_t^{-1})$$

<https://blog.csdn.net/fishmai/article/details/52510826>

$$g_{a_k} = g_{a_k^b} \circ |_{|a_k| \leq 1} \quad (\text{Htanh}(1))$$

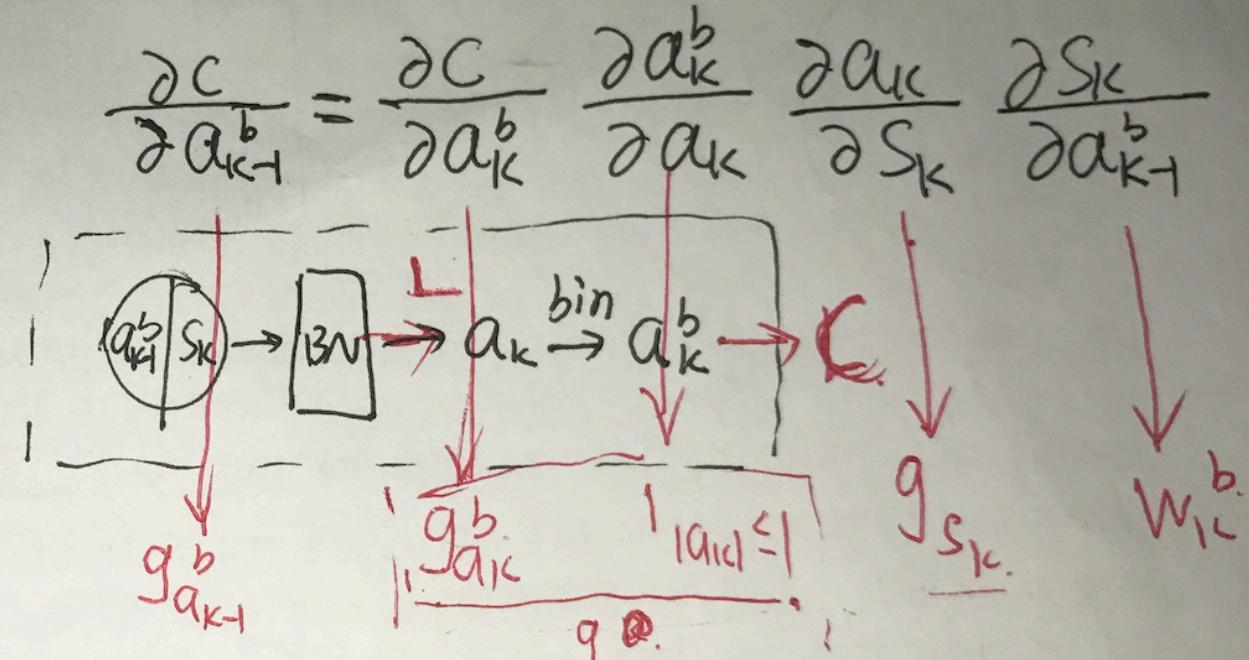
$$\frac{\partial C}{\partial a_k} = \frac{\partial C}{\partial a_k^b} \cdot \frac{\partial a_k^b}{\partial a_k}$$

$$\boxed{a_{k-1}^b W_k^b} \quad \boxed{a_k^b}$$

~~Back~~ BackBN(g_{a_k}, S_k, θ_k).

$$Y = \gamma X + \beta \quad \downarrow S_k \quad \rightarrow g_{\theta_k}$$

γ, β



$$\frac{\partial L}{\partial a_{k-1}^b} = \frac{\partial L}{\partial a_k^b} \cdot \frac{\partial a_k^b}{\partial a_{k-1}^b} = g_{S_k} W_k^b$$

$$\frac{\partial L}{\partial W_k^b} = \frac{\partial L}{\partial S_k^b} \frac{\partial S_k^b}{\partial W_k^b} = g_{S_k^T}^T a_{k-1}^b$$

Algorithm 4 Running a BNN. L is the number of layers.

Require: a vector of 8-bit inputs a_0 , the binary weights W^b and the BatchNorm parameters θ .

Ensure: the MLP output a_L .

{1. First layer:}

$$a_1 \leftarrow 0$$

for $n = 1$ to 8 **do**

$$a_1 \leftarrow a_1 + 2^{n-1} \times \text{XnorDotProduct}(a_0^n, W_1^b)$$

end for

$$a_1^b \leftarrow \text{Sign}(\text{BatchNorm}(a_1, \theta_1))$$

{2. Remaining hidden layers:}

for $k = 2$ to $L - 1$ **do**

$$a_k \leftarrow \text{XnorDotProduct}(a_{k-1}^b, W_k^b)$$

$$a_k^b \leftarrow \text{Sign}(\text{BatchNorm}(a_k, \theta_k))$$

end for

{3. Output layer:}

$$a_L \leftarrow \text{XnorDotProduct}(a_{L-1}^b, W_L^b)$$

$$a_L \leftarrow \text{BatchNorm}(a_L, \theta_L)$$

Image pixels to binary value: it uses 8 bit to represent one pixel.

For example, cifar-10: $32 \times 32 = 1024$
Inputs is a 1024 8 bit vector.

Core trick : **bit manipulation**
talk later

Also can consider it as
Binary activation function

XnorDotProduct

- XNOR ----> ~ XOR

$$a = [1, -1, 1, 1, -1]$$

$$w = [-1, 1, 1, -1, -1]$$

$$a = [1, 0, 1, 1, 0]$$

$$w = [0, 1, 1, 0, 0]$$

$$a \wedge w = [1^0, 0^1, 1^1, 1^0, 0^0] = 11010$$

$$\text{Popcount}(a \wedge w) = 3$$

$$-(2 * \text{Popcount}(a \wedge w) - 5) = -1$$

$$a_1 * w_1 + a_2 * w_2 + a_3 * w_3 + a_4 * w_4 + a_5 * w_5 = 1 * -1 + -1 * 1 + 1 * 1 + 1 * -1 + -1 * -1 = -1$$

$$\vec{s}_x \begin{array}{|c|c|c|c|} \hline +1 & -1 & -1 & +1 \\ \hline \end{array}$$
$$\vec{s}_w \begin{array}{|c|c|c|c|} \hline -1 & +1 & -1 & -1 \\ \hline \end{array}$$

Multiply

$$\begin{array}{|c|c|c|c|} \hline -1 & -1 & +1 & -1 \\ \hline \end{array}$$

Sum

$$-2$$

$$\vec{b}_x \begin{array}{|c|c|c|c|} \hline 1 & 0 & 0 & 1 \\ \hline \end{array}$$
$$\vec{b}_w \begin{array}{|c|c|c|c|} \hline 0 & 1 & 0 & 0 \\ \hline \end{array}$$

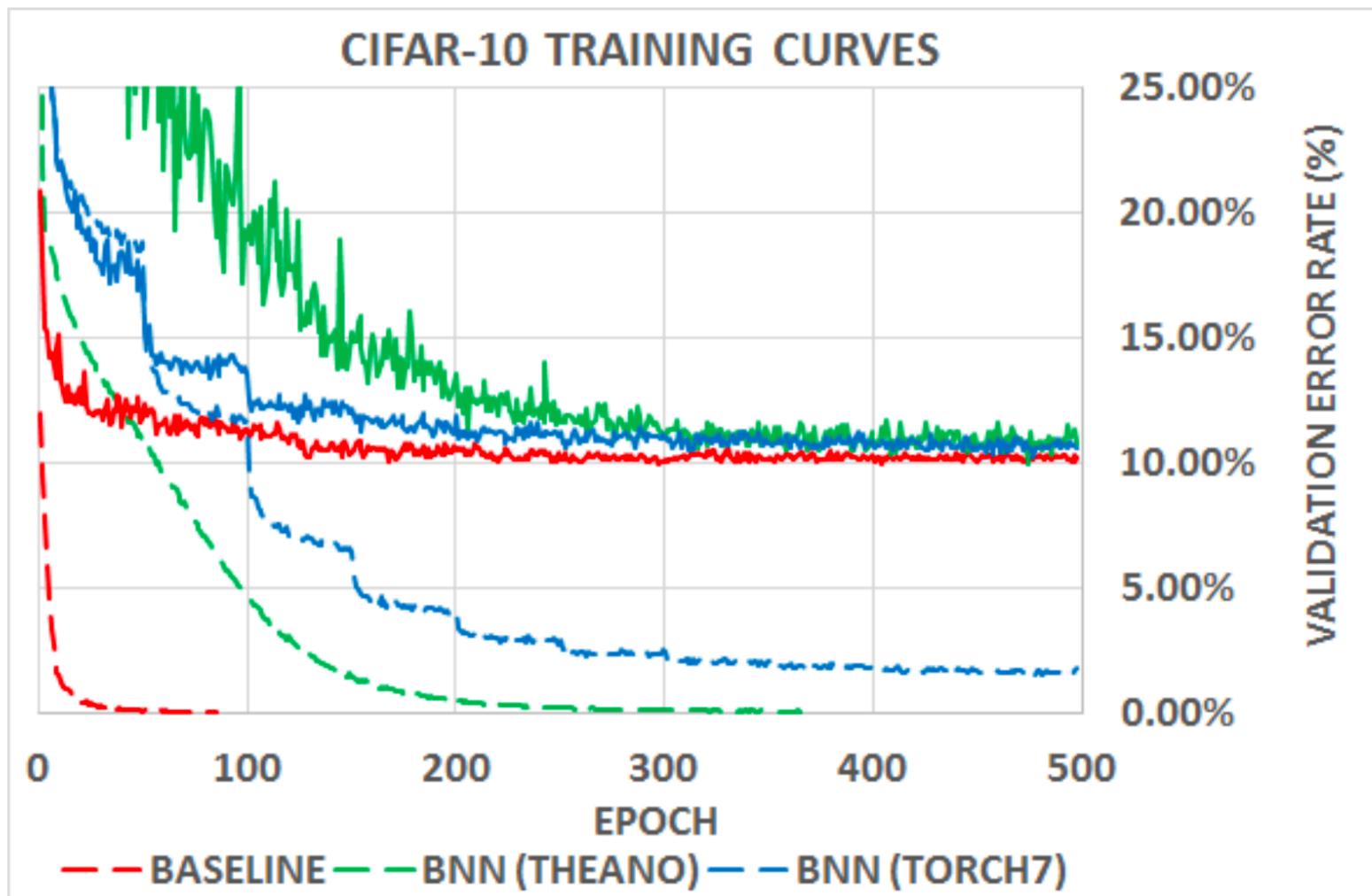
XNOR

$$\begin{array}{|c|c|c|c|} \hline 0 & 0 & 1 & 0 \\ \hline \end{array}$$

popcount
($2p-N$)

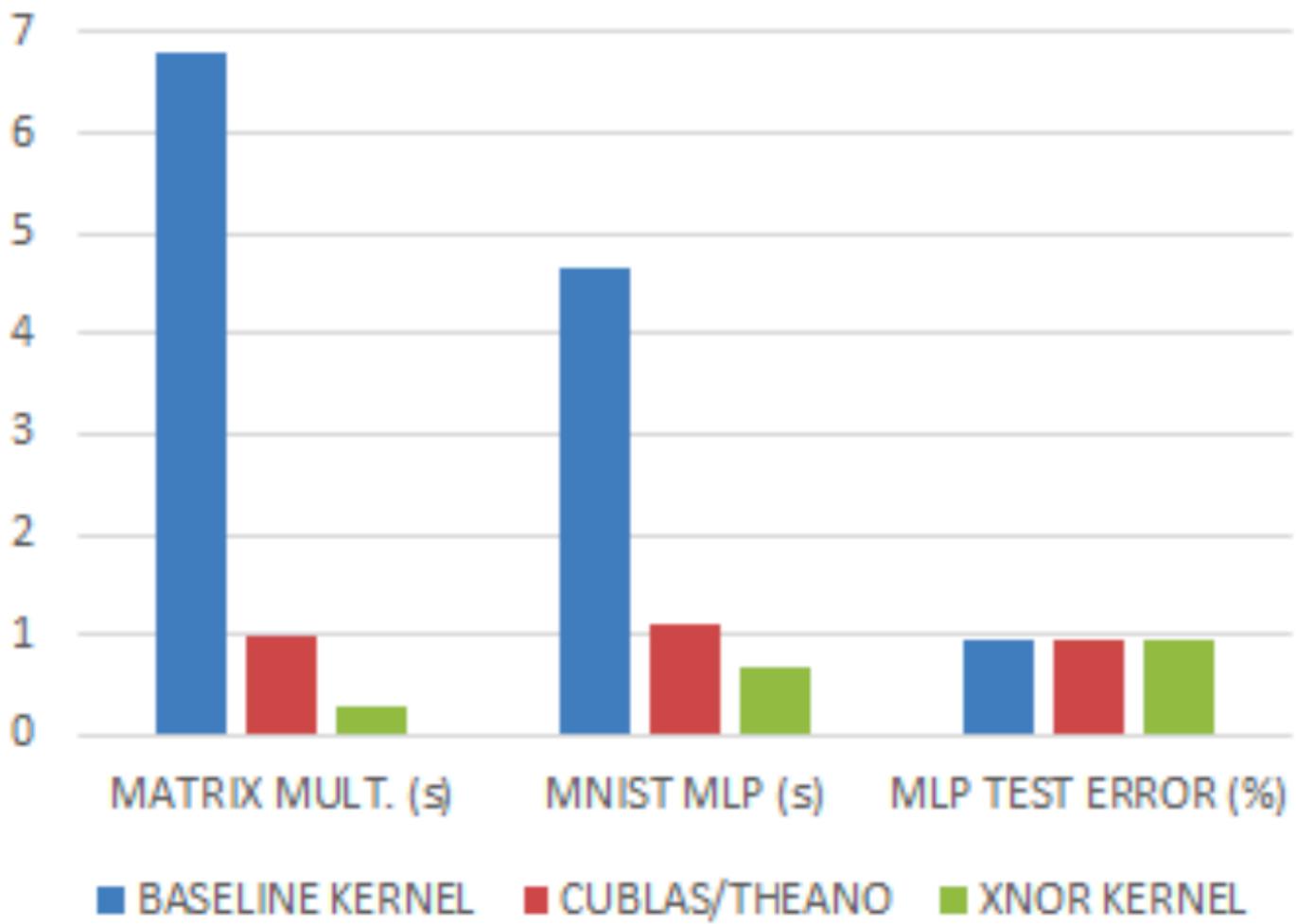
$$-2$$

Experiments: CIFAR-10 Training



Data set	MNIST	SVHN	CIFAR-10
Binarized activations+weights, during training and test			
BNN (Torch7)	1.40%	2.53%	10.15%
BNN (Theano)	0.96%	2.80%	11.40%
Committee Machines' Array (Baldassi et al., 2015)	1.35%	-	-
Binarized weights, during training and test			
BinaryConnect (Courbariaux et al., 2015)	$1.29 \pm 0.08\%$	2.30%	9.90%
Binarized activations+weights, during test			
EBP (Cheng et al., 2015)	$2.2 \pm 0.1\%$	-	-
Bitwise DNNs (Kim & Smaragdis, 2016)	1.33%	-	-
Ternary weights, binary activations, during test			
(Hwang & Sung, 2014)	1.45%	-	-
No binarization (standard results)			
Maxout Networks (Goodfellow et al.)	0.94%	2.47%	11.68%
Network in Network (Lin et al.)	-	2.35%	10.41%
Gated pooling (Lee et al., 2015)	-	1.69%	7.62%

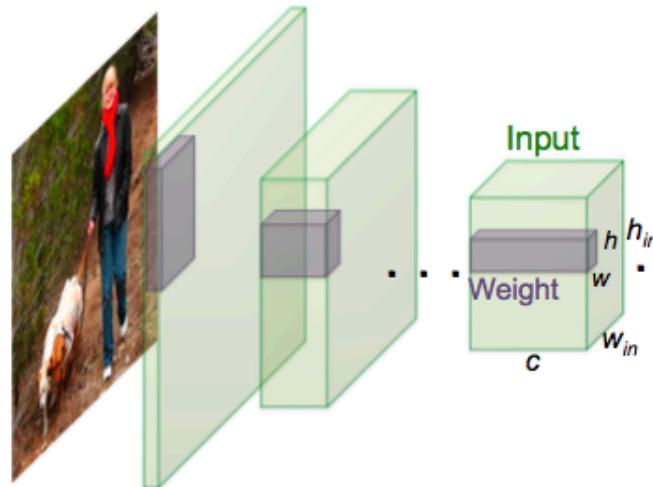
GPU KERNELS' EXECUTION TIMES



Summary

- 1. “fast”
 - 2. low memory, light network
 - 3. mobile
-
- 1. “slow”
 - 2. lower accuracy
 - 3. bottleneck : gradient still needs real high precision

XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks



	Network Variations	Operations used in Convolution	Memory Saving (Inference)	Computation Saving (Inference)	Accuracy on ImageNet (AlexNet)	
Standard Convolution	Real-Value Inputs 	Real-Value Weights 	+ , - , \times	1x	1x	%56.7
Binary Weight	Real-Value Inputs 	Binary Weights 	+ , -	~32x	~2x	%56.8
BinaryWeight Binary Input (XNOR-Net)	Binary Inputs 	Binary Weights 	XNOR , bitcount	~32x	~58x	%44.2

$$\mathbf{W} \approx \alpha \mathbf{B} \quad \alpha \in \mathbb{R}^+$$

$$\mathbf{B} \in \{+1, -1\}^{c \times w \times h}$$

$$\mathbf{I} * \mathbf{W} \approx (\mathbf{I} \oplus \mathbf{B}) \alpha \quad \oplus \text{ indicates a convolution without any multiplication}$$

$$J(\mathbf{B}, \alpha) = \|\mathbf{W} - \alpha \mathbf{B}\|^2$$

$$\alpha^*, \mathbf{B}^* = \operatorname{argmin}_{\alpha, \mathbf{B}} J(\mathbf{B}, \alpha)$$

$$J(\mathbf{B}, \alpha) = \alpha^2 \mathbf{B}^\top \mathbf{B} - 2\alpha \mathbf{W}^\top \mathbf{B} + \mathbf{W}^\top \mathbf{W}$$

$$J(\mathbf{B}, \alpha) = \alpha^2 n - 2\alpha \mathbf{W}^\top \mathbf{B} + \mathbf{c}$$

$$\mathbf{B}^* = \operatorname{argmax}_{\mathbf{B}} \{\mathbf{W}^\top \mathbf{B}\} \quad s.t. \quad \mathbf{B} \in \{+1, -1\}^n \quad \text{vector}$$

$$\mathbf{B} \in \{+1, -1\}^n,$$

$$\mathbf{B}^\top \mathbf{B} = n$$

$$\mathbf{c} = \mathbf{W}^\top \mathbf{W}$$

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

$$\alpha^* = \frac{\mathbf{W}^\top \mathbf{B}^*}{n}$$

$$\alpha^* = \frac{\mathbf{W}^\top \operatorname{sign}(\mathbf{W})}{n} = \frac{\sum |\mathbf{W}_i|}{n} = \frac{1}{n} \|\mathbf{W}\|_{\ell_1}$$

$\mathbf{X}^\top \mathbf{W} \approx \beta \mathbf{H}^\top \alpha \mathbf{B}$, where $\mathbf{H}, \mathbf{B} \in \{+1, -1\}^n$ and $\beta, \alpha \in \mathbb{R}^+$

$$\alpha^*, \mathbf{B}^*, \beta^*, \mathbf{H}^* = \operatorname{argmin}_{\alpha, \mathbf{B}, \beta, \mathbf{H}} \|\mathbf{X} \odot \mathbf{W} - \beta \alpha \mathbf{H} \odot \mathbf{B}\|$$

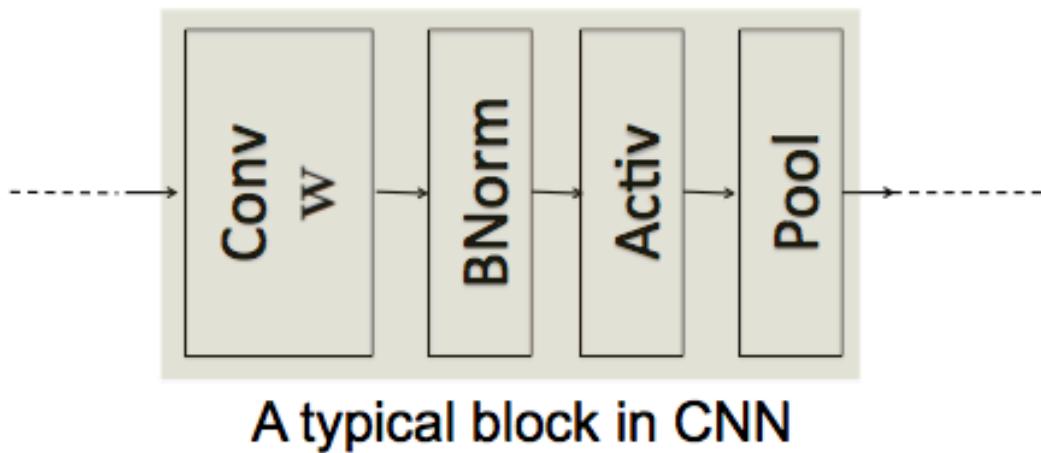
$$\gamma^*, \mathbf{C}^* = \operatorname{argmin}_{\gamma, \mathbf{C}} \|\mathbf{Y} - \gamma \mathbf{C}\|$$

$$\mathbf{C}^* = \text{sign}(\mathbf{Y}) = \text{sign}(\mathbf{X}) \odot \text{sign}(\mathbf{W}) = \mathbf{H}^* \odot \mathbf{B}^*$$

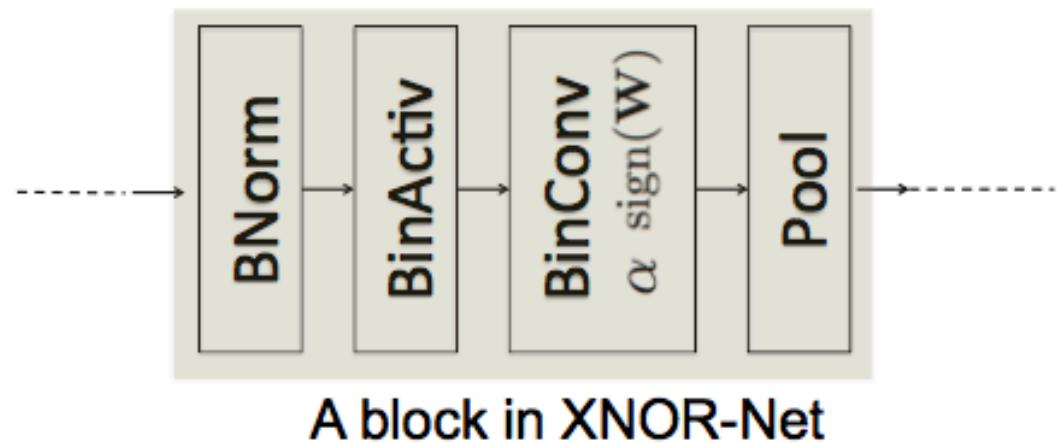
$$\gamma^* = \frac{\sum |\mathbf{Y}_i|}{n} = \frac{\sum |\mathbf{X}_i| |\mathbf{W}_i|}{n} \approx \left(\frac{1}{n} \|\mathbf{X}\|_{\ell_1} \right) \left(\frac{1}{n} \|\mathbf{W}\|_{\ell_1} \right) = \beta^* \alpha^*$$

$$\mathbf{E}[|\mathbf{Y}_i|] = \mathbf{E}[|\mathbf{X}_i| |\mathbf{W}_i|] = \mathbf{E}[|\mathbf{X}_i|] \mathbf{E}[|\mathbf{W}_i|]$$

$$\mathbf{I} * \mathbf{W} \approx (\text{sign}(\mathbf{I}) \circledast \text{sign}(\mathbf{W})) \odot \mathbf{K}\alpha$$

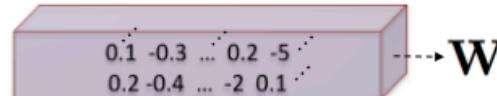


A typical block in CNN

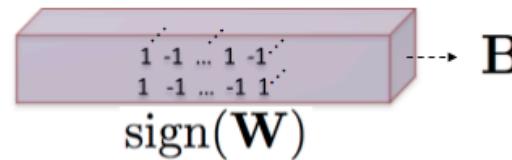


A block in XNOR-Net

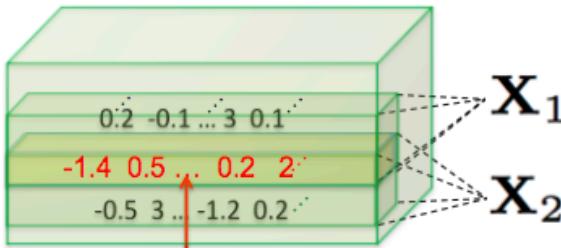
(1) Binarizing Weight



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



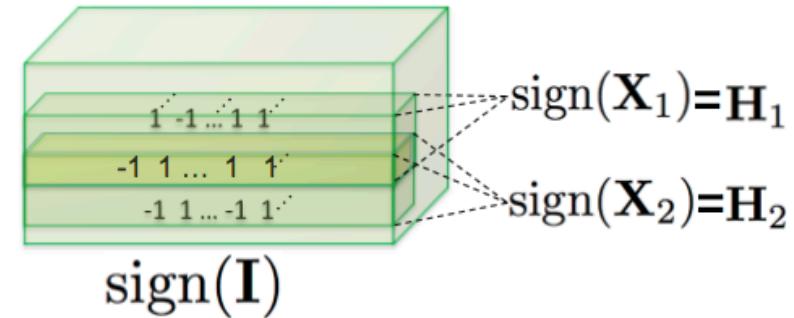
(2) Binarizing Input



Redundant computations in overlapping areas

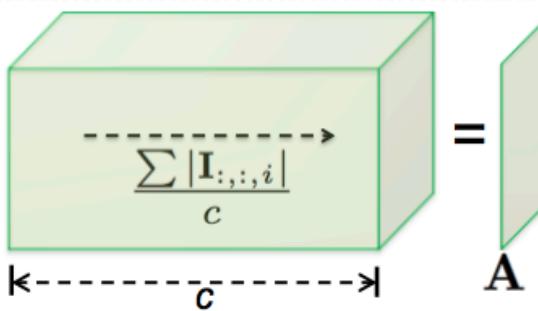
$$\frac{1}{n} \|\mathbf{X}_1\|_{\ell_1} = \beta_1$$

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

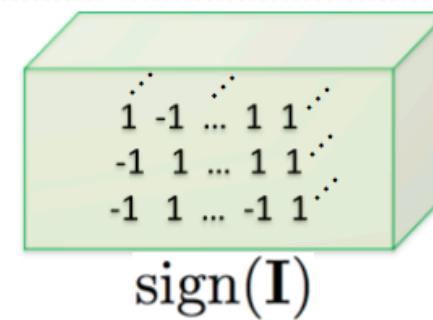


$$\text{sign}(\mathbf{I})$$

(3) Binarizing Input



$$\mathbf{A} * \underset{\mathbf{k}}{\mathbb{1}} = \beta_1$$



(4) Convolution with XNOR-Bitcount

$$\mathbf{I} * \mathbf{W} \approx \text{sign}(\mathbf{I})$$

$$\text{sign}(\mathbf{I}) * \text{sign}(\mathbf{W}) \odot \alpha \odot \mathbf{K}$$



Input

$$\bar{y}_1^{\ell-1} = \begin{bmatrix} 3 & 2 & 1 \\ 0 & -1 & -2 \\ -3 & 1 & -1 \end{bmatrix}, \quad \bar{w}_{1,1}^{\ell} = \begin{bmatrix} 1 & 1 \\ -1 & -2 \end{bmatrix}$$

$$\bar{y}_2^{\ell-1} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & -1 & -2 \\ -3 & 1 & -1 \end{bmatrix}, \quad \bar{w}_{2,1}^{\ell} = \begin{bmatrix} -1 & -1 \\ 2 & 1 \end{bmatrix}, \quad b_1^{\ell} = 1$$



Convolution computation

$$\bar{y}_1^{\ell-1} * \bar{w}_{1,1}^{\ell} = \begin{bmatrix} 7 & 8 \\ 0 & -2 \end{bmatrix}, \quad \bar{y}_2^{\ell-1} * \bar{w}_{2,1}^{\ell} = \begin{bmatrix} -3 & -9 \\ -4 & 4 \end{bmatrix}$$

$$\bar{s}_n^{\ell} = \bar{y}_1^{\ell-1} * \bar{w}_{1,1}^{\ell} + \bar{y}_2^{\ell-1} * \bar{w}_{2,1}^{\ell} + b_n^{\ell} = \begin{bmatrix} 5 & 0 \\ -3 & 3 \end{bmatrix}$$



BWN Convolution

$$\bar{y}_1^{\ell-1} = \begin{bmatrix} 3 & 2 & 1 \\ 0 & -1 & -2 \\ -3 & 1 & -1 \end{bmatrix}, \quad \bar{w}_{1,1}^{\ell} = \begin{bmatrix} 1 & 1 \\ -1 & -2 \end{bmatrix}$$

$$\bar{y}_2^{\ell-1} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & -1 & -2 \\ -3 & 1 & -1 \end{bmatrix}, \quad \bar{w}_{2,1}^{\ell} = \begin{bmatrix} -1 & -1 \\ 2 & 1 \end{bmatrix}, \quad b_1^{\ell} = 1$$

$$\alpha_1^{\ell} = \frac{\sum_{i,j,m \in V_1^{\ell}} |\bar{w}_{m,1}^{\ell}|}{\text{Ch}_1^{\ell} \times w_{\ell} \times h_{\ell}}$$

$$= \frac{1+1+1+2+1+1+1+2+1}{8} = \frac{5}{4}$$

$$\bar{B}_{1,1}^{\ell} = \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix}, \quad \bar{B}_{2,1}^{\ell} = \begin{bmatrix} -1 & -1 \\ 1 & 1 \end{bmatrix}$$

$$\bar{s}_n^{\ell} = \alpha_1^{\ell} (\bar{y}_1^{\ell-1} \otimes \bar{B}_{1,1}^{\ell} + \bar{y}_2^{\ell-1} \otimes \bar{B}_{2,1}^{\ell}) + b_1^{\ell}$$

$$= \frac{5}{4} \left(\begin{bmatrix} 6 & 4 \\ 1 & -3 \end{bmatrix} + \begin{bmatrix} -4 & -8 \\ -1 & 3 \end{bmatrix} \right) + b_1^{\ell} = \begin{bmatrix} \frac{7}{2} & \frac{-3}{2} \\ 1 & 1 \end{bmatrix}$$



0:47 / 3:18





XNOR Convolution

$$\bar{y}_{c,1}^{\ell-1} = \bar{y}_1^{\ell-1} + \bar{y}_2^{\ell-1} = \begin{bmatrix} 4 & 4 & 4 \\ 0 & -2 & -4 \\ -6 & 2 & -2 \end{bmatrix}$$

$$\bar{H}_1^{\ell-1} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & -1 \\ -1 & 1 & -1 \end{bmatrix}, \quad \bar{H}_2^{\ell-1} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & -1 \\ -1 & 1 & -1 \end{bmatrix}$$

BWN

$$\alpha_1^\ell = \frac{\sum_{i,j,m \in V_1^\ell} |\bar{w}_{m,1}^\ell|}{\text{Ch}_1^\ell \times w_\ell \times h_\ell} = \frac{1+1+1+2+1+1+1+2+1}{8} = \frac{5}{4}$$

$$\bar{B}_{1,1}^\ell = \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix}, \quad \bar{B}_{2,1}^\ell = \begin{bmatrix} -1 & -1 \\ 1 & 1 \end{bmatrix}$$

$$\bar{y}_1^{\ell-1} = \begin{bmatrix} 3 & 2 & 1 \\ 0 & -1 & -2 \\ -3 & 1 & -1 \end{bmatrix}, \quad \bar{w}_{1,1}^\ell = \begin{bmatrix} 1 & 1 \\ -1 & -2 \end{bmatrix}$$

$$\bar{y}_2^{\ell-1} = \begin{bmatrix} 1 & 2 & 3 \\ 0 & -1 & -2 \\ -3 & 1 & -1 \end{bmatrix}, \quad \bar{w}_{2,1}^\ell = \begin{bmatrix} -1 & -1 \\ 2 & 1 \end{bmatrix}, \quad b_1^\ell = 1$$

$$\beta_1^{\ell-1} = \bar{y}_{c,1}^{\ell-1} \otimes \text{ones}(h_\ell, w_\ell) \times \frac{1}{\text{Ch}_1^\ell \times h_\ell \times w_\ell}$$

$$= \frac{1}{2 \times 2 \times 2} \begin{bmatrix} 6 & 2 \\ -6 & -6 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 3 & 1 \\ -3 & -3 \end{bmatrix}$$

$$\bar{s}_n^\ell = \alpha_1^\ell \times \beta_1^{\ell-1} (\bar{H}_1^{\ell-1} \otimes \bar{B}_{1,1}^\ell + \bar{H}_2^{\ell-1} \otimes \bar{B}_{2,1}^\ell) + b_1^\ell$$

$$= \frac{5}{4} \times \frac{1}{4} \begin{bmatrix} 3 & 1 \\ -3 & -3 \end{bmatrix} \odot \left(\begin{bmatrix} 2 & 4 \\ 0 & -2 \end{bmatrix} + \begin{bmatrix} -2 & -4 \\ 0 & 2 \end{bmatrix} \right) + 1$$

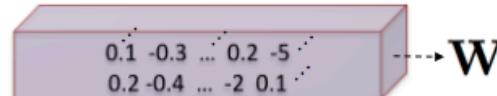
$$= \frac{5}{16} \begin{bmatrix} 3 & 1 \\ -3 & -3 \end{bmatrix} \odot \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} + 1 = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$



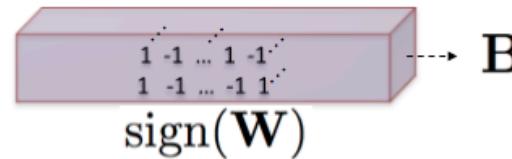
1:53 / 3:18



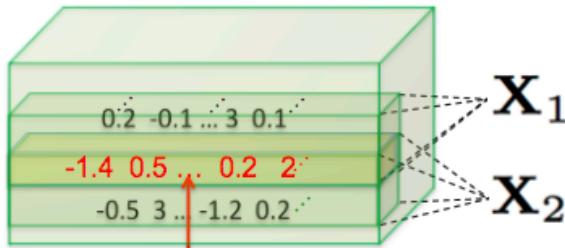
(1) Binarizing Weight



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

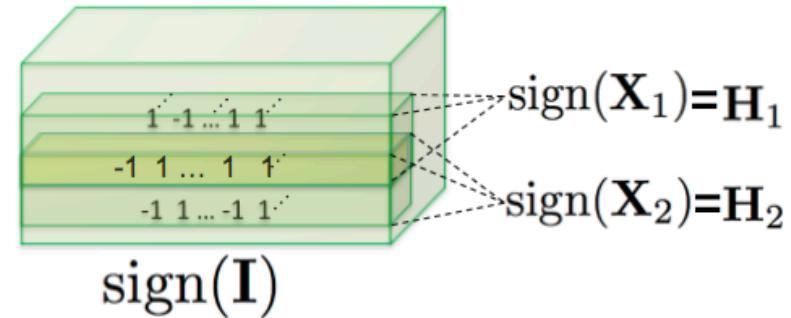


(2) Binarizing Input

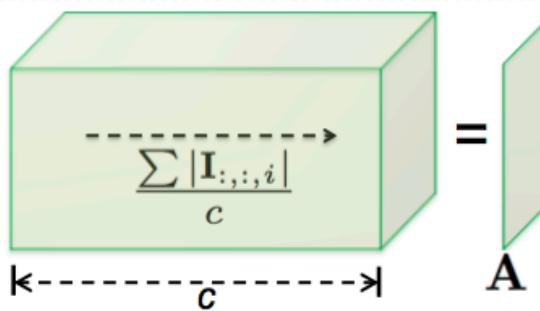


$$\frac{1}{n} \|\mathbf{X}_1\|_{\ell_1} = \beta_1$$

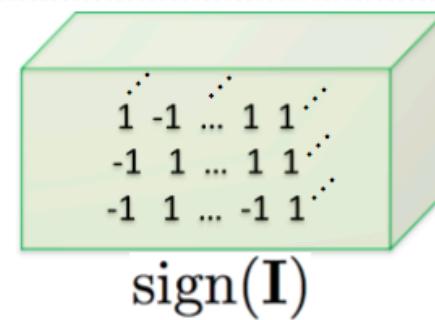
$$\frac{1}{n} \|\mathbf{X}_2\|_{\ell_1} = \beta_2$$



(3) Binarizing Input



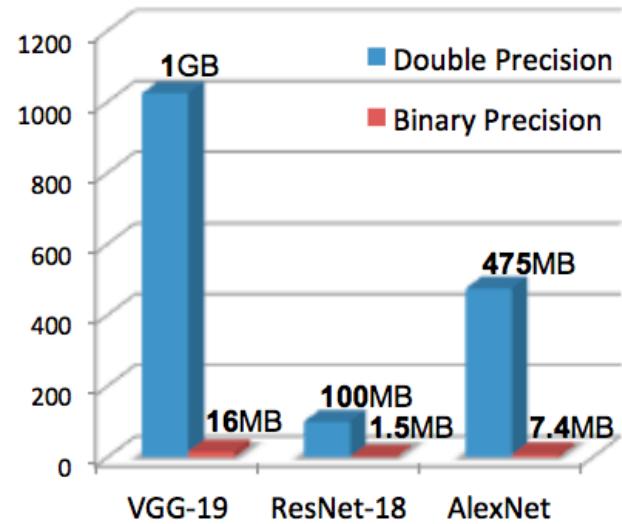
$$\mathbf{A} * \underset{k}{\mathbb{K}} = \beta_1$$
$$\beta_2$$



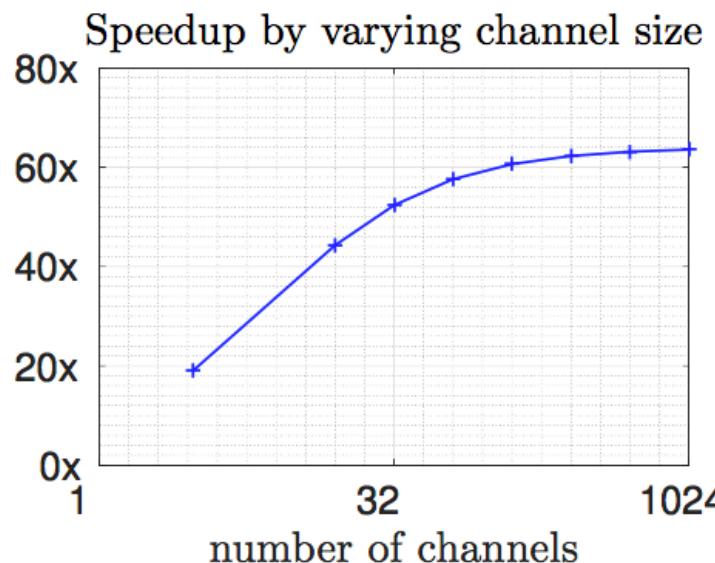
(4) Convolution with XNOR-Bitcount

$$\approx \left[\text{sign}(\mathbf{I}) * \text{sign}(\mathbf{W}) \right] * \mathbb{K} \odot \alpha$$

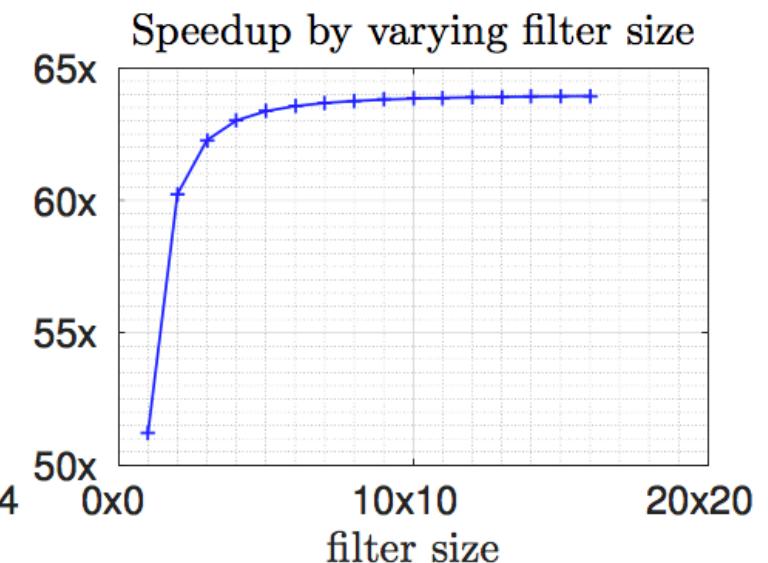
experiments



(a)



(b)



(c)

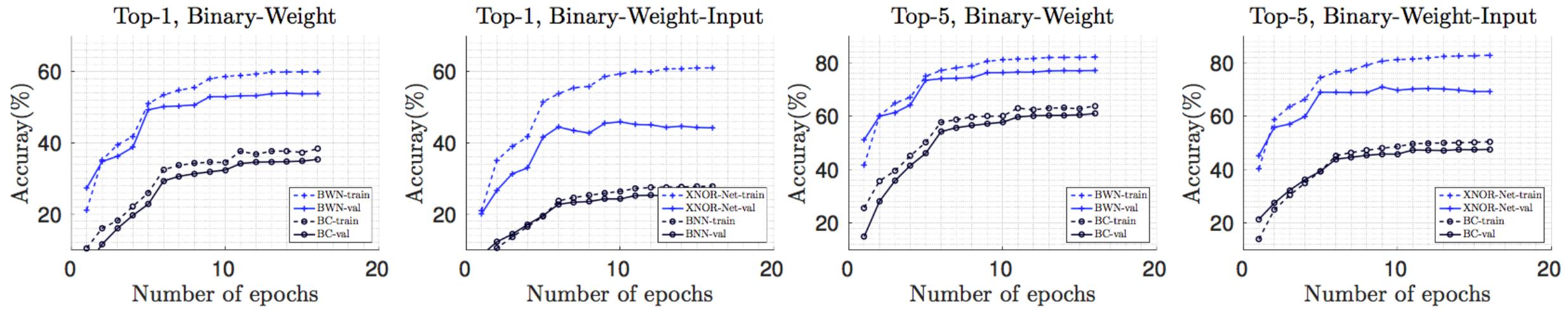
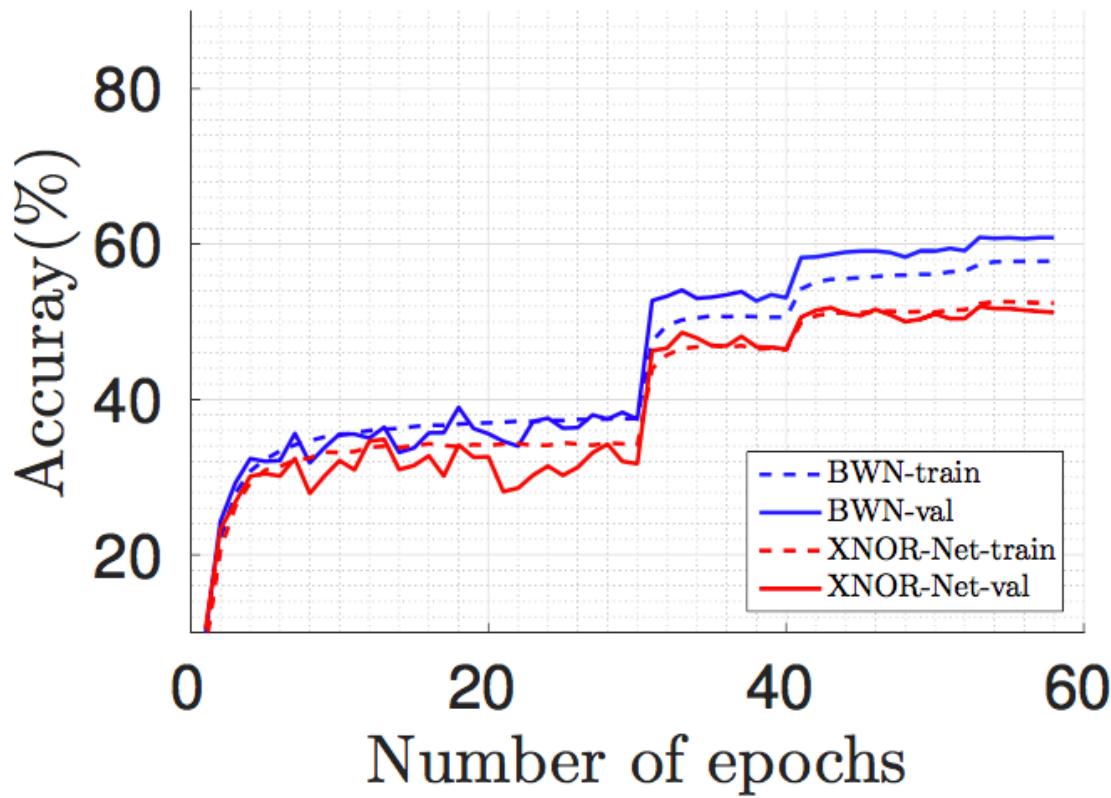


Fig. 5: This figure compares the imagenet classification accuracy on Top-1 and Top-5 across training epochs. Our approaches BWN and XNOR-Net outperform BinaryConnect(BC) and BinaryNet(BNN) in all the epochs by large margin($\sim 17\%$).

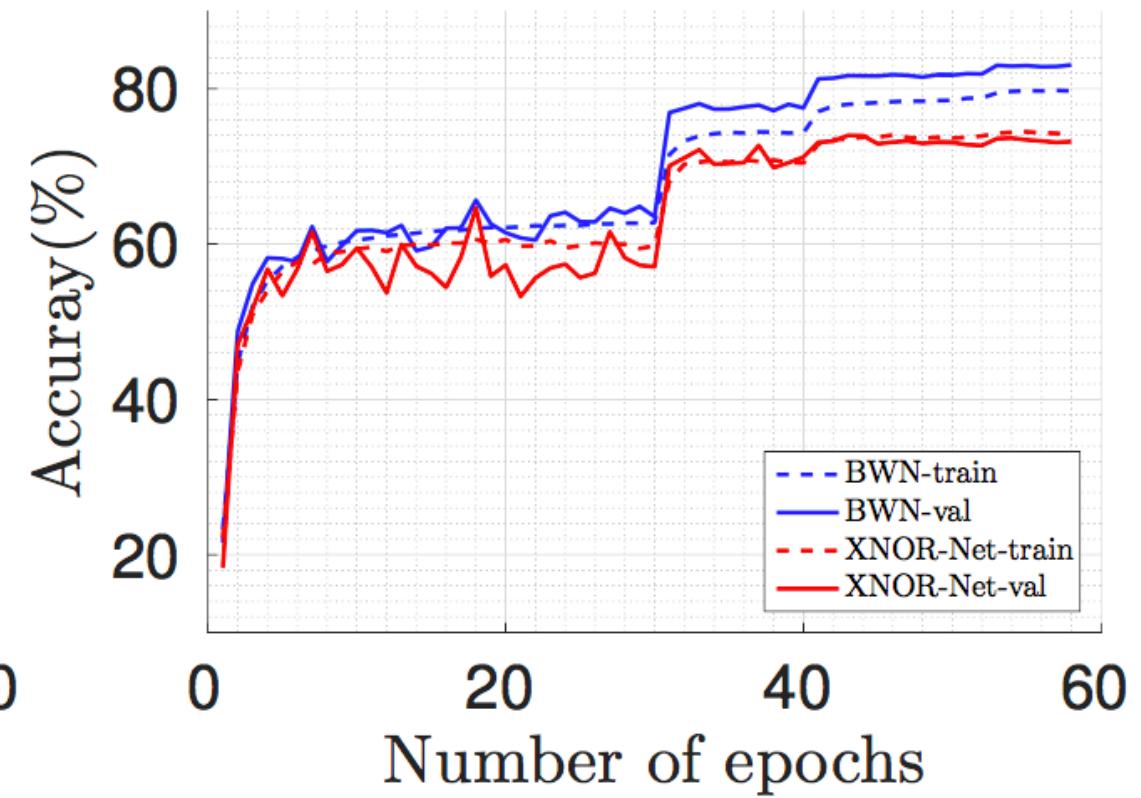
Classification Accuracy(%)											
Binary-Weight				Binary-Input-Binary-Weight				Full-Precision			
BWN		BC[11]		XNOR-Net		BNN[11]		AlexNet[1]			
Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
56.8	79.4	35.4	61.0	44.2	69.2	27.9	50.42	56.6	80.2		

ResNet, Top-1



(a)

ResNet, Top-5



(b)

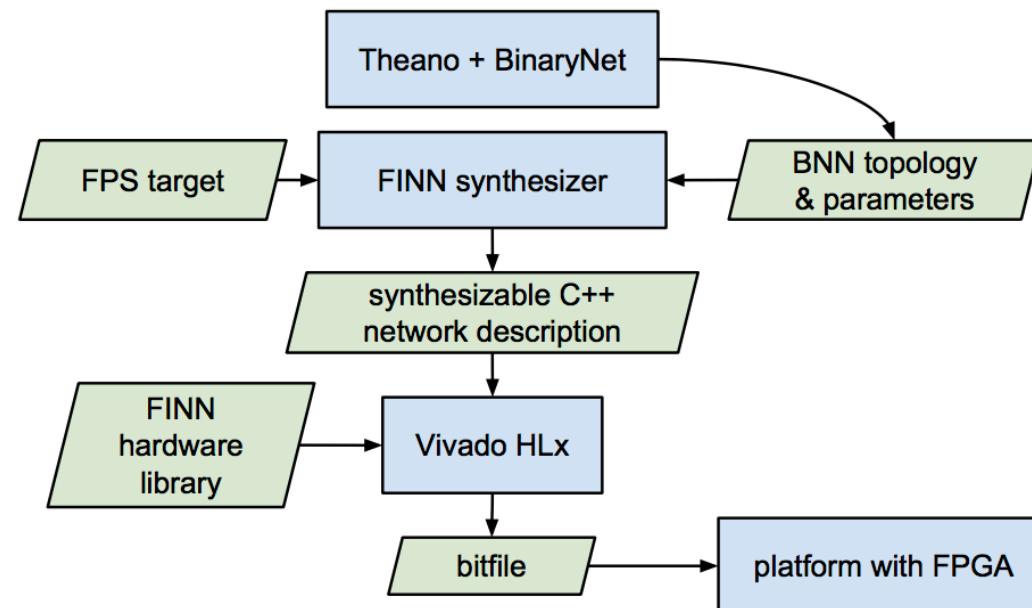
	ResNet-18		GoogLenet	
Network Variations	top-1	top-5	top-1	top-5
Binary-Weight-Network	60.8	83.0	65.5	86.1
XNOR-Network	51.2	73.2	N/A	N/A
Full-Precision-Network	69.3	89.2	71.3	90.0

Binary-Weight-Network			XNOR-Network		
Strategy for computing α	top-1	top-5	Block Structure	top-1	top-5
Using equation 6	56.8	79.4	C-B-A-P	30.3	57.5
Using a separate layer	46.2	69.5	B-A-C-P	44.2	69.2

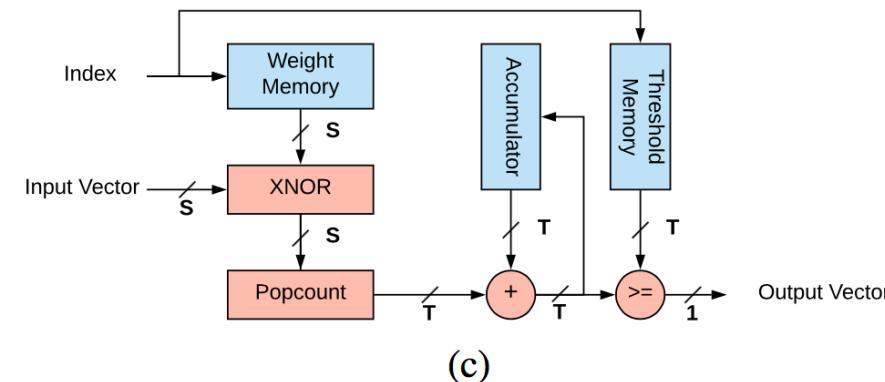
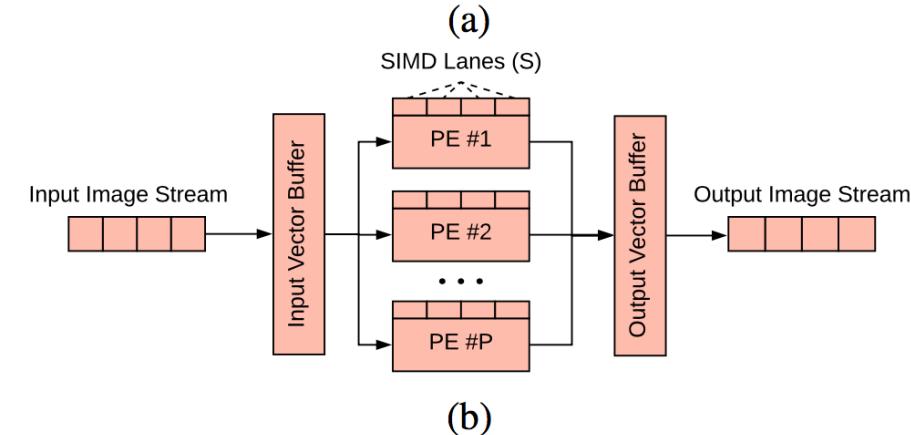
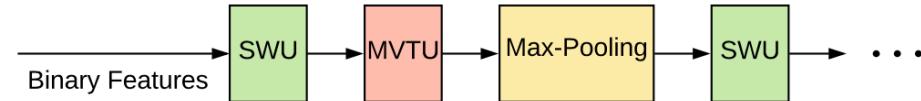
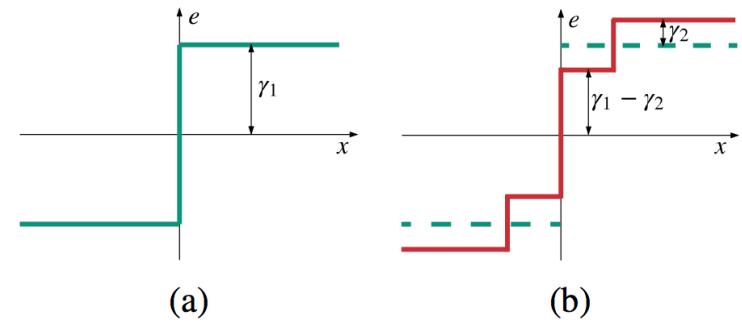
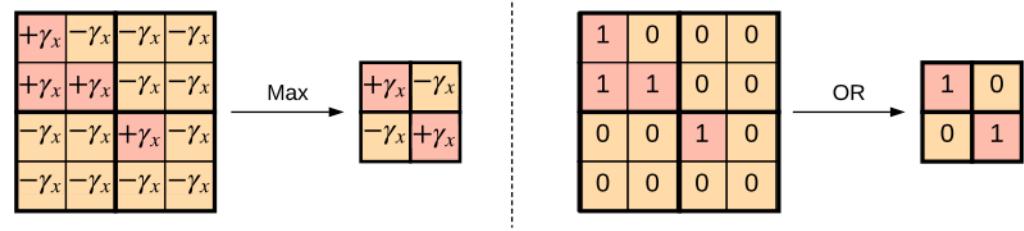
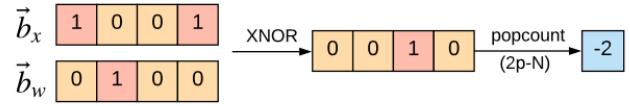
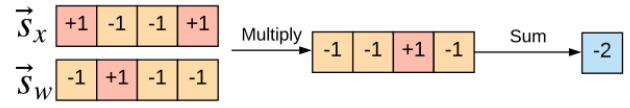
FINN: A Framework for Fast, Scalable Binarized Neural Network Inference

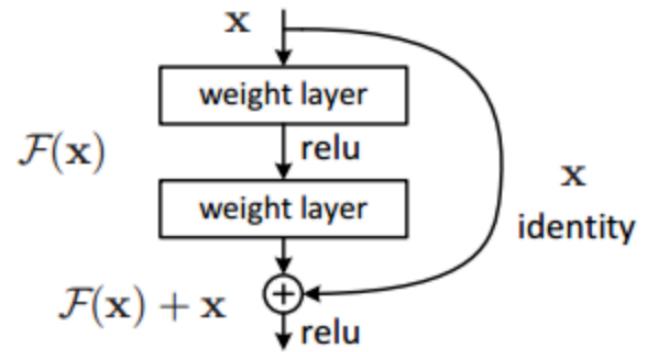
Yaman Umuroglu^{*†}, Nicholas J. Fraser^{*‡}, Giulio Gambardella*, Michaela Blott*,
Philip Leong[‡], Magnus Jähre[†] and Kees Vissers*

*Xilinx Research Labs; [†]Norwegian University of Science and Technology; [‡]University of Sydney
yamanu@idi.ntnu.no



ReBNet: Residual Binarized Neural Network





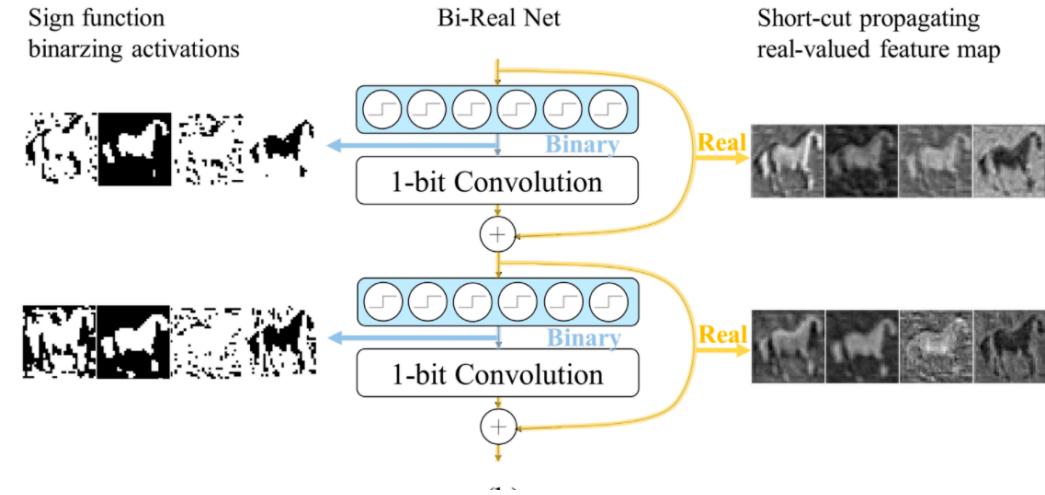
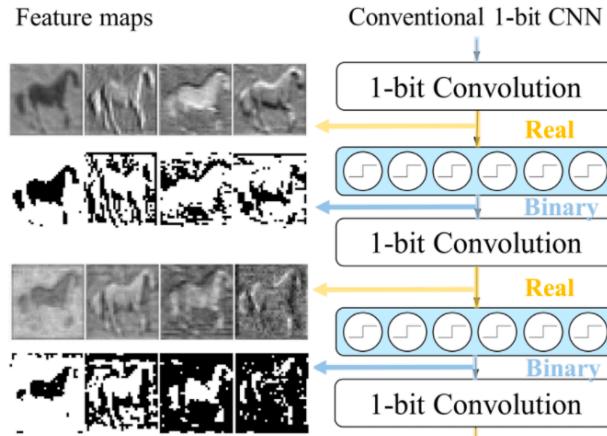
$$F'(5)=5.1 \quad H(5)=5.1 \quad H(5)=F(5)+5, F(5)=0.1$$

$$1/51=2\% \quad (0.2-0.1)/0.1 = 1$$

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \{W_i\}) + \mathbf{x}.$$

			MNIST	CIFAR10	SVHN	Imagenet
FINN		Accuracy (%)	95.83	80.1	94.9	27.9
Ours	Throughput (samples/sec)	6.3×10^5	6×10^3	6×10^3	5.2×10^2	
	M=1	Accuracy (%)	97.87	80.59	95.46	40.89
		Throughput (samples/sec)	6.4×10^5	6×10^3	6×10^3	5.3×10^2
	M=2	Accuracy (%)	98.29	85.94	96.82	41.37
		Throughput (samples/sec)	3.3×10^5	3×10^3	3×10^3	2.6×10^2
	M=3	Accuracy (%)	98.25	86.98	97.00	41.43
		Throughput (samples/sec)	2.3×10^5	2×10^3	2×10^3	1.7×10^2

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



	Map size	Value range	Different value	\mathbb{R}
A_b^l	14x14x32	{-1, 1}	2	$2^{14 \times 14 \times 32}$
A_m^l	14x14x32	{-288, -286, ..., 288}	289	$289^{14 \times 14 \times 32}$
A_r^{l+1}	14x14x32	{a certain range}	289	$289^{14 \times 14 \times 32}$
A_b^{l+1}	14x14x32	{-1, 1}	2	$2^{14 \times 14 \times 32}$

(a)

	Map size	Value range	Different value	\mathbb{R}
A_r^l	14x14x32	{a certain range}	289	$289^{14 \times 14 \times 32}$
A_b^l	14x14x32	{-1, 1}	2	$2^{14 \times 14 \times 32}$
A_m^l	14x14x32	{-288, -286, ..., 288}	289	$289^{14 \times 14 \times 32}$
A_r^{l+1}	14x14x32	{a certain range}	289	$289^{14 \times 14 \times 32}$
A_{add}^{l+1}	14x14x32	{a certain range}	$289 * 289$	$289^{2 \times 14 \times 14 \times 32}$

(b)

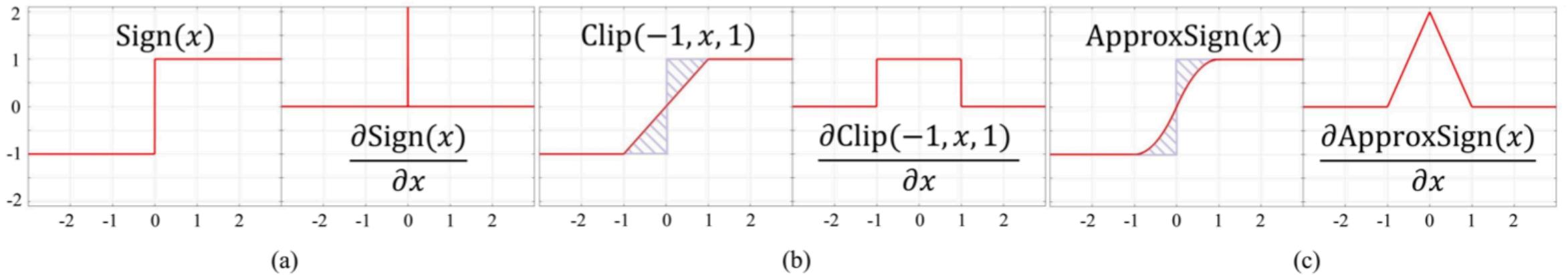


Fig. 5. (a) Sign function and its derivative, (b) Clip function and its derivative for approximating the derivative of the sign function, proposed in [7], (c) Proposed differentiable piecewise polynomial function and its triangle-shaped derivative for approximating the derivative of the sign function in gradients computation.

$$\mathbf{x} \cdot \mathbf{y} = N - 2 \times \text{bitcount}(\text{xnor}(\mathbf{x}, \mathbf{y})), x_i, y_i \in \{-1, 1\} \forall i.$$

DoReFa-NET: TRAINING LOW BITWIDTH CONVOLUTIONAL NEURAL NETWORKS WITH LOW BITWIDTH GRADIENTS

Binarized LSTM Language Model