



# XNOR-NET: IMAGENET CLASSIFICATION USING BINARY CONVOLUTIONAL NEURAL NETWORKS

RYAN MCCAMPBELL

2019 SPRING @ [HTTPS://QDATA.GITHUB.IO/DEEP2READ](https://qdata.github.io/deep2read)

# PROBLEM

- Ordinary neural networks are expensive to store and evaluate
- Not good for embedded and mobile platforms
- How to make them more efficient?

# APPROACHES

- More compact networks
- Compressing pre-trained networks
- Quantizing parameters
- Binarizing parameters/activations

# BINARY-WEIGHT-NETWORK

- Idea: store weights as binary vectors,  $\pm 1$ , plus scale

$$\mathbf{W} \approx \alpha \mathbf{B}, \alpha > 0, \mathbf{B} \in \{-1, 1\}^{\uparrow n}$$

$$\mathbf{I} * \mathbf{W} \approx (\mathbf{I} \oplus \mathbf{B}) \alpha$$

- $\oplus$ : convolution using only addition/subtraction
- Eliminates most multiplications
- $\sim 32x$  storage reduction

# BINARY-WEIGHT-NETWORK

- Minimize  $J(B, \alpha) = \|W - \alpha B\|_2^2$

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

$$\alpha^* = \frac{\sum |W_{ij}|}{n} = \frac{1}{n} \|W\|_1$$

# TRAINING

- Binarize weights during forward pass and backwards pass
- Use full-parameter weights for update
- Approximate  $\partial/\partial x \text{sign}(x) = 1[|r| \leq 1]$

# XNOR-NET

- Binarize both weights and inputs
- Convolutions using efficient binary operations: shift, XNOR & bit-count

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

- $K \downarrow ij = \beta = 1/n \|\text{subtensor of } X \text{ at } ij\| \downarrow 1$
- $K = A * k, A = \sum |I \downarrow :, :, i| / c, k \downarrow ij = 1 / wh$

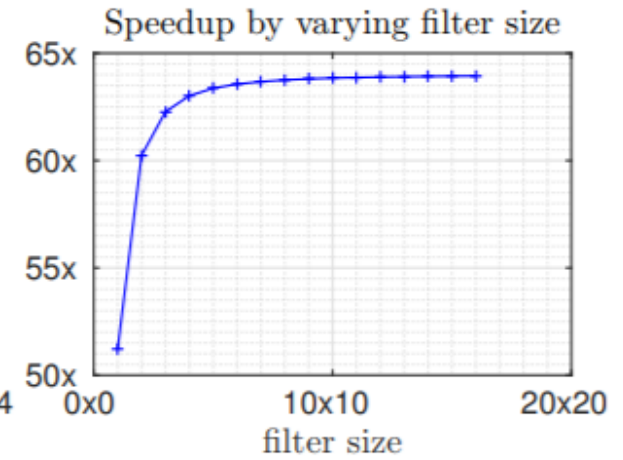
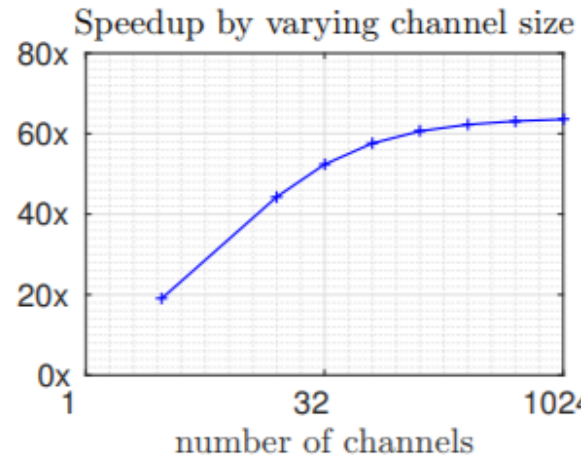
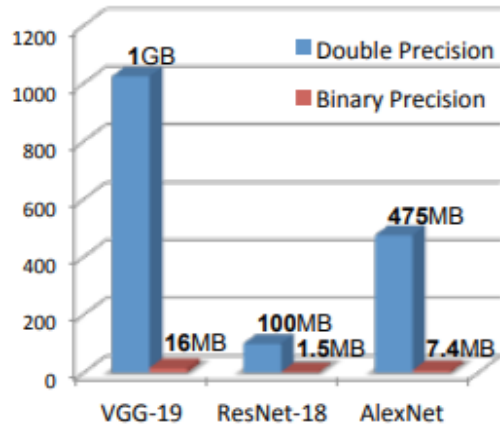
# TRAINING XNOR-NET

- Batch normalization first
- Binary activation: compute  $K$  and  $\text{sign}(I)$
- Binary convolution
- Pool after convolution





# RESULTS



# RESULTS

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
<b>56.8</b>	<b>79.4</b>	35.4	61.0	<b>44.2</b>	<b>69.2</b>	27.9	50.42	56.6	80.2