

PopBin: Popcount Binarization for Lightweight Binary Neural Networks

Hyungdong Park, Inguk Yeo
Department of Computer Engineering

Contents

1. Baseline Model

- ReActNet-18

2. Implementing binarized counting in hardware architectures

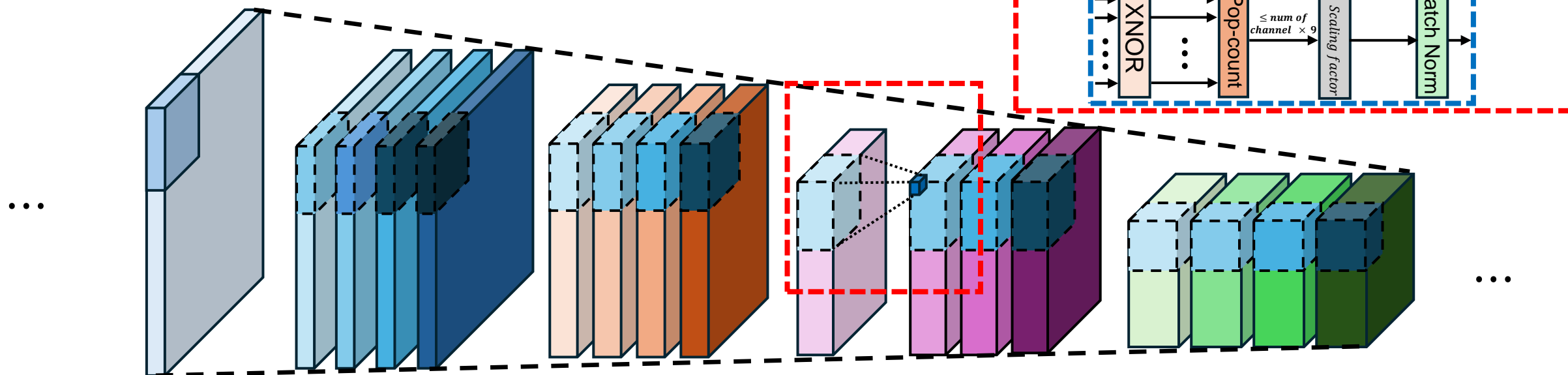
- Architectures and Benefits of popcount binarization



3. Popcount binarization techniques

- PTQ-popcount binarization
- Simple QAT-popcount binarization
- QAT-popcount binarization

Progress

- ReActNet-18 with CIFAR-10 using xnor & popcount

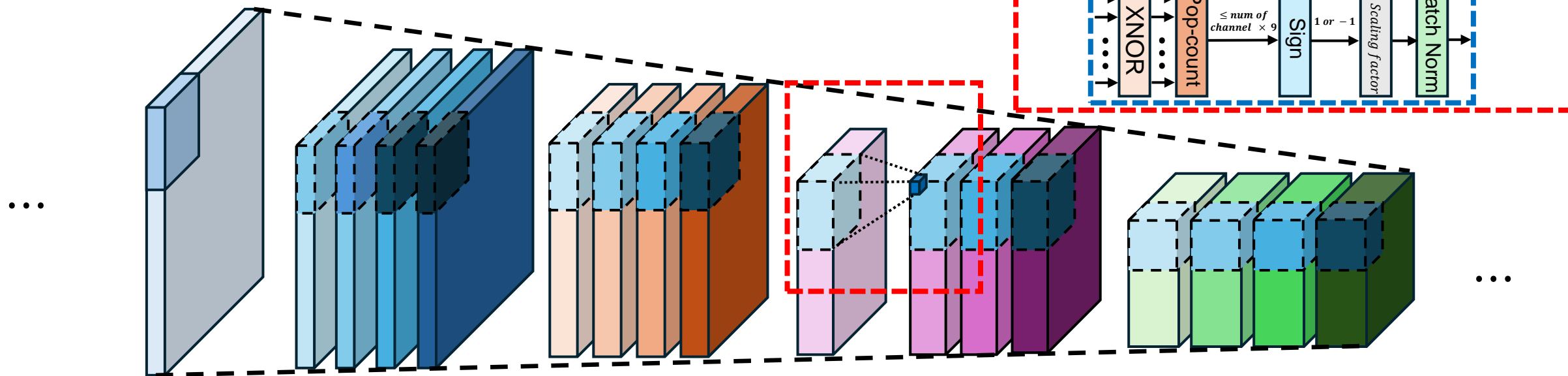




Num of channels	64	64	128	256	512	Pooling & FC
Num of layers	1	4	4	4	4	
Image size	$64 \times 32 \times 32$	$64 \times 32 \times 32$	$128 \times 16 \times 16$	$256 \times 8 \times 8$	$512 \times 4 \times 4$	
Operations	\otimes	XNOR & Pop-Count & Multiplication & Batch Norm				
Activations and weights	\mathbb{R}	\mathbb{R} (After BN) & Binarized values (1 or -1)				\mathbb{R}
# units of Mul		$c_o \times h_o \times w_o$				

Implementing popcount binarization on baseline model

Progress

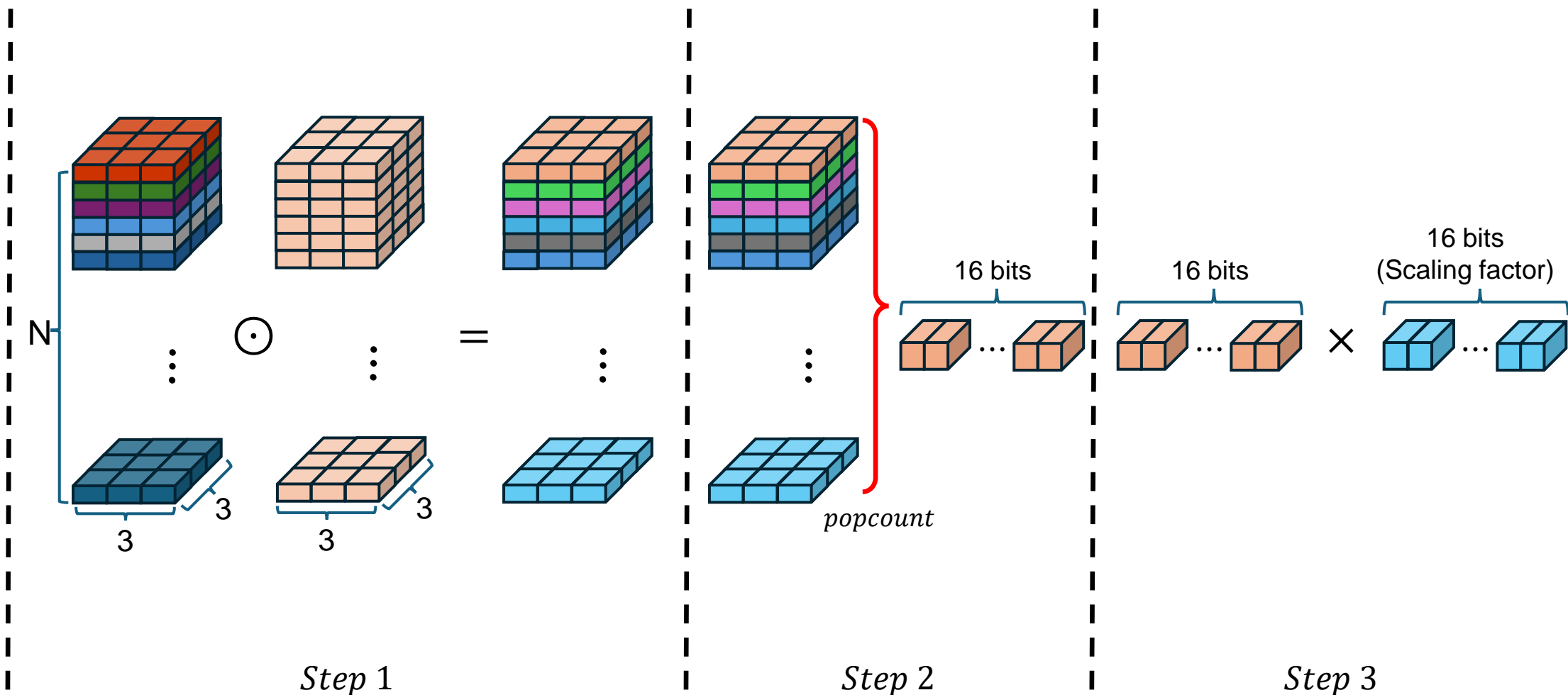
- Software Implementation of PopBin Network with CIFAR-10



Num of channels	64	64	128	256	512	Pooling & FC
Num of layers	1	4	4	4	4	
Image size	$64 \times 32 \times 32$	$64 \times 32 \times 32$	$128 \times 16 \times 16$	$256 \times 8 \times 8$	$512 \times 4 \times 4$	
Operations	\otimes	XNOR & Pop-Count & Multiplication & Batch Norm				
Activations and weights	\mathbb{R}	\mathbb{R} (After BN) & Binarized values (1 or -1)				\mathbb{R}
# units of Mul		None				

Progress

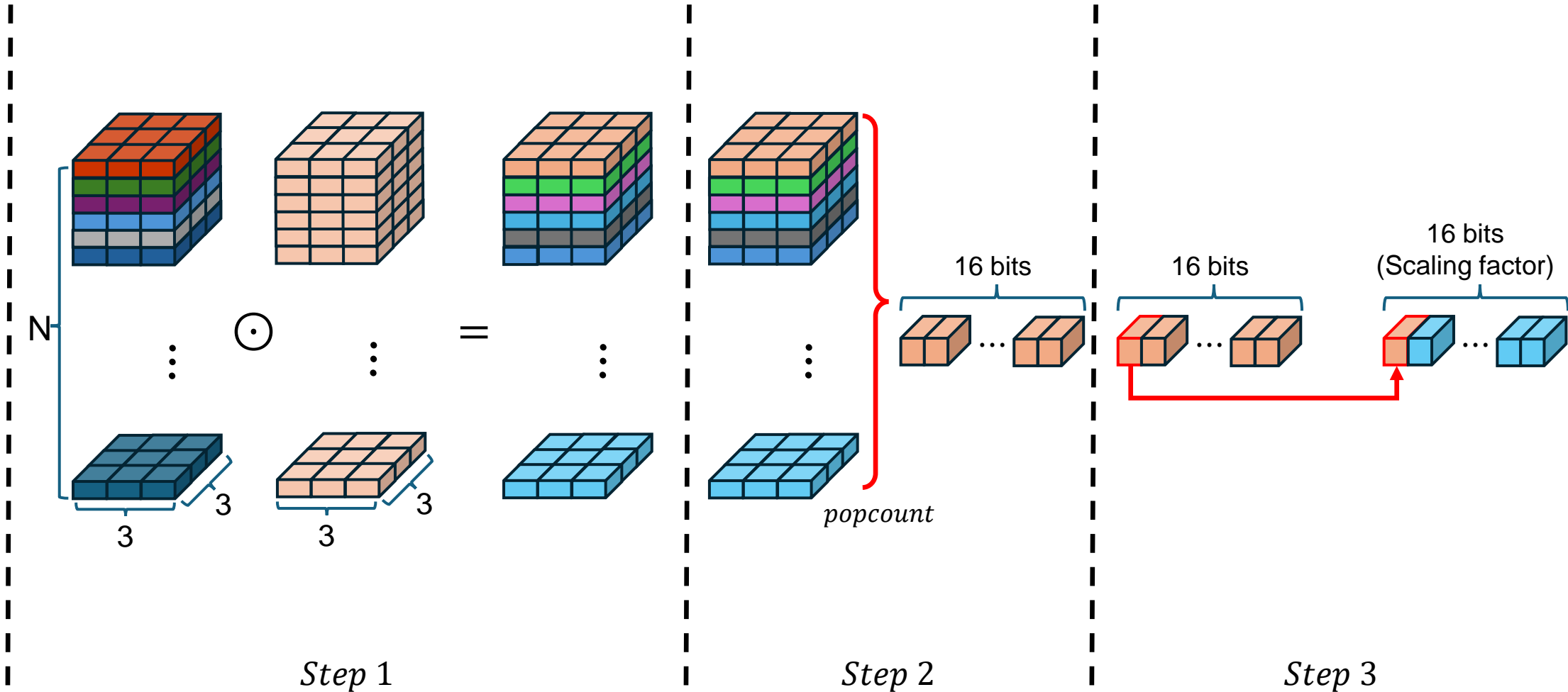
- Hardware operation process per a xnor & popcount for ReActNet-18 (our baseline)



Operations	$Xnor (\odot)$	$Popcount$	$Integer\ multiplication\ \&\ Bit\ shift$
# operations	$N \times 3^2$	$(N \times 3^2) - 1$	1 $Integer\ multiplication\ \&\ 1\ Bit\ shift$

Progress

- Hardware operation process per a xnor & popcount for our model

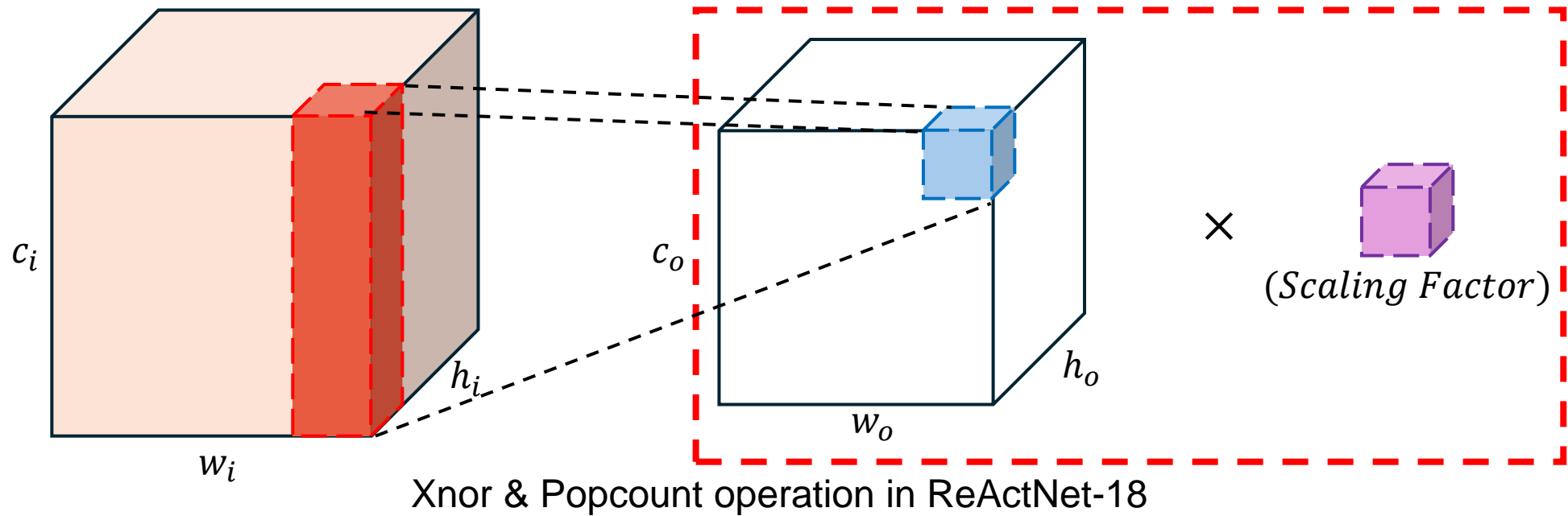


Operations	$Xnor (\odot)$	$Popcount$	$None$
# operations	$N \times 3^2$	$(N \times 3^2) - 1$	$None$

Benefits of popcount binarization

Progress

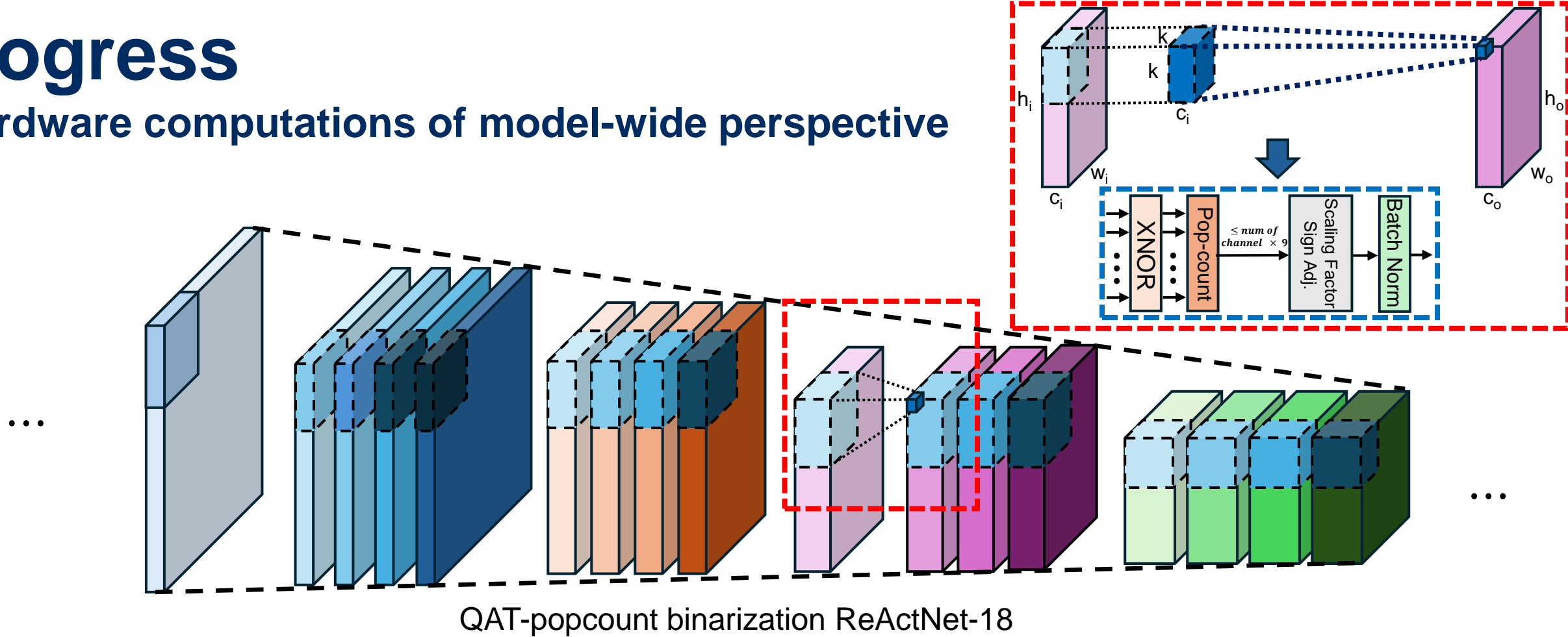
- Hardware computations per a layer



Models	Operations	# Operations per a layer
ReActNet-18	Integer Multiplication & Bit shift	$c_o \times w_o \times h_o$ Integer Multiplications & Bit shifts
QAT-popcount binarization ReActNet-18	None	None

Progress

- Hardware computations of model-wide perspective

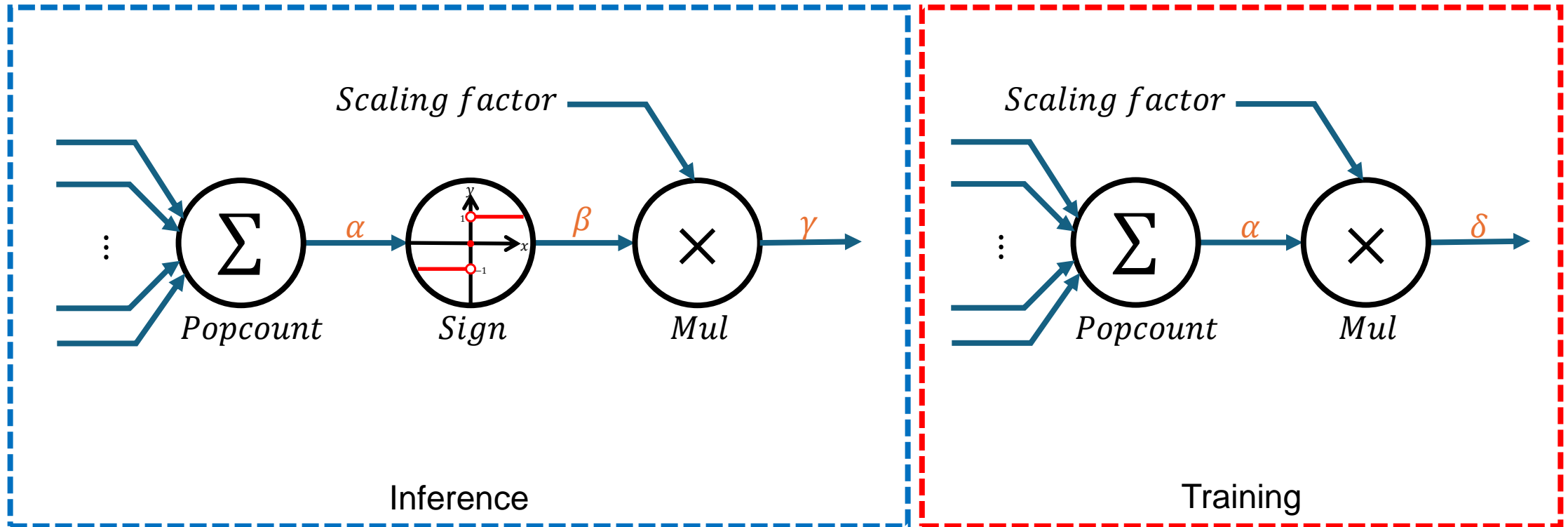


Models	Operations	# Operations
ReActNet-18	Integer Multiplication & Bit shift	557,056 Integer Multiplications & Bit shifts
QAT-popcount binarization ReActNet-18	None	None

Popcount binarization techniques

Progress

- Structure for PTQ-popcount binarization



$$|\alpha| \leq (\text{channel_num} \times \text{kernel}^2)$$

$$\beta = \pm 1$$

$$\gamma = \pm \text{scaling factor}$$

$$\delta = \pm (\text{scaling factor} \times \text{channel_num} \times \text{kernel}^2)$$

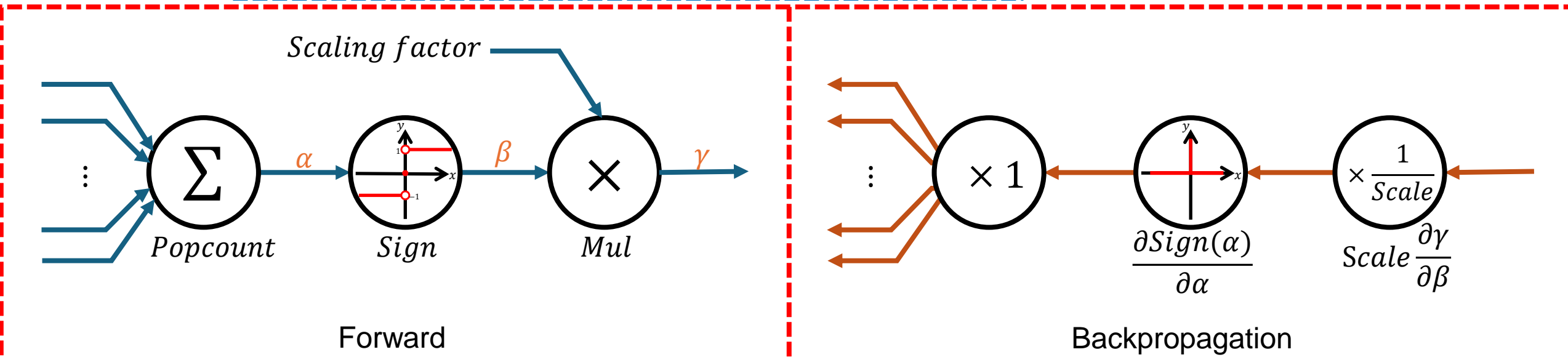
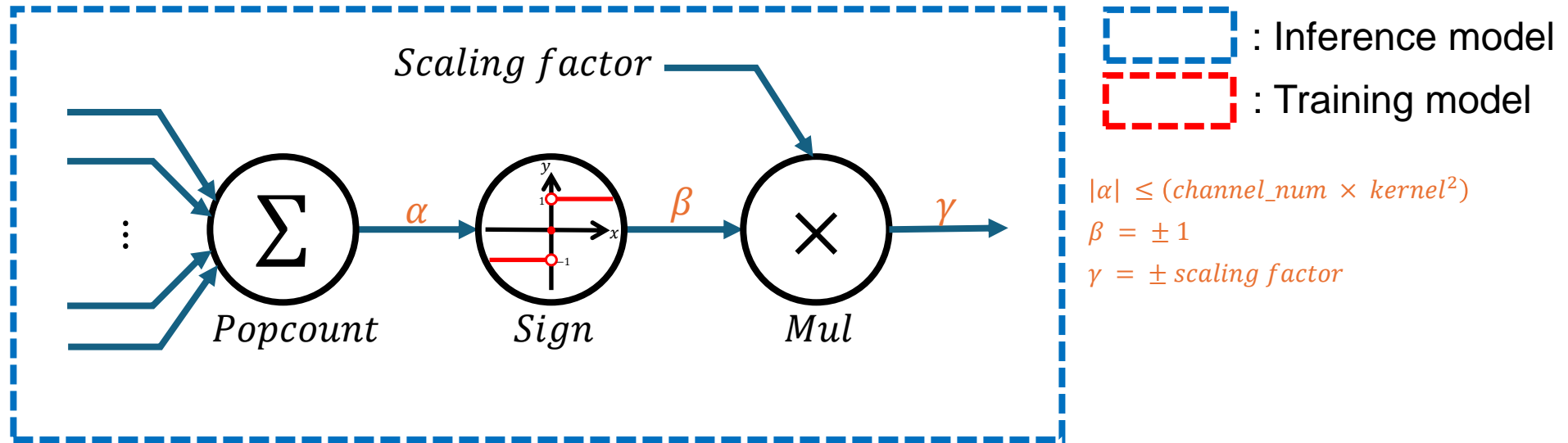
Progress

- PTQ-popcount binarization's results with CIFAR-10

Models	Top-1 Accuracy (%)	Top-5 Accuracy (%)
ReActNet-18	93.380	99.800
PTQ-popcount binarization ReActNet-18	10.000	52.040
Bi-RealNet-18	88.770	98.250
PTQ-popcount binarization Bi-RealNet-18	10.000	50.000

Progress

- Structure for Simple QAT-popcount binarization



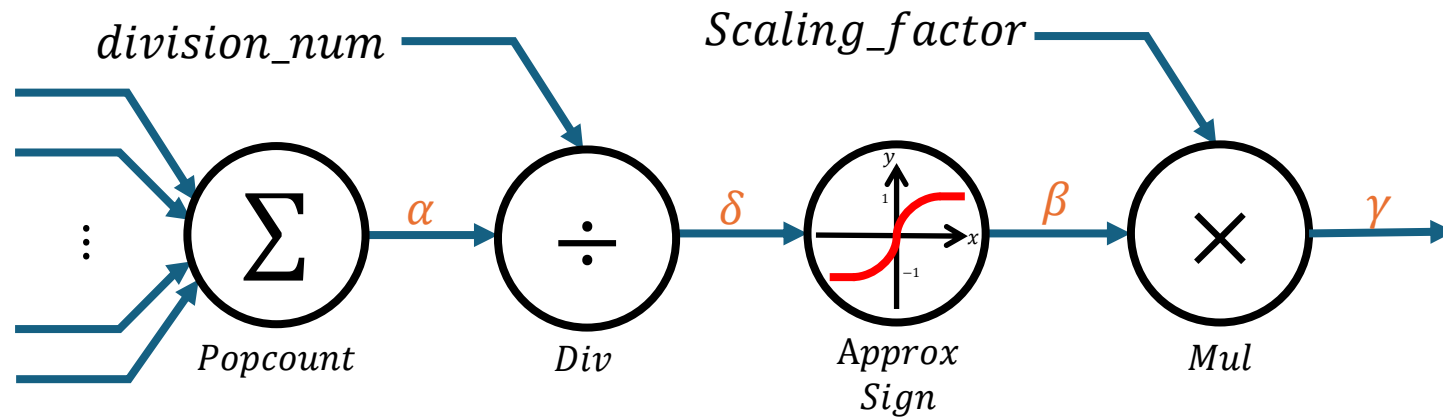
Progress

- Simple QAT-popcount binarization's results with CIFAR-10

Models	Top-1 Accuracy (%)	Top-5 Accuracy (%)
ReActNet-18	93.380	99.800
Simple QAT-popcount binarization ReActNet-18	84.930	99.250
Bi-RealNet-18	88.770	98.250
Simple QAT-popcount binarization Bi-RealNet-18	30.070	79.690

Progress

- Structure for QAT-popcount binarization



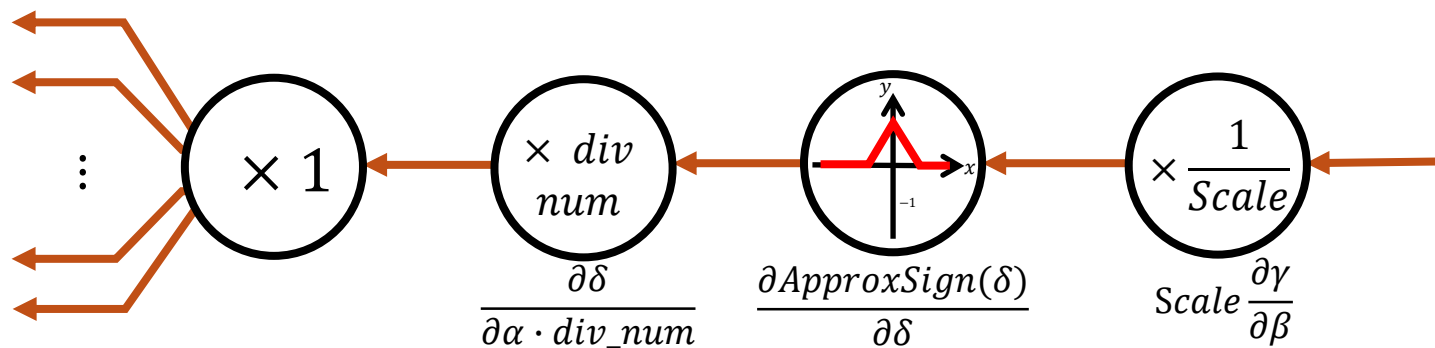
Forward in training

$$|\alpha| \leq (\text{channel_num} \times \text{kernel}^2)$$

$$\beta = \pm 1$$

$$\gamma = \pm \text{scaling_factor}$$

$$\delta = \pm (\text{channel_num} \times \text{kernel}^2 \div \text{division_num})$$



Backpropagation in training

Progress

- QAT-popcount binarizationReActNet-18's results with CIFAR-10 along with division num

Division num	Top-1 Accuracy (%)	Top-5 Accuracy (%)
$channel\ num + \alpha$	92.150	99.640
$(channel\ num \times kernel^2) + \alpha$	89.580	99.460
$channel\ num \times \alpha$	92.510	99.640
$(channel\ num \times kernel^2) \times \alpha$	92.160	99.660
Min-Max Normalization $(channel\ num \times kernel^2)$	89.230	99.390

Progress

- QAT-popcount binarization's results with CIFAR-10

Models	Top-1 Accuracy (%)	Top-5 Accuracy (%)
ReActNet-18	93.380	99.800
Simple QAT-popcount binarization ResNet-18	84.930	99.250
QAT-popcount binarization ReActNet-18 (PopBin)	92.510	99.640
Bi-RealNet-18	88.770	98.250
Simple QAT-popcount binarization Bi-RealNet-18	30.070	79.690
QAT-popcount binarization Bi-RealNet-18 (PopBin)	87.660	98.720

Thank you