

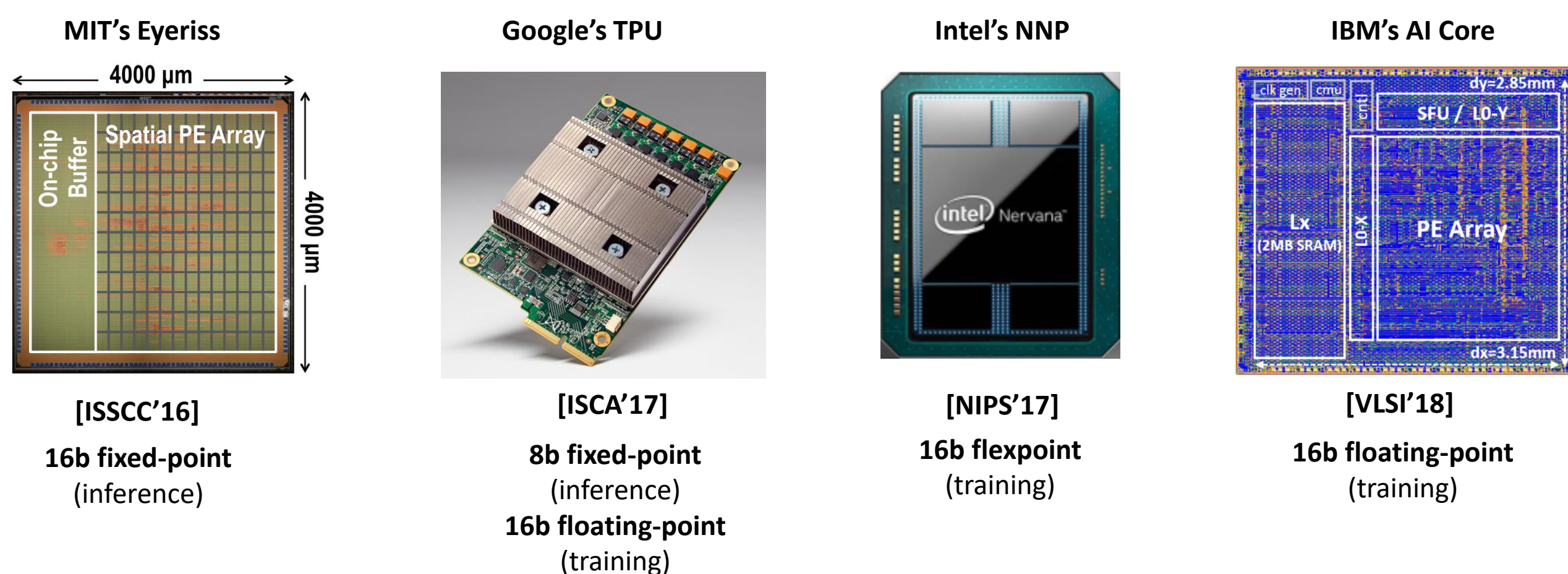
# Per-Tensor Fixed-Point Quantization of the Back-Propagation Algorithm



Charbel Sakr & Naresh Shanbhag  
University of Illinois at Urbana-Champaign  
{sakr2,shanbhag}@illinois.edu

## Motivation

### Machine Learning in Reduced Precision



Are these the minimum precisions required?  
Can minimum precision requirements be determined analytically?  
Specifically for training

## Current Approaches

### Reducing Complexity of Inference

**Forward path quantization:**

- Fixed-point via SQNR analysis [Lin et al., ICML'17]
- Extreme quantization, e.g., BinaryNet, via training [Courbariaux et al., NIPS'15 – Rastegari et al., ECCV'16 – Hubara et al., NIPS'16]

**Structural methods:**

- Pruning [Han et al., NIPS'15]
- Parameter clustering [Wu et al., ICML'18]

No theoretical guarantees on accuracy

### Reducing Complexity of Training

**Quantized back-propagation:**

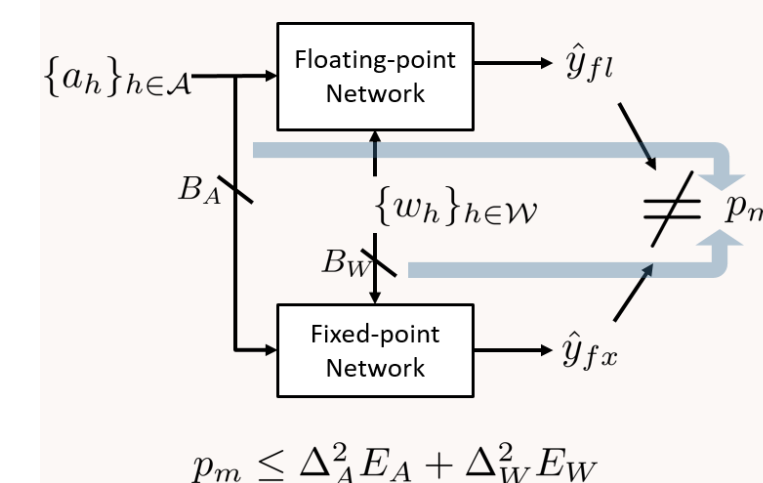
- Fixed-point via stochastic rounding [Gupta et al., ICML'15]
- Fixed-point/floating-point hybrid [Koester et al., NIPS'17]
- Finite precision floating-point [Wang et al., NIPS'18]

**Gradient Compression:**

- Extreme gradient quantization, e.g., TernGrad [Wen et al., NIPS'17]
- Gradient sparsification [Lin et al., ICLR'18]

Hybrid fixed/floating-point training

### Fixed-point inference with theoretical guarantees

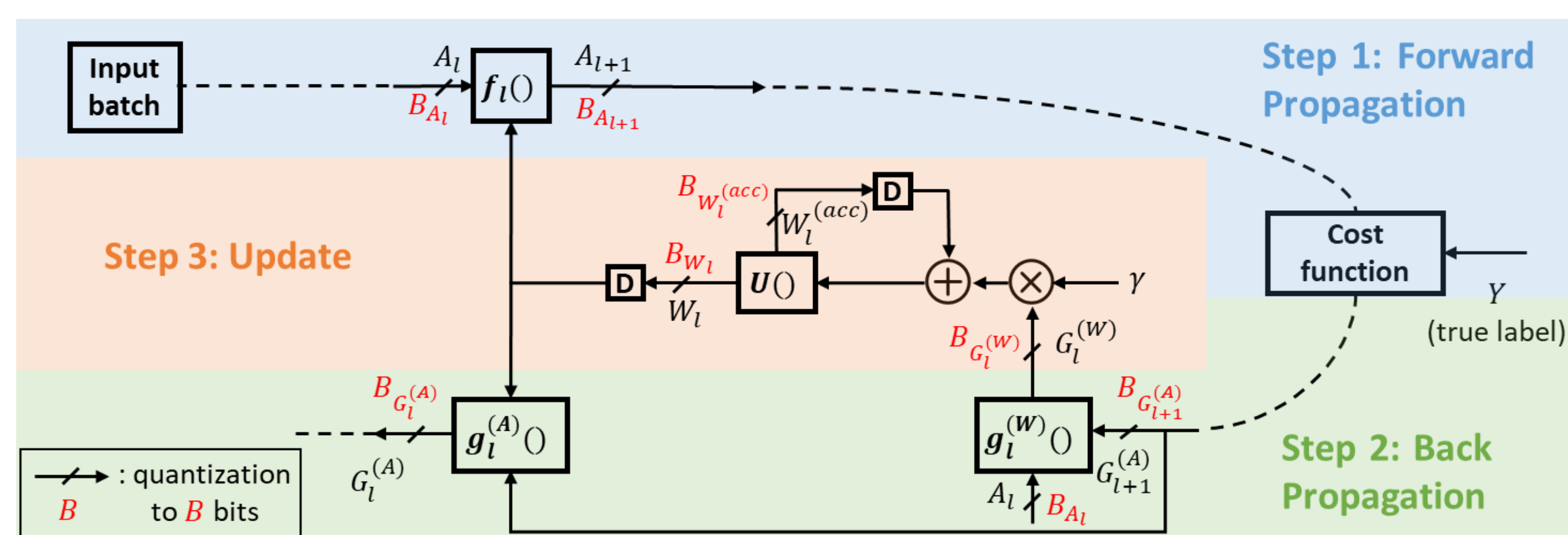


Sakr, Kim, Shanbhag ICML 2017  
Sakr & Shanbhag ICASSP 2018

Largely based on heuristics

What about training?

## Problem Setup and Challenges



- multiple forward quantization noise sources
- unknown gradient dynamic range
- instability due to quantization noise bias in updates
- back-propagation of quantization noise in activation gradients
- risk of premature stoppage of convergence

## Criteria-based Approach

Criterion 1: equalization of quantization noise gains

$$B_{W_i} = \text{rnd} \left( \log_2 \left( \sqrt{\frac{E_{W_i \rightarrow p_m}}{E^{(\min)}}} \right) \right) + B^{(\min)}$$

$$B_{A_i} = \text{rnd} \left( \log_2 \left( \sqrt{\frac{E_{A_i \rightarrow p_m}}{E^{(\min)}}} \right) \right) + B^{(\min)}$$

Criterion 2: proper gradient clipping

$$r_{G_i^{(W)}} \geq 2\sigma_{G_i^{(W)}}^{(\max)}$$

$$r_{G_{i+1}^{(A)}} \geq 4\sigma_{G_{i+1}^{(A)}}^{(\max)}$$

Criterion 3: quantization bias elimination

$$\Delta_{G_i^{(W)}} < \frac{\sigma_{G_i^{(W)}}^{(\min)}}{4}$$

Criterion 4: back-propagated noise bound

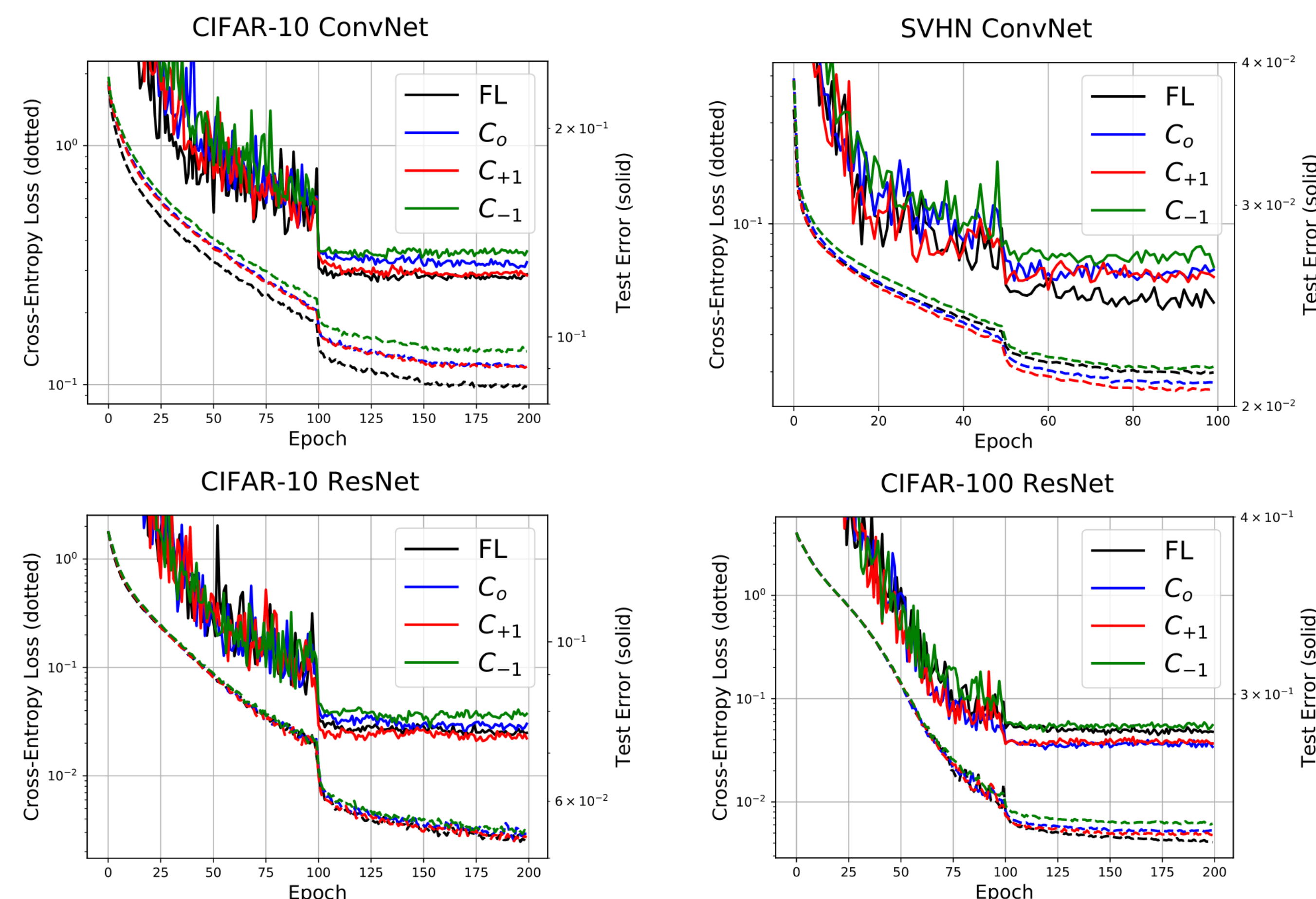
$$\Delta_{G_{i+1}^{(A)}} < \frac{\Delta_{G_i^{(W)}}}{\sqrt{\lambda_{G_{i+1}^{(A)} \rightarrow G_i^{(W)}}^{(\max)}}} \left( \frac{|G_i^{(W)}|}{|G_{i+1}^{(A)}|} \right)^{1/4}$$

Criterion 5: accumulator stopping condition

$$r_{W_i^{(acc)}} \geq 2^{-B_{W_i}}$$

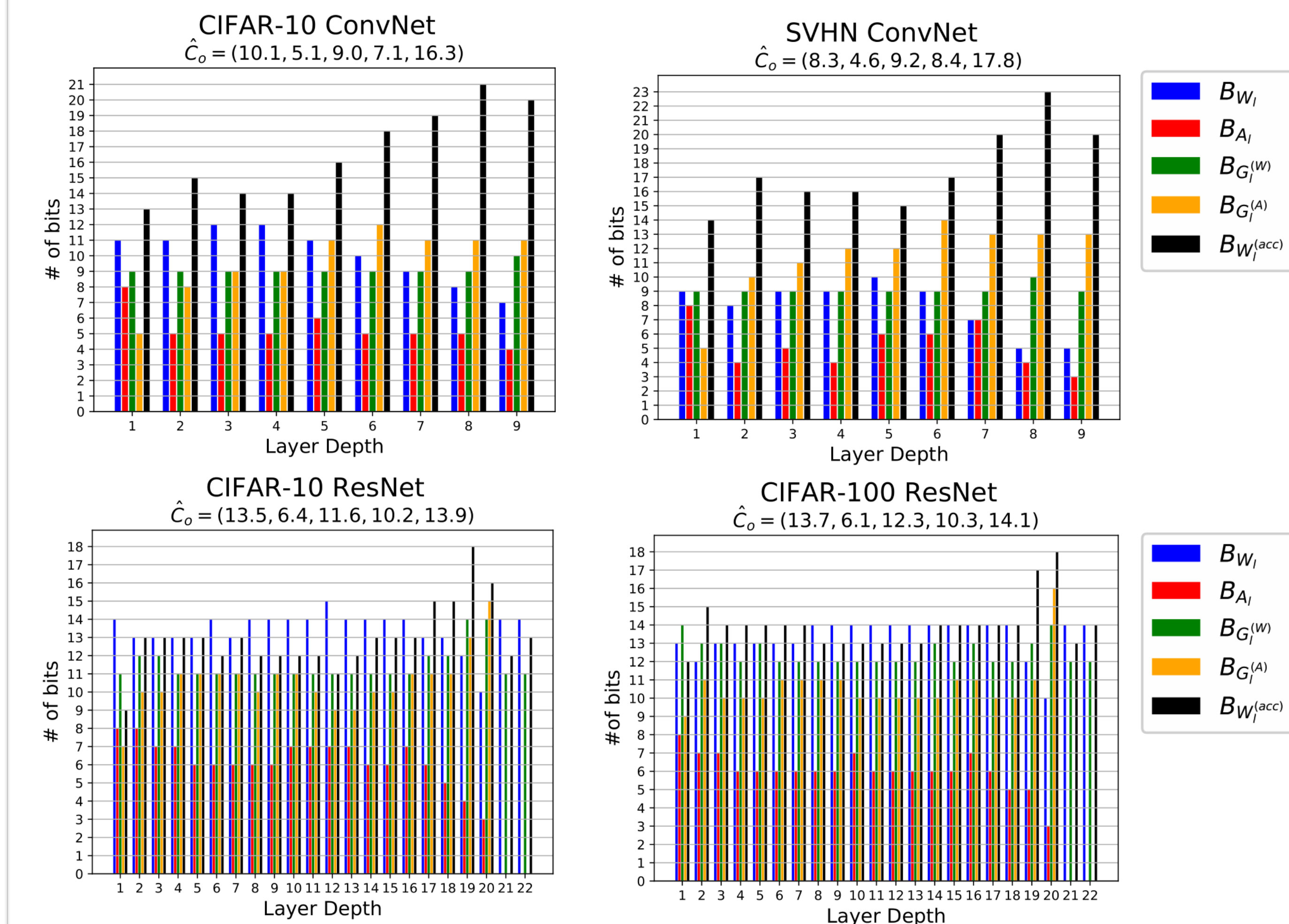
$$\Delta_{W_i^{(acc)}} < \gamma^{(\min)} \Delta_{G_i^{(W)}}$$

## Convergence with Close-to-Minimal Precision



- FX training was believed to be impossible due to dynamic range issues [Koester et al. – NIPS'2017]
- proposed FX training is able to match FL training accuracy
- precision assignment found to be nearly minimal

## Per-Layer Precision Trends



- weight precision decreases from network input to output
- precisions of activation gradients and weight accumulators increase
- ResNets have uniform precision requirements per tensor type

## Hyper-Precision Reduction is Inefficient

	$C_W$ (10 <sup>6</sup> b)	$C_A$ (10 <sup>6</sup> b)	$C_M$ (10 <sup>9</sup> FA)	$C_C$ (10 <sup>6</sup> b)	Test Error	$C_W$ (10 <sup>6</sup> b)	$C_A$ (10 <sup>6</sup> b)	$C_M$ (10 <sup>9</sup> FA)	$C_C$ (10 <sup>6</sup> b)	Test Error
<b>CIFAR-10 ConvNet</b>						<b>SVHN ConvNet</b>				
FL	148	9.3	94.4	49	12.02%	148	9.3	94.4	49	2.43%
FX ( $C_0$ )	<b>56.5</b>	<b>1.7</b>	11.9	14	12.58%	<b>54.3</b>	<b>1.9</b>	10.5	14	2.58%
BN	100	4.7	<b>2.8</b>	49	18.50%	100	4.7	<b>2.8</b>	49	3.60%
SQ	78.8	<b>1.7</b>	11.9	14	<b>11.32%</b>	76.3	<b>1.9</b>	10.5	14	2.73%
TG	102	9.3	94.4	<b>3.1</b>	12.49%	102	9.3	94.4	<b>3.1</b>	3.65%
<b>CIFAR-10 ResNet</b>						<b>CIFAR-100 ResNet</b>				
FL	1784	96	4319	596	7.42%	1789	97	4319	597	28.06%
FX ( $C_0$ )	<b>726</b>	<b>25</b>	785	216	7.51%	<b>750</b>	<b>25</b>	776	216	<b>27.43%</b>
BN	1208	50	<b>128</b>	596	<b>7.24%</b>	1211	50	<b>128</b>	597	29.35%
SQ	1062	<b>25</b>	785	216	7.42%	1081	<b>25</b>	776	216	28.03%
TG	1227	96	4319	<b>37.3</b>	7.94%	1230	97	4319	<b>37.3</b>	30.62%

- feedforward binarization (BN) and gradient ternarization (TG) fail to match FL accuracy for same topology
- stochastic quantization (SQ) provides marginal returns
- BN, TG, SQ do not address the fundamental problem of realizing true FX training

## Acknowledgement

This work was supported in part by C-BRIC, one of six centers in JUMP, a Semiconductor Research Corporation (SRC) program sponsored by DARPA.