



Trained Quantization Thresholds (TQT)

for Accurate and Efficient Fixed-Point Inference of Deep Neural Networks

Sambhav Jain^{^*}, Albert Gural^{#*}, Michael Wu[^], Chris Dick[^]

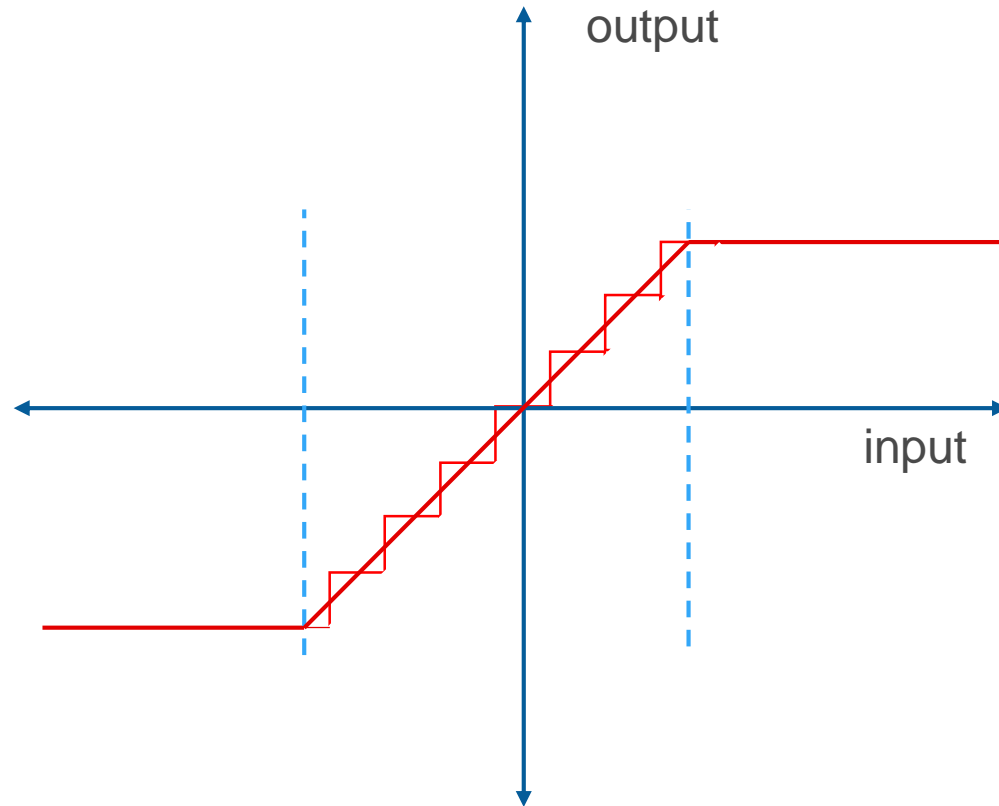
[^]Xilinx Inc., [#]Stanford University (*equal contribution)

March 3, 2020

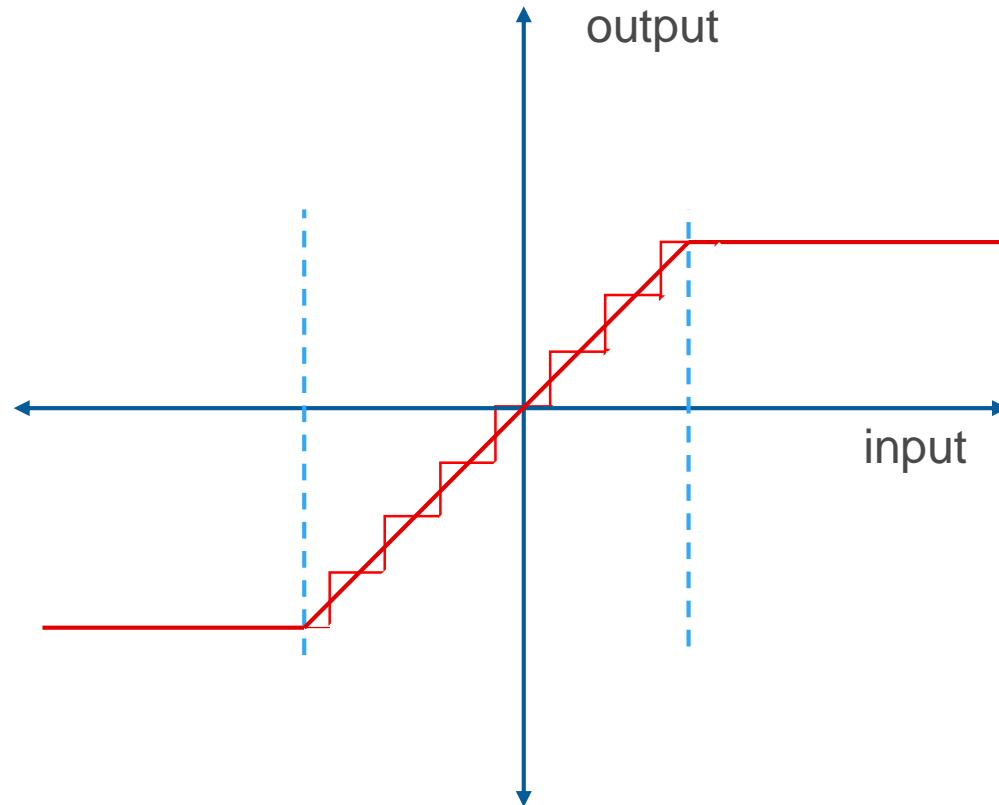
Background & Motivation



Uniform Quantization



Uniform Quantization

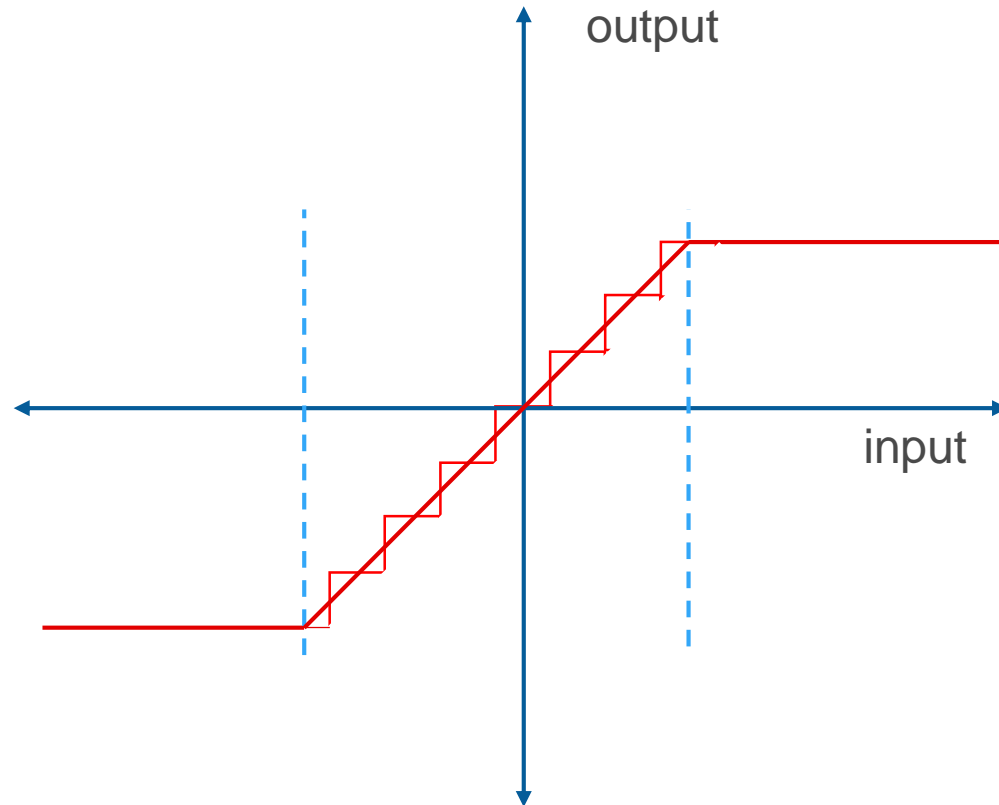


Affine Mapping

$$r = s \cdot (q - z)$$

real domain quantized domain

Uniform Quantization



Affine Mapping

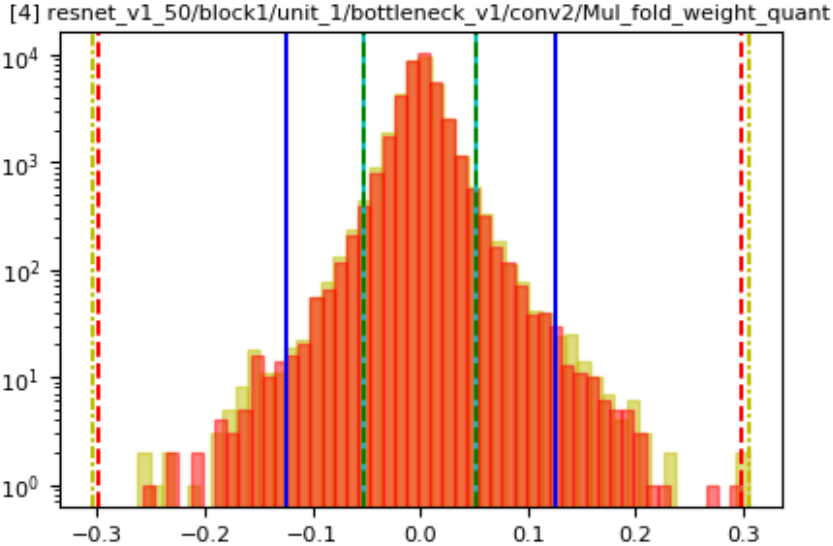
quantization parameters

$$r = s \cdot (q - z)$$

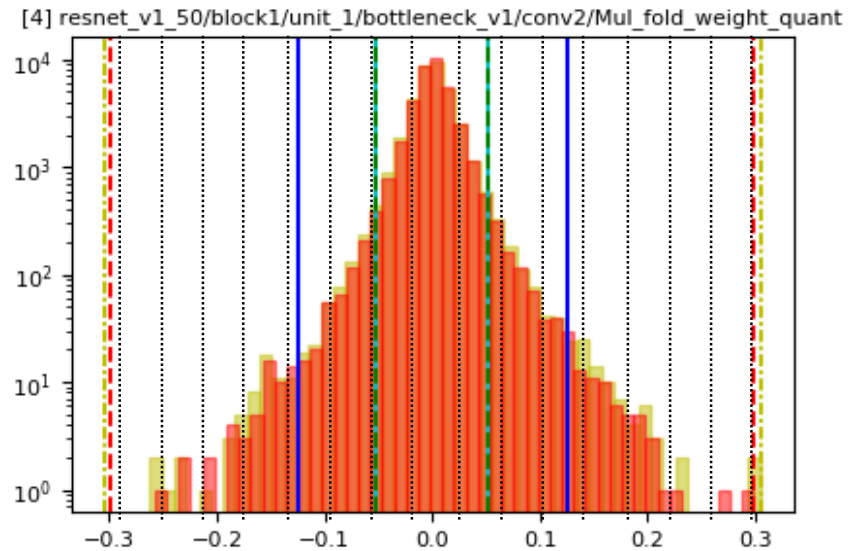
real domain quantized domain

The diagram shows the affine mapping equation $r = s \cdot (q - z)$. Red arrows point from the labels 'real domain' and 'quantized domain' to the variables r and z respectively. Black arrows point from the label 'quantization parameters' to the variables s and q .

Threshold Selection

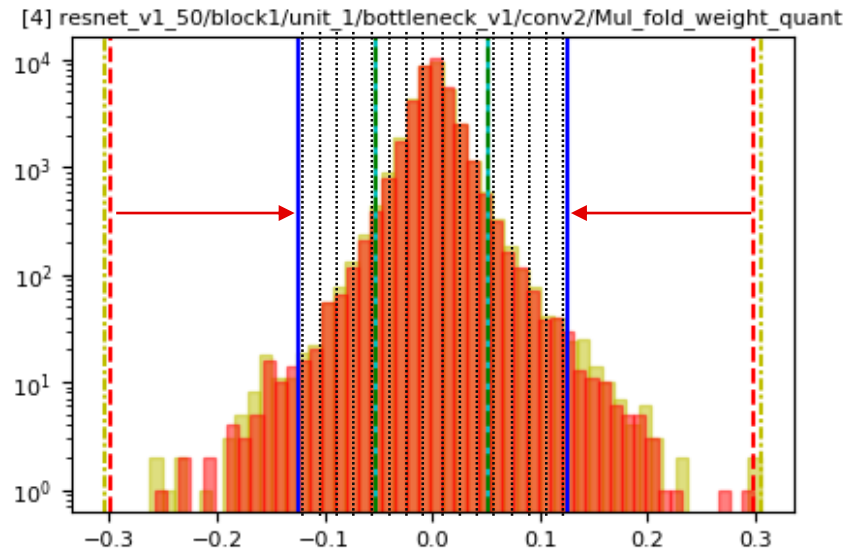


Threshold Selection



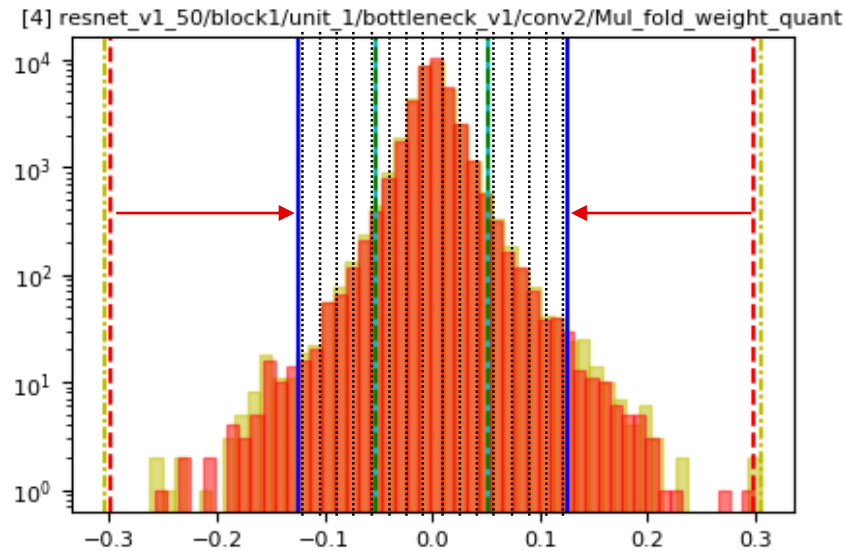
- ▶ Say 4-bit quantizer
 - 16 quantization levels
 - Clipping thresholds: dotted red lines (min, max)
 - Poor utilization of available precision

Threshold Selection



- ▶ Say 4-bit quantizer
 - 16 quantization levels
 - Clipping thresholds: blue lines
 - Better utilization of available precision

Threshold Selection



1

Statistical Methods

- Calibration
- KL divergence minimization
- SQNR maximization
- Percentile / nSD initialization
- ...

2

Gradient Descent Methods

- Google's QAT (Jacob et al., 2017)
- IBM's PACT (Choi et al., 2018)
- Xilinx's TQT (Jain et al., 2019)
- ...

Quantizer “Degrees of Freedom”

General case

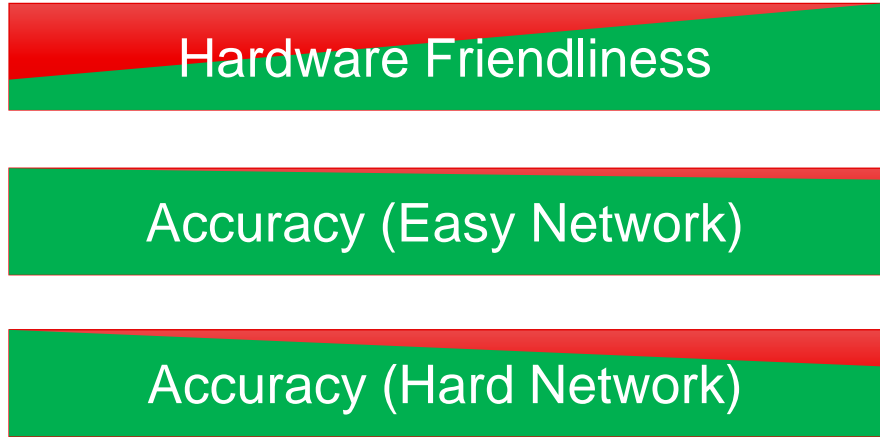
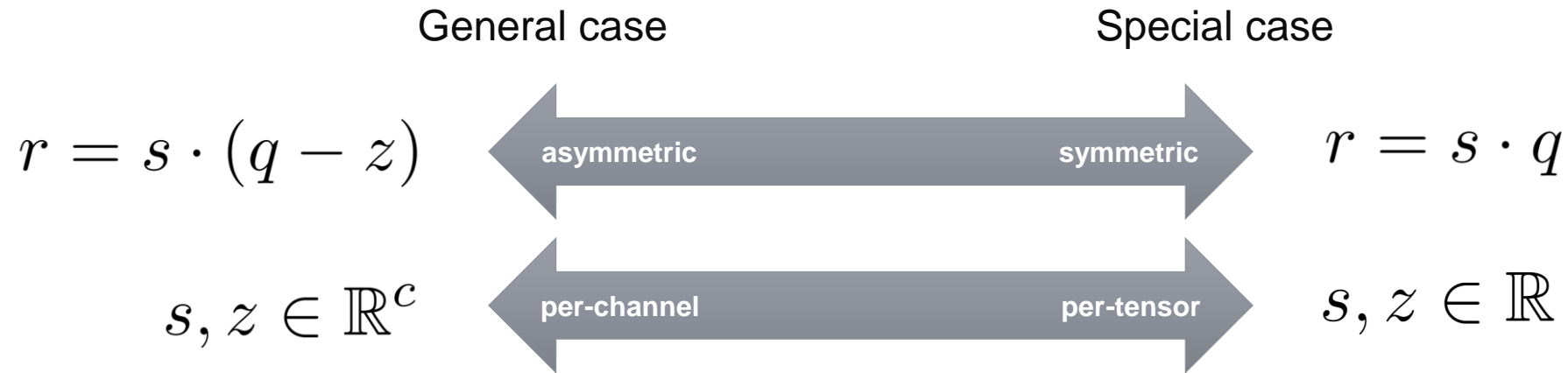
Special case



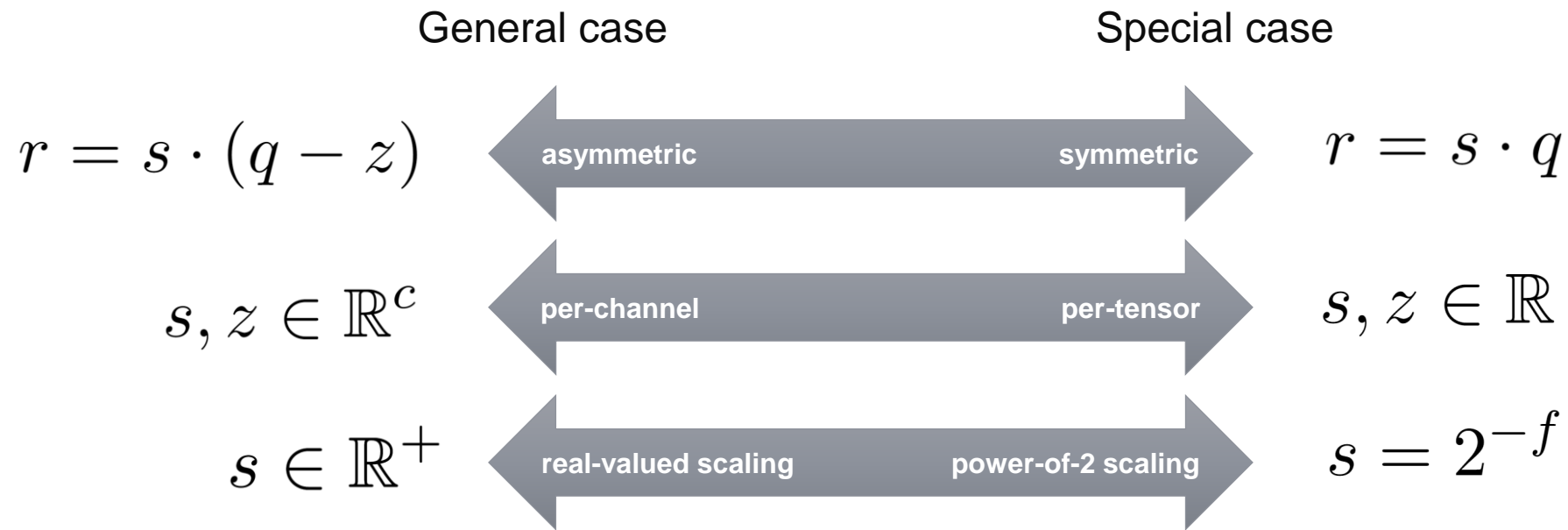
Quantizer “Degrees of Freedom”



Quantizer “Degrees of Freedom”



Quantizer “Degrees of Freedom”




Hardware Friendliness

Accuracy (Easy Network)

Accuracy (Hard Network)

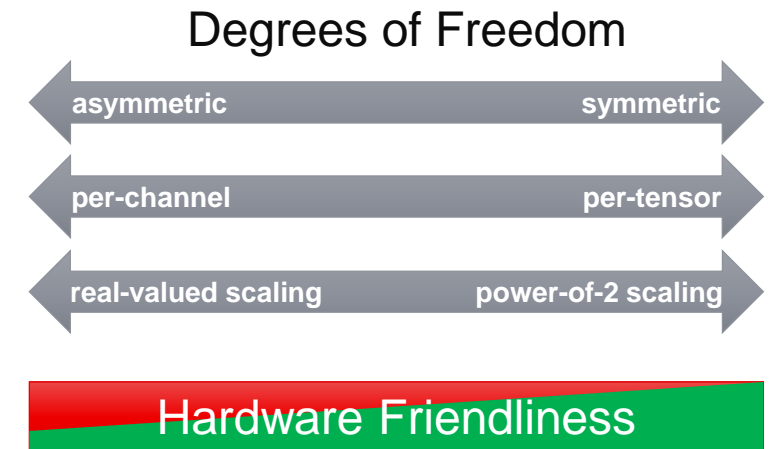
MobileNets are hard to quantize

real-valued scaling




Network	Asymmetric, per-layer (Post Training Quantization)	Symmetric, per-channel (Post Training Quantization)	Asymmetric, per-layer (Quantization Aware Training)	Symmetric, per-channel (Quantization Aware Training)	Floating Point
Mobilenet-v1_1_224	0.001	0.591	0.70	0.707	0.709
Mobilenet-v2_1_224	0.001	0.698	0.709	0.711	0.719
Nasnet-Mobile	0.722	0.721	0.73	0.73	0.74
Mobilenet-v2_1.4_224	0.004	0.74	0.735	0.745	0.749
Inception-v3	0.78	0.78	0.78	0.78	0.78
Resnet-v1_50	0.75	0.751	0.75	0.75	0.752
Resnet-v2_50	0.75	0.75	0.75	0.75	0.756
Resnet-v1_152	0.766	0.762	0.765	0.762	0.768
Resnet-v2_152	0.761	0.76	0.76	0.76	0.778

(Krishnamoorthi, 2018)



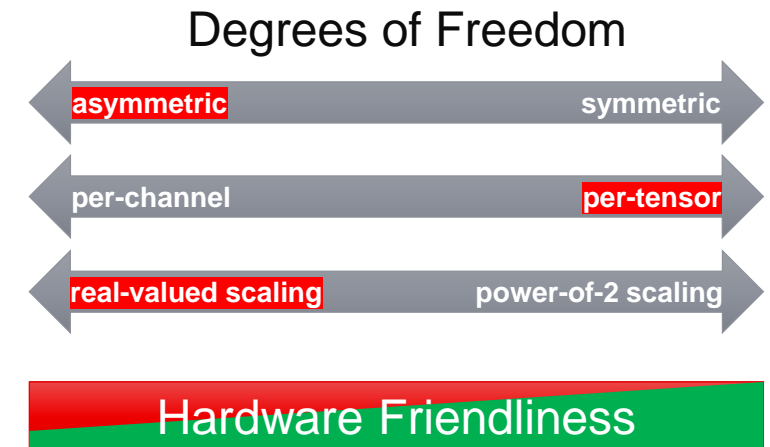
MobileNets are hard to quantize

real-valued scaling



Network	Asymmetric, per-layer (Post Training Quantization)	Symmetric, per-channel (Post Training Quantization)	Asymmetric, per-layer (Quantization Aware Training)	Symmetric, per-channel (Quantization Aware Training)	Floating Point
Mobilenet-v1_1_224	0.001	0.591	0.70	0.707	0.709
Mobilenet-v2_1_224	0.001	0.698	0.709	0.711	0.719
Nasnet-Mobile	0.722	0.721	0.73	0.73	0.74
Mobilenet-v2_1.4_224	0.004	0.74	0.735	0.745	0.749
Inception-v3	0.78	0.78	0.78	0.78	0.78
Resnet-v1_50	0.75	0.751	0.75	0.75	0.752
Resnet-v2_50	0.75	0.75	0.75	0.75	0.756
Resnet-v1_152	0.766	0.762	0.765	0.762	0.768
Resnet-v2_152	0.761	0.76	0.76	0.76	0.778

(Krishnamoorthi, 2018)

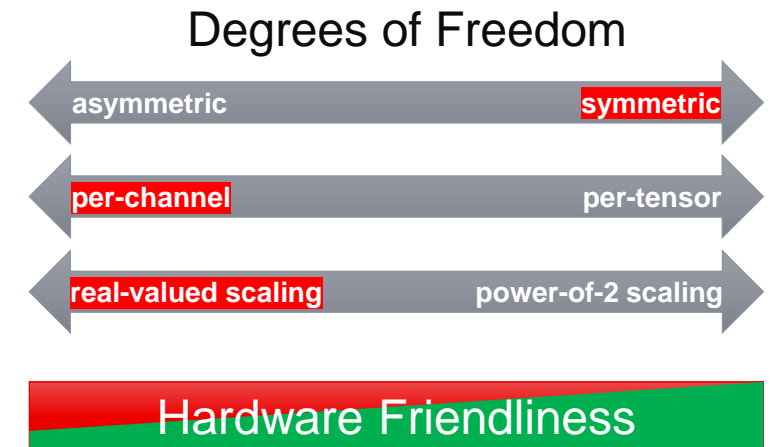


MobileNets are hard to quantize

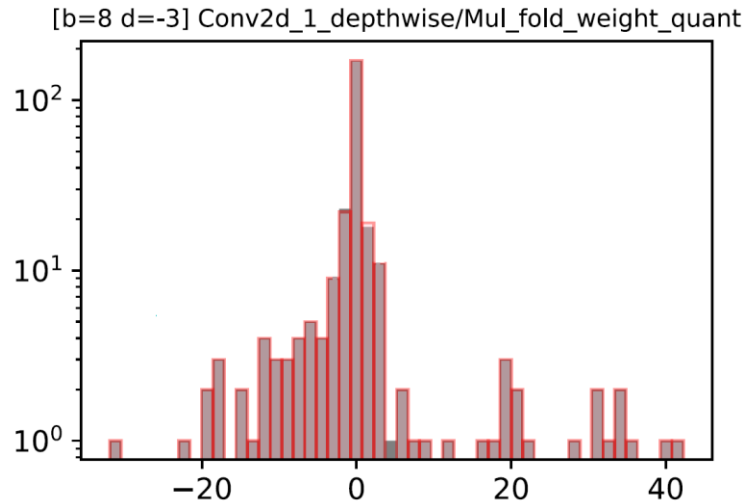
← real-valued scaling →

Network	Asymmetric, per-layer (Post Training Quantization)	Symmetric, per-channel (Post Training Quantization)	Asymmetric, per-layer (Quantization Aware Training)	Symmetric, per-channel (Quantization Aware Training)	Floating Point
Mobilenet-v1_1_224	0.001	0.591	0.70	0.707	0.709
Mobilenet-v2_1_224	0.001	0.698	0.709	0.711	0.719
Nasnet-Mobile	0.722	0.721	0.73	0.73	0.74
Mobilenet-v2_1.4_224	0.004	0.74	0.735	0.745	0.749
Inception-v3	0.78	0.78	0.78	0.78	0.78
Resnet-v1_50	0.75	0.751	0.75	0.75	0.752
Resnet-v2_50	0.75	0.75	0.75	0.75	0.756
Resnet-v1_152	0.766	0.762	0.765	0.762	0.768
Resnet-v2_152	0.761	0.76	0.76	0.76	0.778

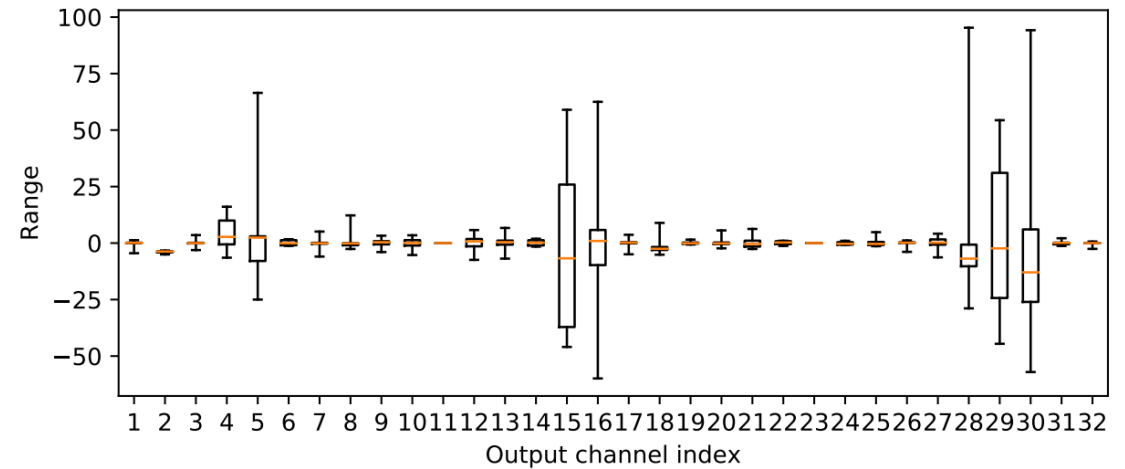
(Krishnamoorthi, 2018)



MobileNets are hard to quantize – Why?



Weight distribution in first depthwise separable layer of MobileNet v1

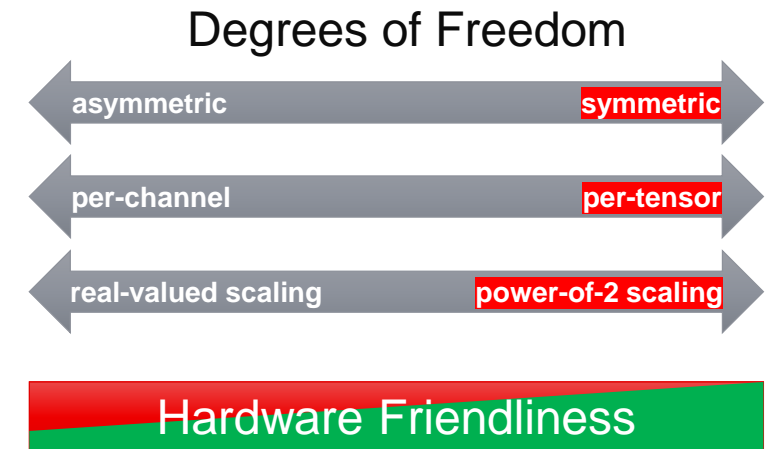


Dynamic range of weights (per-channel) in first depthwise separable layer of MobileNet v2 (Nagel et al., 2019)

With TQT: MobileNets can be quantized well

Method	Precision	Quantization Scheme	Top-1
MobileNet v1 1.0 224			
	FP32		70.9
QAT	INT8	per-channel, symmetric, real scaling	70.7
	INT8	per-tensor, asymmetric, real scaling	70.0
TQT	FP32		71.1
	INT8	per-tensor, symmetric, p-of-2 scaling	71.1
MobileNet v2 1.0 224			
	FP32		71.9
QAT	INT8	per-channel, symmetric, real scaling	71.1
	INT8	per-tensor, asymmetric, real scaling	70.9
TQT	FP32		71.7
	INT8	per-tensor, symmetric, p-of-2 scaling	71.8

ours





Trained Quantization Thresholds

Implementation

Forward Pass

$$q(x; s) := \text{clip} \left(\left\lfloor \frac{x}{s} \right\rfloor ; n, p \right) \cdot s$$

where $n = -2^{b-1}$, $p = 2^{b-1} - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^{b-1}}$ for signed data; $n = 0$, $p = 2^b - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^b}$ for unsigned data.

Backward Pass

Implementation

Forward Pass

$$q(x; s) := \begin{cases} \left\lfloor \frac{x}{s} \right\rfloor \cdot s & \text{if } n \leq \left\lfloor \frac{x}{s} \right\rfloor \leq p, \\ n \cdot s & \text{if } \left\lfloor \frac{x}{s} \right\rfloor < n, \\ p \cdot s & \text{if } \left\lfloor \frac{x}{s} \right\rfloor > p. \end{cases}$$

where $n = -2^{b-1}$, $p = 2^{b-1} - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^{b-1}}$ for signed data; $n = 0$, $p = 2^b - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^b}$ for unsigned data.

Backward Pass

Implementation

Forward Pass

$$q(x; s) := \begin{cases} \left\lfloor \frac{x}{s} \right\rfloor \cdot s & \text{if } n \leq \left\lfloor \frac{x}{s} \right\rfloor \leq p, \\ n \cdot s & \text{if } \left\lfloor \frac{x}{s} \right\rfloor < n, \\ p \cdot s & \text{if } \left\lfloor \frac{x}{s} \right\rfloor > p. \end{cases}$$

where $n = -2^{b-1}$, $p = 2^{b-1} - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^{b-1}}$ for signed data; $n = 0$, $p = 2^b - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^b}$ for unsigned data.

Backward Pass

$\frac{\partial}{\partial x} \lfloor x \rfloor = \frac{\partial}{\partial x} \lceil x \rceil = 1$, but $\lfloor x \rfloor \neq x$ and $\lceil x \rceil \neq x$ (Straight-Through Estimator)

In the backward pass, approximate gradients of round/ceil to 1, without approximating round/ceil to be identity

Implementation

Forward Pass

$$q(x; s) := \begin{cases} \left\lfloor \frac{x}{s} \right\rfloor \cdot s & \text{if } n \leq \left\lfloor \frac{x}{s} \right\rfloor \leq p, \\ n \cdot s & \text{if } \left\lfloor \frac{x}{s} \right\rfloor < n, \\ p \cdot s & \text{if } \left\lfloor \frac{x}{s} \right\rfloor > p. \end{cases}$$

where $n = -2^{b-1}$, $p = 2^{b-1} - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^{b-1}}$ for signed data; $n = 0$, $p = 2^b - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^b}$ for unsigned data.

Backward Pass

$$\nabla_s q(x; s) := \begin{cases} \left\lfloor \frac{x}{s} \right\rfloor - \frac{x}{s} & \text{if } n \leq \left\lfloor \frac{x}{s} \right\rfloor \leq p \\ n & \text{if } \left\lfloor \frac{x}{s} \right\rfloor < n \\ p & \text{if } \left\lfloor \frac{x}{s} \right\rfloor > p \end{cases}$$

non-zero!

$\frac{\partial}{\partial x} \lfloor x \rfloor = \frac{\partial}{\partial x} \lceil x \rceil = 1$, but $\lfloor x \rfloor \neq x$ and $\lceil x \rceil \neq x$ (Straight-Through Estimator)

Implementation

Forward Pass

$$q(x; s) := \begin{cases} \left\lfloor \frac{x}{s} \right\rfloor \cdot s & \text{if } n \leq \left\lfloor \frac{x}{s} \right\rfloor \leq p, \\ n \cdot s & \text{if } \left\lfloor \frac{x}{s} \right\rfloor < n, \\ p \cdot s & \text{if } \left\lfloor \frac{x}{s} \right\rfloor > p. \end{cases}$$

where $n = -2^{b-1}$, $p = 2^{b-1} - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^{b-1}}$ for signed data; $n = 0$, $p = 2^b - 1$ and $s = \frac{2^{\lceil \log_2 t \rceil}}{2^b}$ for unsigned data.

Backward Pass

$$\nabla_s q(x; s) := \begin{cases} \left\lfloor \frac{x}{s} \right\rfloor - \frac{x}{s} & \text{if } n \leq \left\lfloor \frac{x}{s} \right\rfloor \leq p \\ n & \text{if } \left\lfloor \frac{x}{s} \right\rfloor < n \\ p & \text{if } \left\lfloor \frac{x}{s} \right\rfloor > p \end{cases}$$

non-zero!

$$\nabla_x q(x; s) := \begin{cases} 1 & \text{if } n \leq \left\lfloor \frac{x}{s} \right\rfloor \leq p, \\ 0 & \text{otherwise.} \end{cases}$$

$\frac{\partial}{\partial x} \lfloor x \rfloor = \frac{\partial}{\partial x} \lceil x \rceil = 1$, but $\lfloor x \rfloor \neq x$ and $\lceil x \rceil \neq x$ (Straight-Through Estimator)

Transfer Curves

$$\nabla_{(\log_2 t)} q(x; s) := s \ln(2) \cdot \begin{cases} \left\lfloor \frac{x}{s} \right\rfloor - \frac{x}{s} & \text{if } n \leq \left\lfloor \frac{x}{s} \right\rfloor \leq p, \\ n & \text{if } \left\lfloor \frac{x}{s} \right\rfloor < n, \\ p & \text{if } \left\lfloor \frac{x}{s} \right\rfloor > p \end{cases}$$

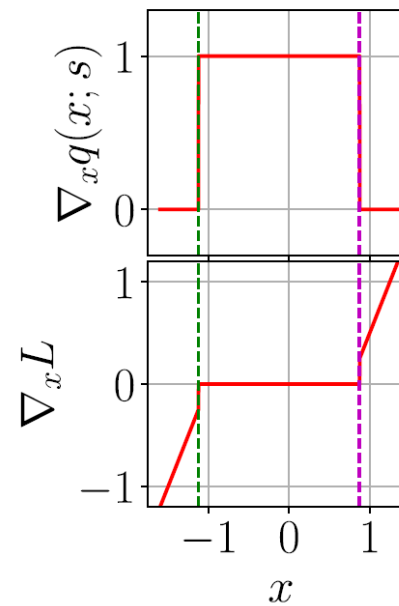
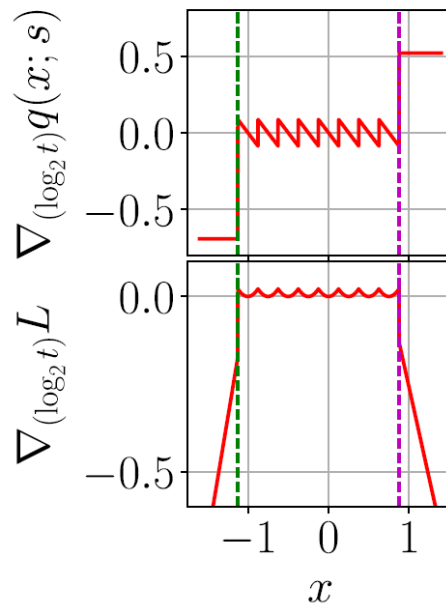
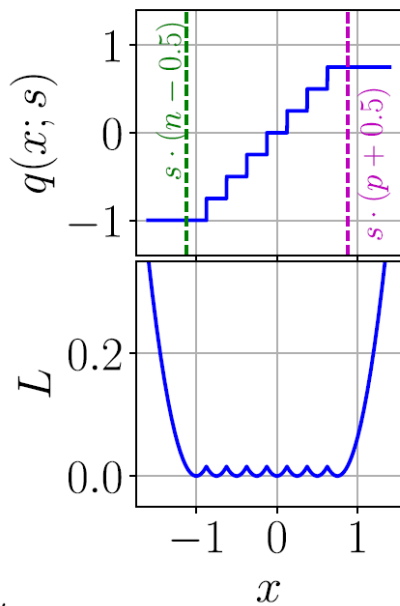
$$\nabla_x q(x; s) := \begin{cases} 1 & \text{if } n \leq \left\lfloor \frac{x}{s} \right\rfloor \leq p, \\ 0 & \text{otherwise.} \end{cases}$$

Threshold Gradient

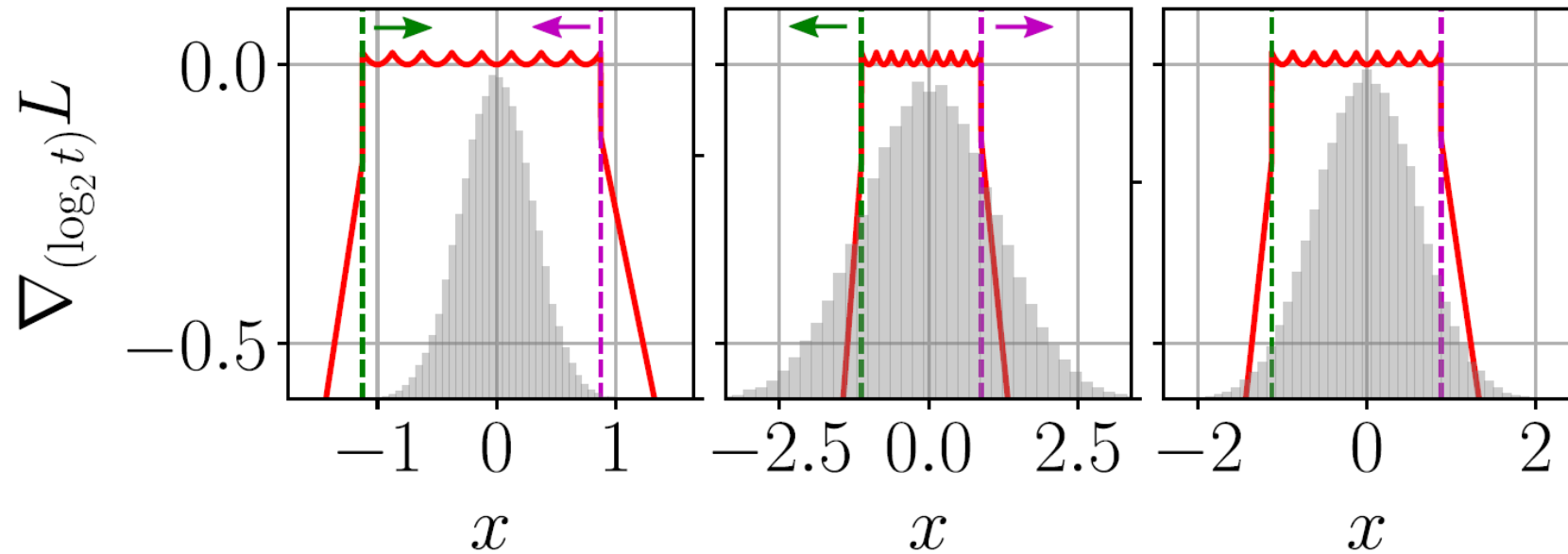
Input Gradient

Toy L₂ loss

$$L = (q(x; s) - x)^2 / 2$$



Range Precision Trade-off



$$\log_2 t := \log_2 t - \alpha \cdot \nabla_{\log_2 t} L \quad (\text{Update Rule})$$

Clipped Threshold Gradients

PACT's threshold gradients:

$$\frac{\partial y_q(x; \alpha)}{\partial \alpha} = \begin{cases} 0 & x \in (-\infty, \alpha) \\ 1 & x \in [\alpha, +\infty) \end{cases}$$

Clipped Threshold Gradients

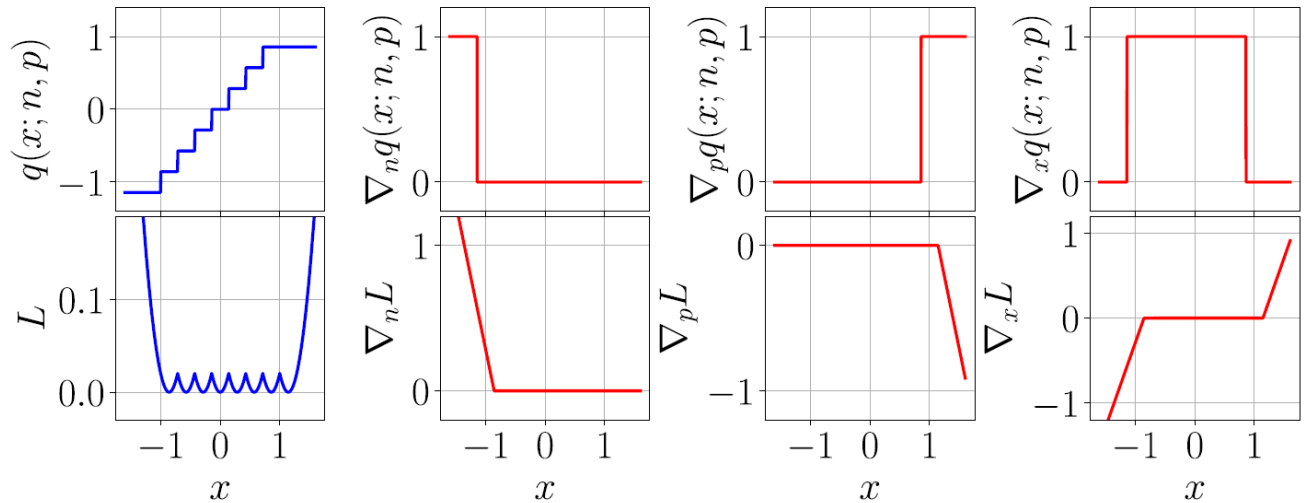
PACT's threshold gradients

$$\frac{\partial y_q(x; \alpha)}{\partial \alpha} = \begin{cases} 0 & x \in (-\infty, \alpha) \\ 1 & x \in [\alpha, +\infty) \end{cases}$$

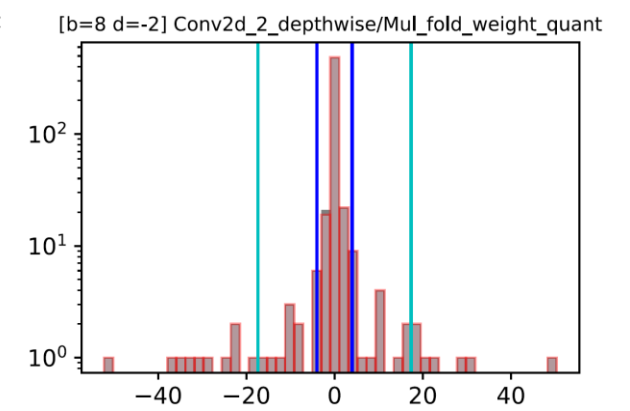
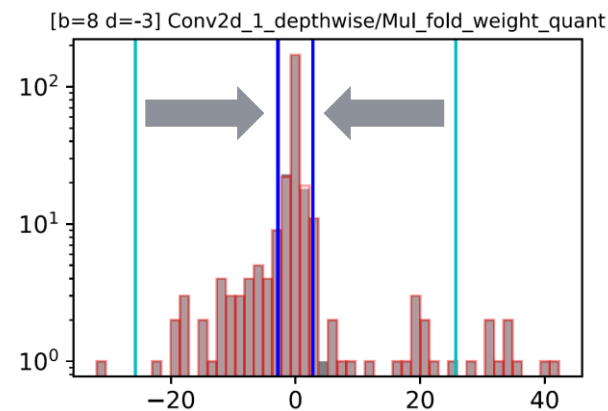
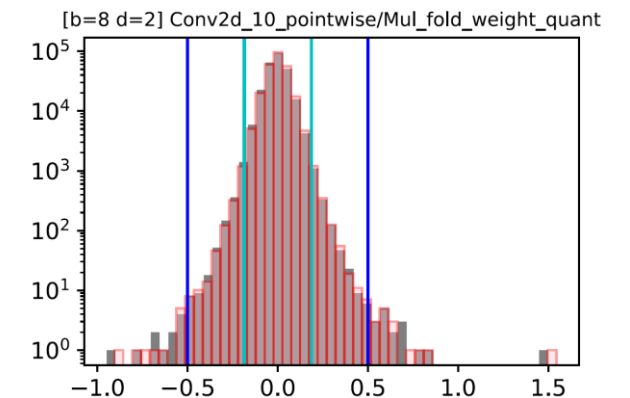
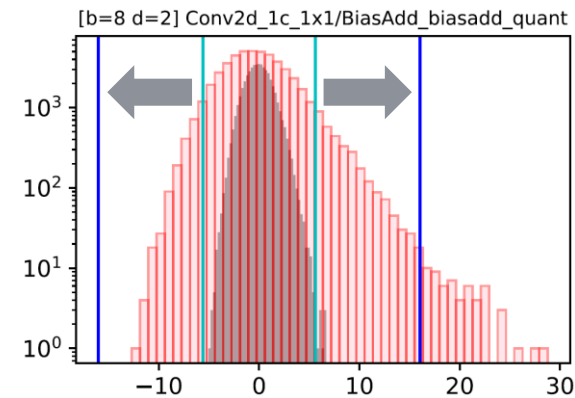
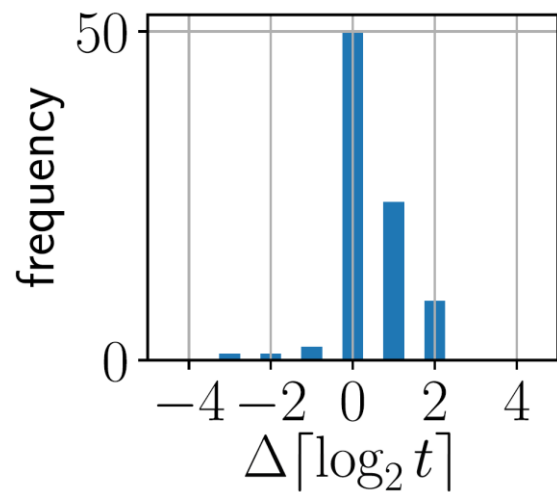
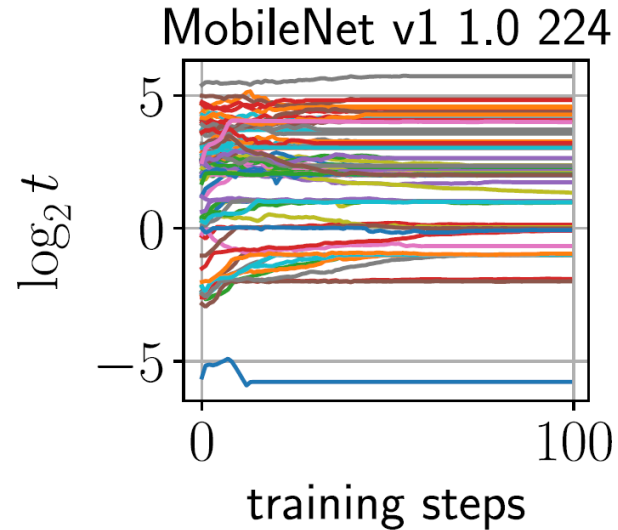
QAT's threshold gradients (FakeQuant)

$$q(x; n, p) := \left\lfloor \frac{\text{clip}(x; n, p) - n}{\frac{p - n}{2^b - 1}} \right\rfloor \cdot \frac{p - n}{2^b - 1} + n$$

$$\frac{\partial q(x; n, p)}{\partial p} = \begin{cases} 0 & x \in (-\infty, p) \\ 1 & x \in [p, +\infty) \end{cases}$$

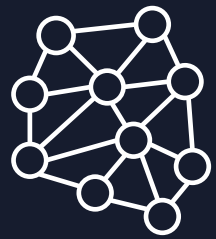


Distributions after TQT retraining



Results

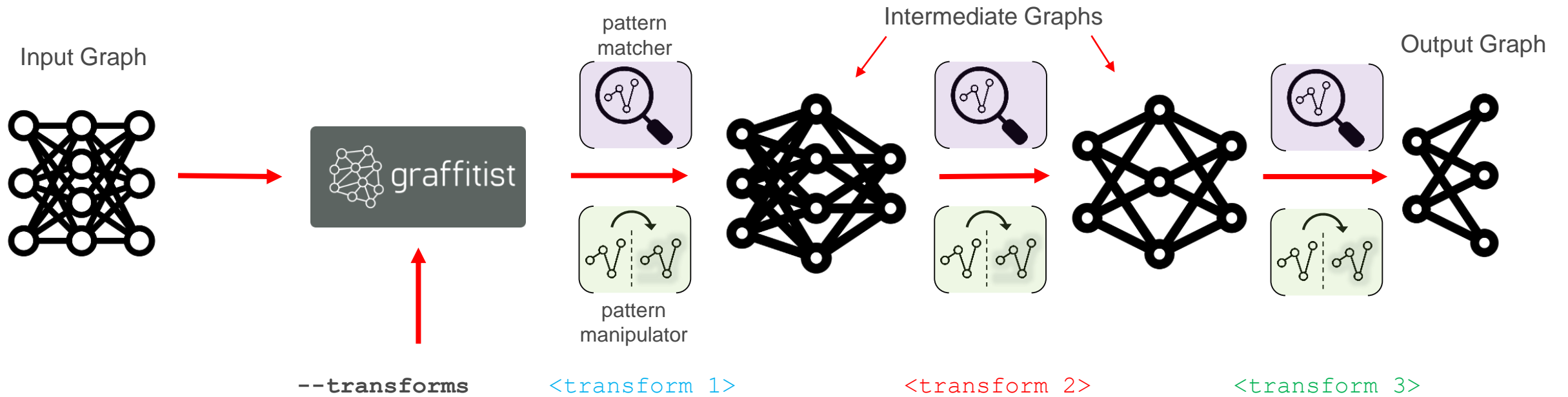
Mode	Precision	Bit-width	Accuracy (%)		Epochs
		(W/A)	top-1	top-5	
MobileNet v1 1.0 224					
Static	FP32	32/32	71.0	90.0	
	INT8	8/8	0.6	3.6	
Retrain	wt	FP32	71.1	90.0	3.4
	wt	INT8	67.0	87.9	4.6
	wt,th	INT8	71.1	90.0	2.1
	wt,th	INT4	4/8	–	–
MobileNet v2 1.0 224					
Static	FP32	32/32	70.1	89.5	
	INT8	8/8	0.3	1.2	
Retrain	wt	FP32	71.7	90.7	3.2
	wt	INT8	68.2	89.0	2.7
	wt,th	INT8	71.8	90.6	2.2
	wt,th	INT4	4/8	–	–



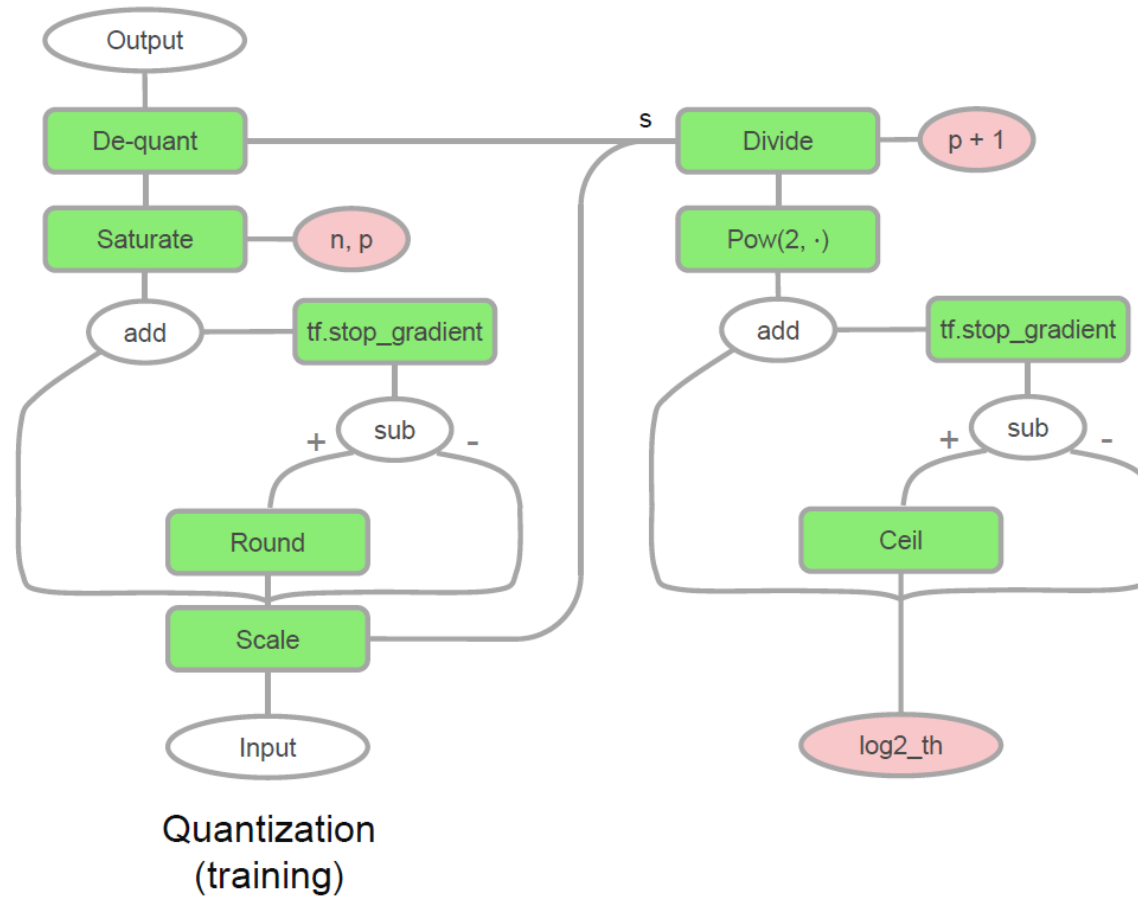
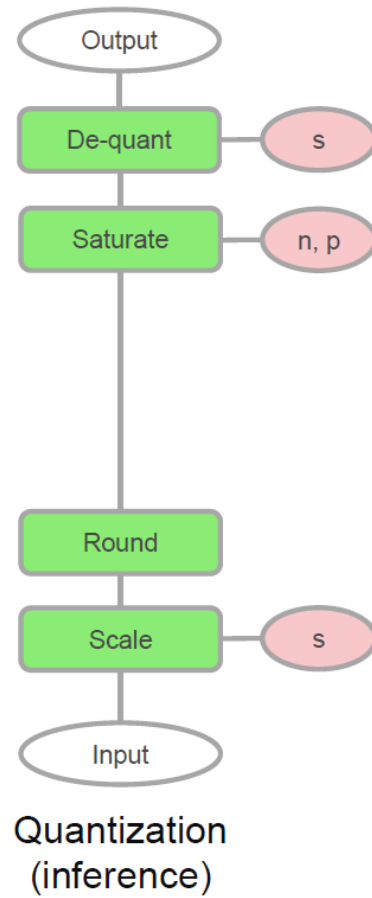
graffitist

github.com/Xilinx/graffitist

Tool for Neural Net Optimizations



Quantization Layer for TQT (unfused)



Layer Precisions

Conv/FC

$$q_8 \left(q'_{16} \left(\sum (q_{8/4}(w) \cdot q_8(x)) \right) + q'_{16}(b) \right)$$

Avgpool

$$q_8 \left(\sum (q_8(r) \cdot q_8(x)) \right)$$

Eltwise Add

$$q_8 (q'_8(x) + q'_8(y))$$

Concat

$$\text{concat}(q'_8(x), q'_8(y), q'_8(z))$$

Graph Optimizations

- ▶ BatchNorm folding (adopt best practices from Jacob et al., 2017)
 - Ensure folded batch norms in training and inference graphs are mathematically equivalent
 - Apply batch norm corrections (reduce training jitter by switching between batch and moving average statistics)
 - Freeze batch norm moving mean and variance updates post convergence for improved accuracy
- ▶ Explicitly merging input scales for scale preserving ops such as concat, bias-add, eltwise-add, and maximum (for leaky relu)
- ▶ Collapsing concat-of-concat layers into single concat, splicing identity nodes
- ▶ Modeling average pool layers as depthwise conv layers with reciprocal multiplier as weights to enable quantization



Thank You



Backup



TQT - Components

Input Graph

Training Platform (TF)

Dataset

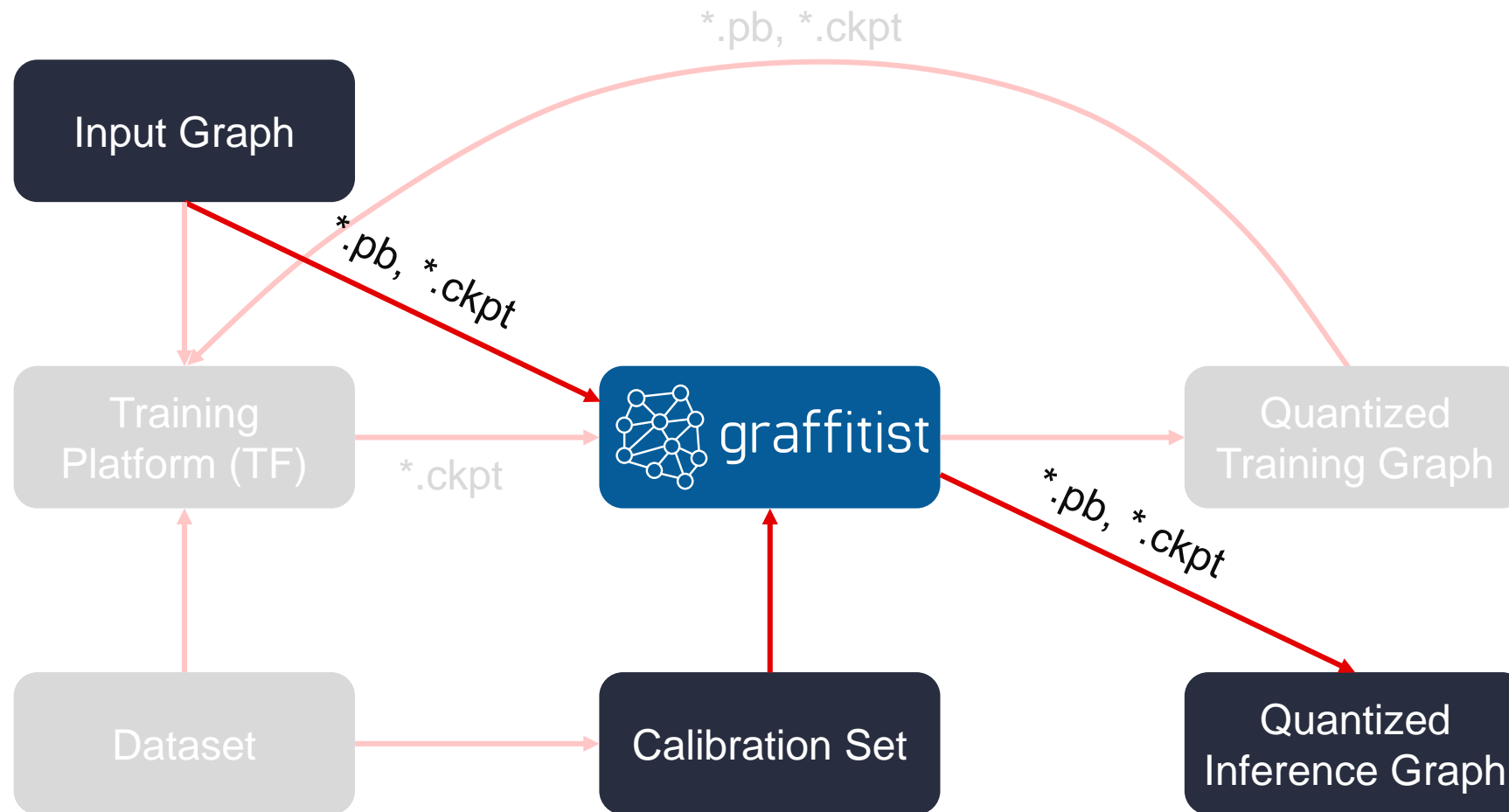


Calibration Set

Quantized Training Graph

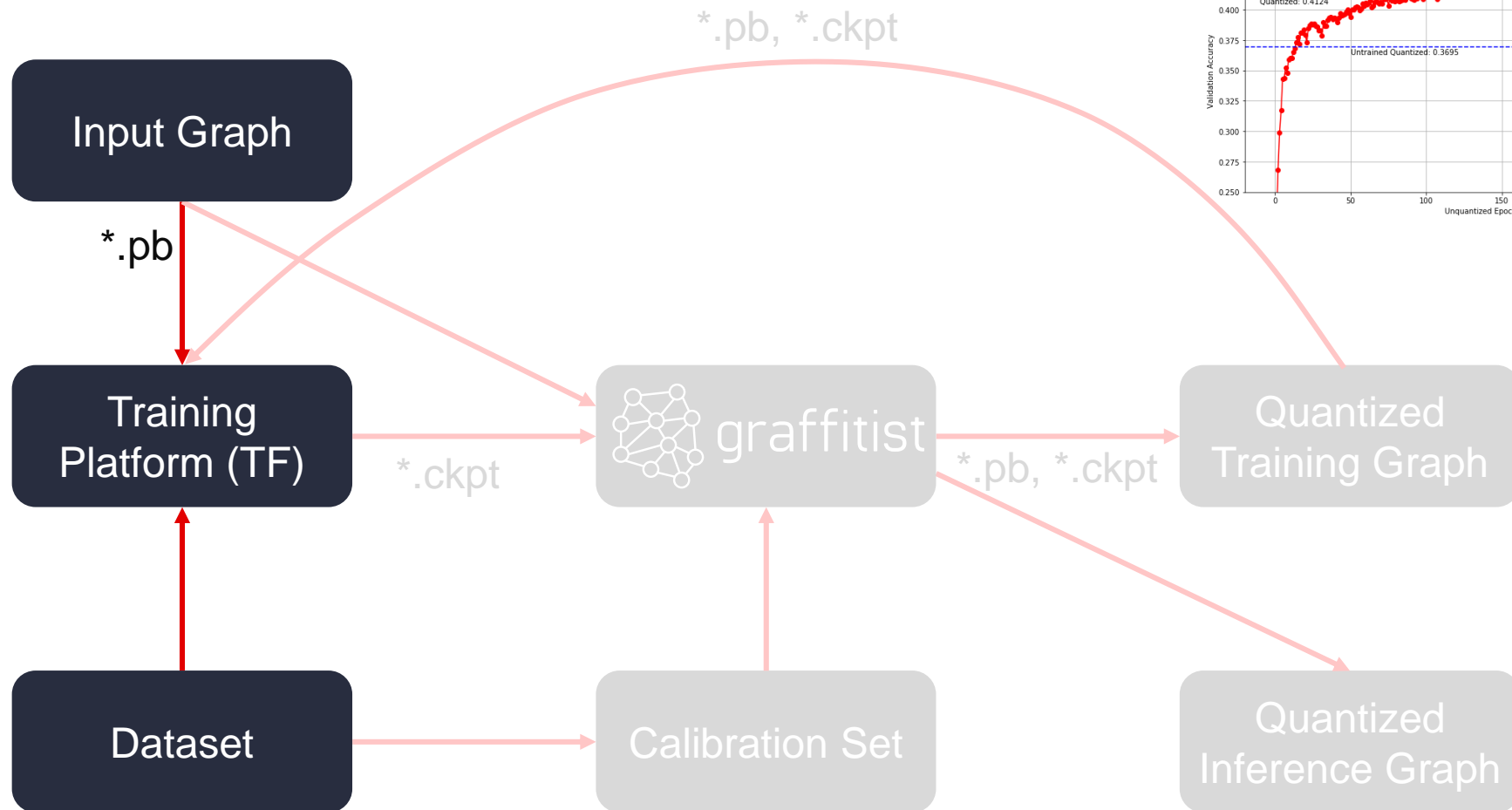
Quantized Inference Graph

Static Mode



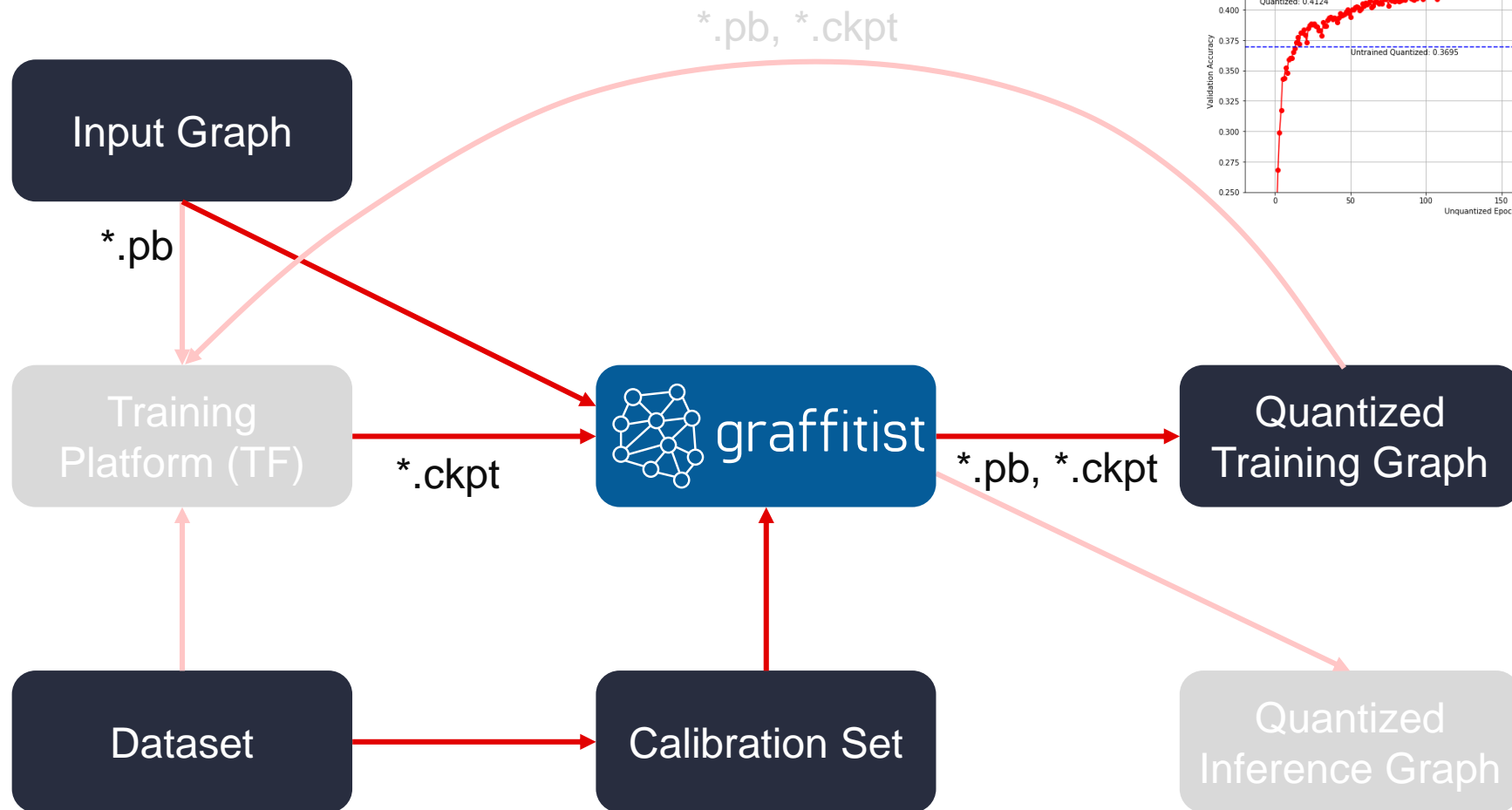
Retrain Mode

Step 1: Train on the original input graph (or use pre-trained weights)



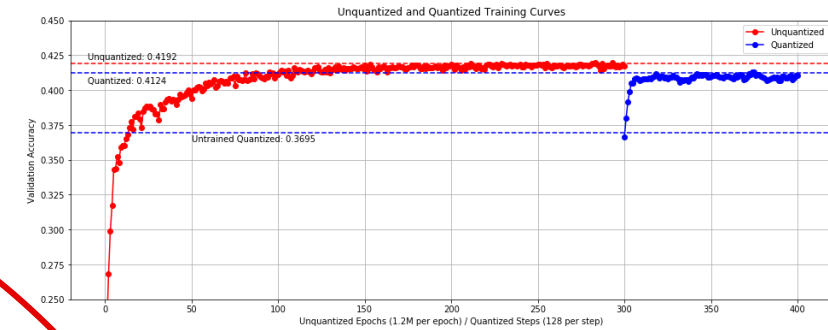
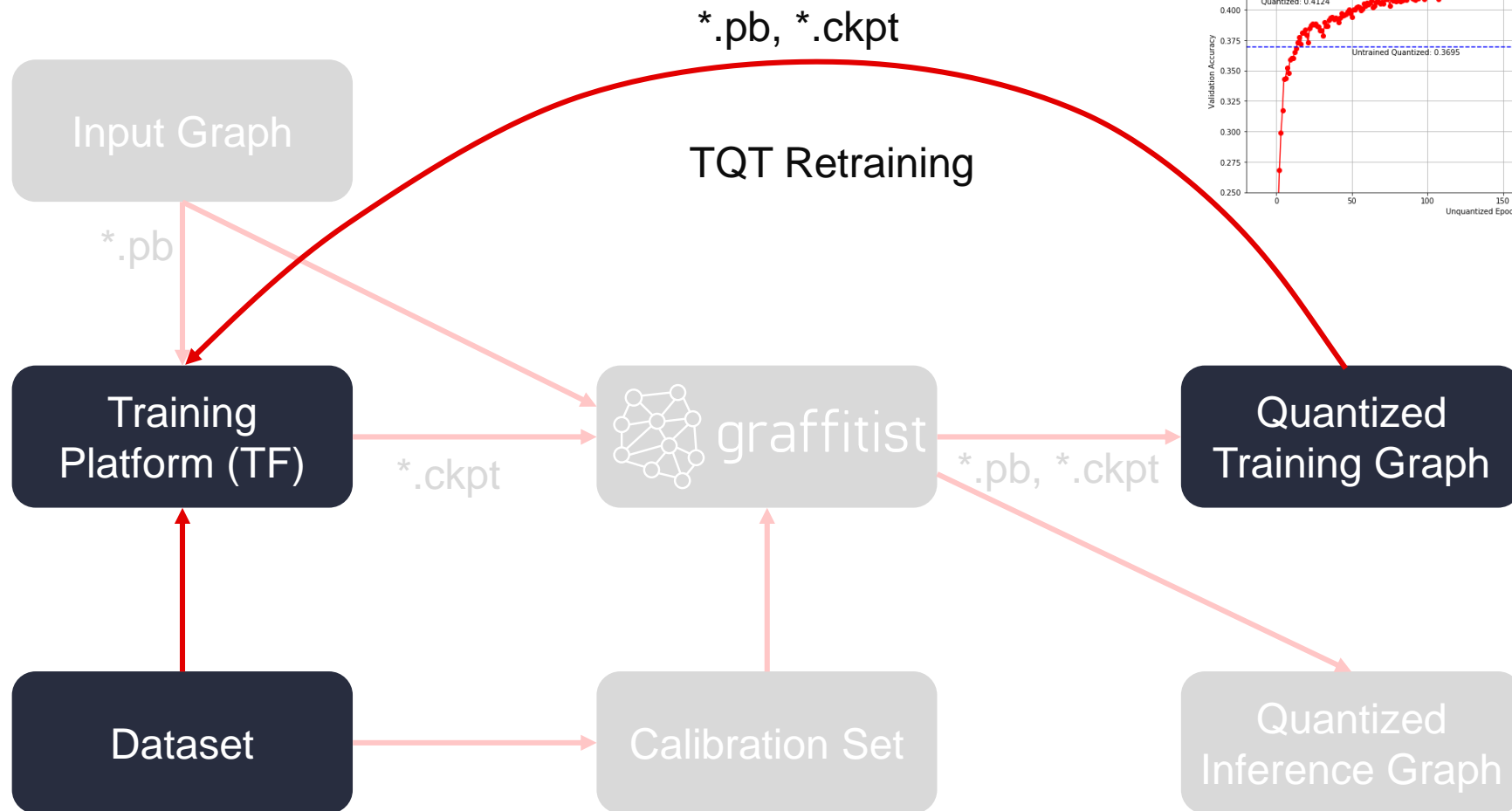
Retrain Mode

Step 2: Generate quantized training graph; calibrate thresholds



Retrain Mode

Step 3: Retrain quantized training graph (learn weights & thresholds)



Retrain Mode

Step 4: Generate quantized inference graph for compiler

