

METAQUANT: LEARNING TO QUANTIZE BY LEARNING TO PENETRATE NON-DIFFERENTIABLE QUANTIZATION

Shangyu Chen (schen025@e.ntu.edu.sg), Wenya Wang (wangwy@ntu.edu.sg), Sinno Jialin Pan(sinnopan@ntu.edu.sg)

School of Computer Science and Engineering, Nanyang Technological University, Singapore



Motivation

Existing training-based quantization methods rely on Straight-Through-Estimator (STE) to enable **non-differentiable** quantization training:

$$\ell = \text{Loss}(f(Q(\mathbf{W}); \mathbf{x}), y), \quad \frac{\partial Q(r)}{\partial r} = \begin{cases} 1 & \text{if } |r| \leq 1 \\ 0 & \text{else.} \end{cases} \quad (1)$$

However, it inevitably brings the problem:

- **Gradient Mismatch**: the gradients of the weights are not generated using the value of weights, but rather its quantized value.
- **Poor Gradient**: STE fails at investigating better gradients for quantization training.

We propose to learn $\frac{\partial Q(r)}{\partial r}$ by a neural network (\mathcal{M}) during quantization training. Such neural network is called *meta quantizer* and is trained together with the base quantized model.

Contribution

- Proper gradients is learned in quantization training without any manually designed knowledge.
- Generated gradient is loss-aware, contributing to the training of base model.
- Meta quantizer is removed and consumes no space in deployment.

Problem Statement

- Pre-trained full-precision base model f with L -layer parameterized as $\mathbf{W} = [\mathbf{W}_1, \dots, \mathbf{W}_L]$.
- Pre-processing function $\mathcal{A}(\cdot)$, pre-quantized weights: $\tilde{\mathbf{W}}$.
- Quantization function $Q(\cdot)$, quantized weights: $\hat{\mathbf{W}}$.
- Forward quantization:

$$\text{dorefa} : \tilde{\mathbf{W}} = \mathcal{A}(\mathbf{W}) = \frac{\tanh(\mathbf{W})}{2\max(|\tanh(\mathbf{W})|)} + \frac{1}{2}, \quad \hat{\mathbf{W}} = Q(\tilde{\mathbf{W}}) = 2 \frac{\text{round}[(2^k - 1)\tilde{\mathbf{W}}]}{2^k - 1} - 1.$$

$$\text{BWN} : \tilde{\mathbf{W}} = \mathcal{A}(\mathbf{W}) = \mathbf{W}, \quad \hat{\mathbf{W}} = Q(\tilde{\mathbf{W}}) = \frac{1}{n} \|\tilde{\mathbf{W}}\|_{l_1} \times \text{sign}(\tilde{\mathbf{W}}).$$

Generation of Meta Gradient

- Gradient of quantized weights: $g_{\tilde{\mathbf{W}}} = \frac{\partial \ell}{\partial \tilde{\mathbf{W}}}$.
- Meta Gradient / Gradient of pre-quantized weights:

$$g_{\tilde{\mathbf{W}}} = \mathcal{M}_{\phi}(g_{\hat{\mathbf{W}}}, \tilde{\mathbf{W}}) \quad (2)$$

- Gradient of full-precision weights:

$$g_{\mathbf{W}} = \frac{\partial \ell}{\partial \tilde{\mathbf{W}}} \frac{\partial \tilde{\mathbf{W}}}{\partial \mathbf{W}} = g_{\tilde{\mathbf{W}}} \frac{\partial \tilde{\mathbf{W}}}{\partial \mathbf{W}} = \mathcal{M}_{\phi}(g_{\hat{\mathbf{W}}}, \tilde{\mathbf{W}}) \frac{\partial \tilde{\mathbf{W}}}{\partial \mathbf{W}} \quad (3)$$

- Calibration: **dorefa**: $\frac{\partial \tilde{\mathbf{W}}}{\partial \mathbf{W}} = \frac{1 - \tanh^2(\mathbf{W})}{\max(|\tanh(\mathbf{W})|)}$, **BWN**: $\frac{\partial \tilde{\mathbf{W}}}{\partial \mathbf{W}} = 1$
- Gradient Refinement: **SGD**: $\pi(g_{\mathbf{W}}) = g_{\mathbf{W}}$, **Adam**: $\pi(g_{\mathbf{W}}) = g_{\mathbf{W}} + \text{residual}$
- Update of full-precision:

$$\mathbf{W}^{t+1} = \mathbf{W}^t - \alpha \pi(g_{\mathbf{W}}^t) \quad (4)$$

Overflow of MetaQuant

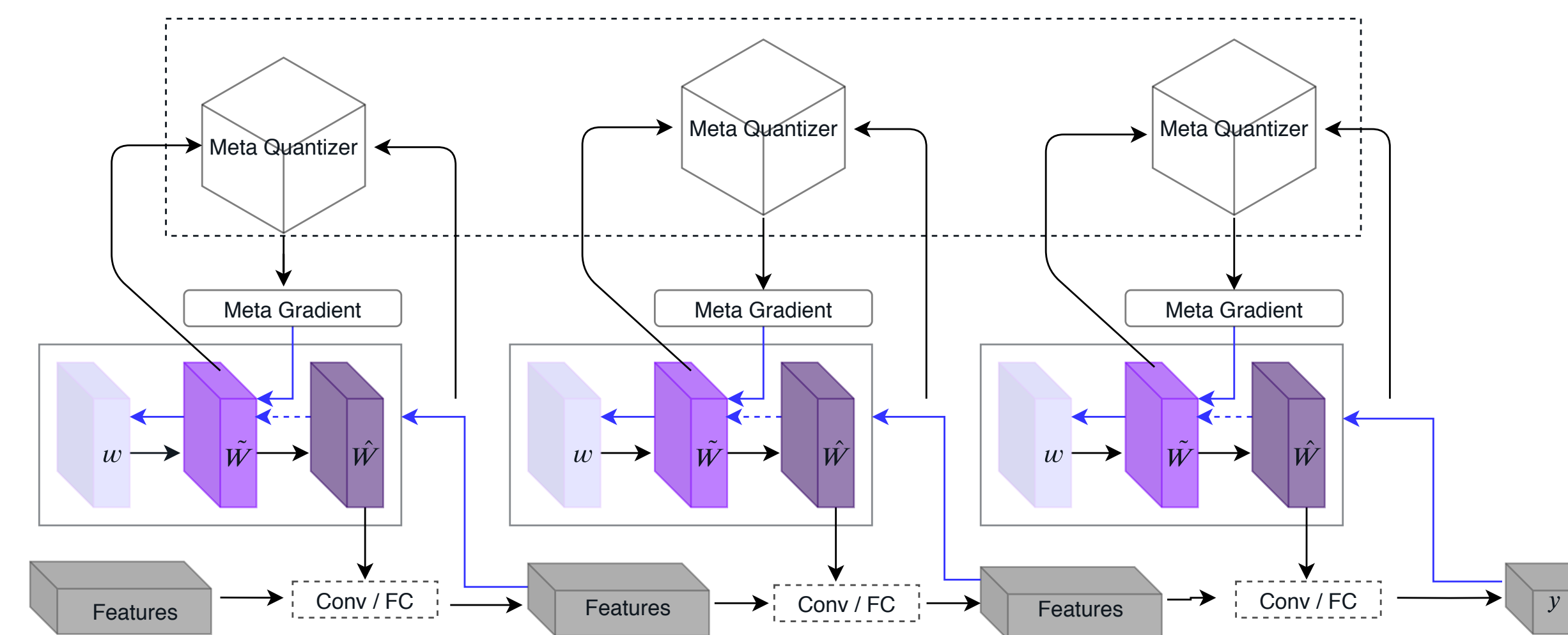


Fig. 1: Overflow of MetaQuant.

Training of Meta Quantizer

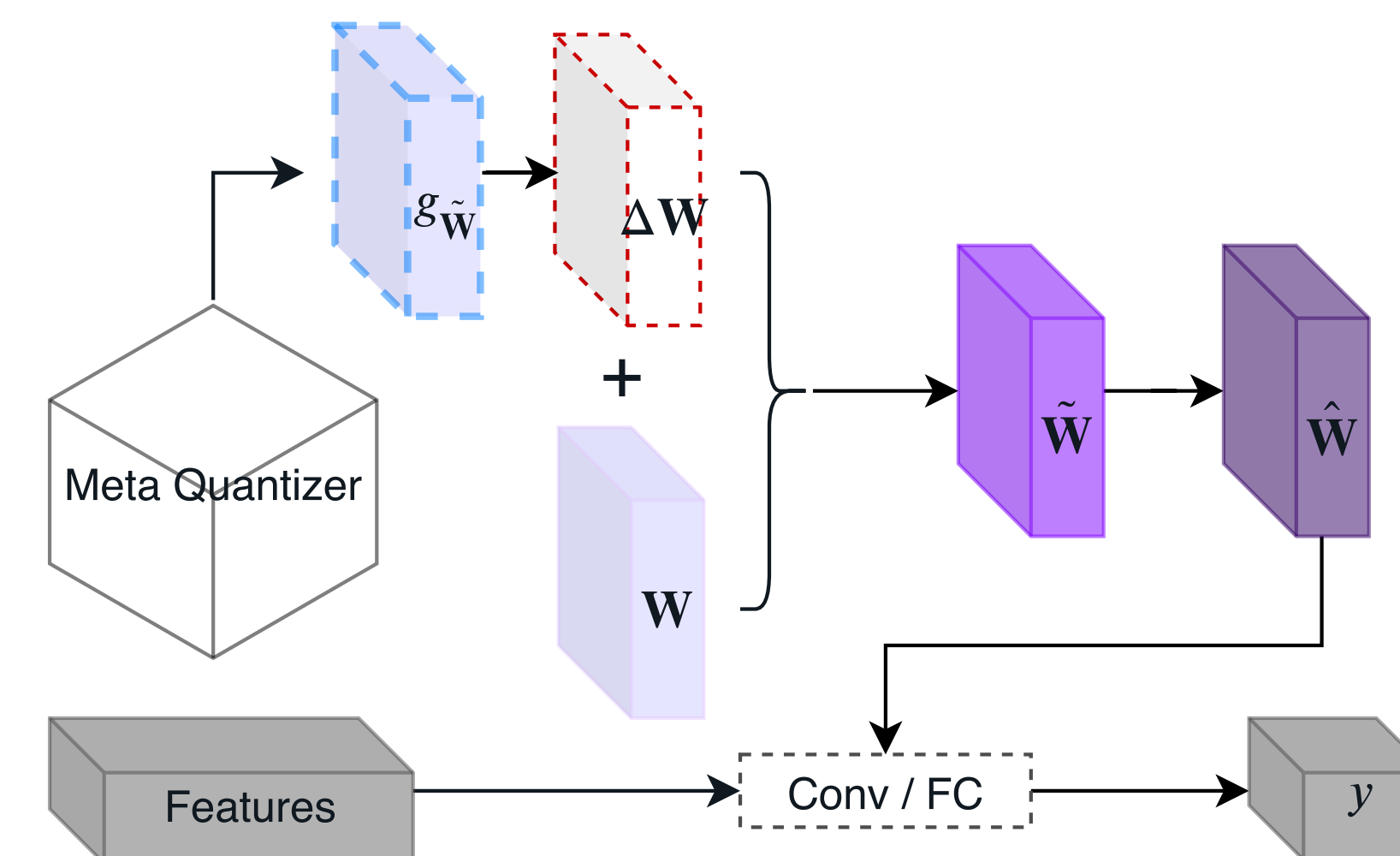


Fig. 2: Incorporation of meta quantizer into quantization training. Red dash box is composed of calibration, gradient refinement and multiplication of learning rate α . Output of meta quantizer is involved in \mathbf{W} 's update and contributes to final loss, constructing a differential path from loss to ϕ -parameterized meta quantizer.

- Meta quantizer is a **coordinate-wise** neural network: each weight in the base model is processed independently.
- During each inference, the inputs in (4) are arranged as batches with size 1: \mathbf{W} comes from a convolution layer with shape $\mathbb{R}^{o \times i \times k \times k} \rightarrow \mathbb{R}^{(o \times i \times k^2) \times 1}$.

$$\text{Forward: } \tilde{\mathbf{W}}^t = \mathcal{A}(\mathbf{W}^t) = \mathcal{A} \left[\mathbf{W}^{t-1} - \alpha \times \pi(\mathcal{M}_{\phi}(g_{\tilde{\mathbf{W}}}^{t-1}, \tilde{\mathbf{W}}^{t-1}) \frac{\partial \tilde{\mathbf{W}}^{t-1}}{\partial \mathbf{W}^{t-1}}) \right],$$

$$\ell = \text{Loss} \left\{ f \left[Q(\tilde{\mathbf{W}}^t); \mathbf{x} \right], y \right\}, \quad (5)$$

$$\text{Backward: } \frac{\partial \ell}{\partial \phi^t} = \frac{\partial \ell}{\partial \tilde{\mathbf{W}}^t} \frac{\partial \tilde{\mathbf{W}}^t}{\partial \phi^t} = \mathcal{M}_{\phi}(g_{\tilde{\mathbf{W}}^t}, \tilde{\mathbf{W}}^t) \frac{\partial \tilde{\mathbf{W}}^t}{\partial \phi^t}. \quad (6)$$

Design of Meta Quantizer

$$\begin{aligned} \text{FCGrad} : \quad \mathcal{M}_{\phi}(g_{\tilde{\mathbf{W}}}) &= \text{FCs}(\phi, \sigma, g_{\tilde{\mathbf{W}}}), \\ \text{MultiFC} : \quad \mathcal{M}_{\phi}(g_{\tilde{\mathbf{W}}}, \tilde{\mathbf{W}}) &= g_{\tilde{\mathbf{W}}} \cdot \text{FCs}(\phi, \sigma, \tilde{\mathbf{W}}), \\ \text{LSTMFC} : \quad \mathcal{M}_{\phi}(g_{\tilde{\mathbf{W}}}, \tilde{\mathbf{W}}) &= g_{\tilde{\mathbf{W}}} \cdot \text{FCs}(\phi_{\text{FCs}}, \sigma, (\text{LSTM}(\phi_{\text{LSTM}}, \tilde{\mathbf{W}}))). \end{aligned}$$

ϕ : parameters in meta quantizer, σ : nonlinear activation.

Experiment

Network-Dataset	Forward	Backward	Optimization	Test Acc (%)	FP Acc (%)
ResNet20-CIFAR10	dorefa	STE	SGD	80.745(2.113)	91.5
		MultiFC	SGD	88.942(0.466)	
		LSTMFC	SGD	88.305(0.810)	
	BWN	STE	SGD	88.840(0.291)	
		MultiFC	SGD	89.782(0.172)	
		LSTMFC	SGD	89.941(0.068)	
ResNet56-CIFAR100	dorefa	STE	SGD	89.979(0.103)	71.22
		MultiFC	SGD	89.962(0.068)	
		LSTMFC	SGD	89.962(0.068)	
	BWN	STE	SGD	75.913(3.495)	
		MultiFC	SGD	89.289(0.212)	
		LSTMFC	SGD	88.949(0.231)	
ResNet18-ImageNet	dorefa	STE	SGD	42.265(8.143)	69.76/89.08
		MultiFC	SGD	65.791(0.415)	
		LSTMFC	SGD	63.645(2.183)	
	BWN	STE	SGD	64.351(0.935)	
		MultiFC	SGD	66.419(0.533)	
		LSTMFC	SGD	66.483(0.793)	

Fig. 3: Overall experiments of MetaQuant.

Network	Method	Acc Drop (%)	Network	Method	Acc Drop (%)
ResNet20	ProxQuant	1.29	ResNet32	ProxQuant	1.28
ResNet20	MetaQuant	0.7	ResNet32	MetaQuant	0.39
ResNet44	ProxQuant	0.99	LABNet	LAB	1.4
ResNet44	MetaQuant	0.08	LABNet	MetaQuant	-0.2
ResNet18	ELQ	3.55/2.65	ResNet18-2bits	TTQ	3.00/2.00
ResNet18	MetaQuant	6.32/4.31	ResNet18-2bits	MetaQuant	5.17/3.59

Fig. 4: Experimental result of MetaQuant V.S Non-STE training-based quantization: ProxQuant, LAB, ELQ, TTQ.

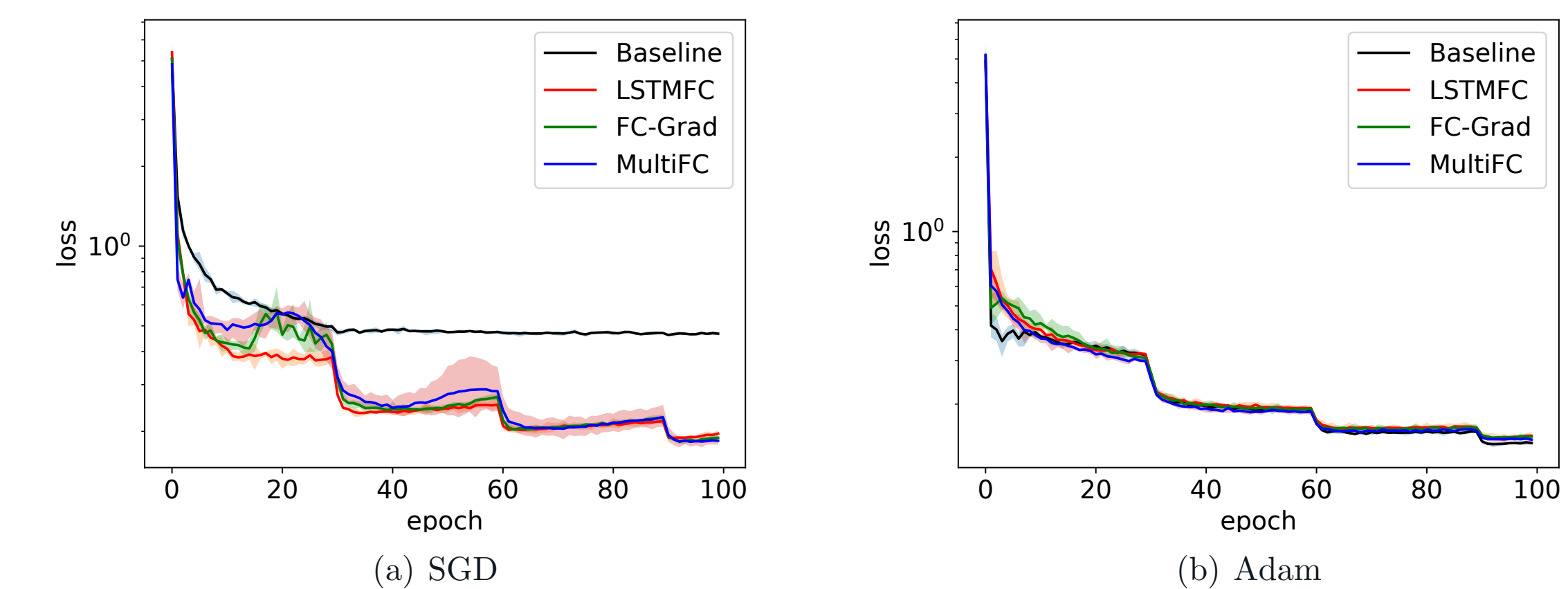


Fig. 5: Convergence analysis of MetaQuant using ResNet20, CIFAR10, dorefa, SGD/Adam.

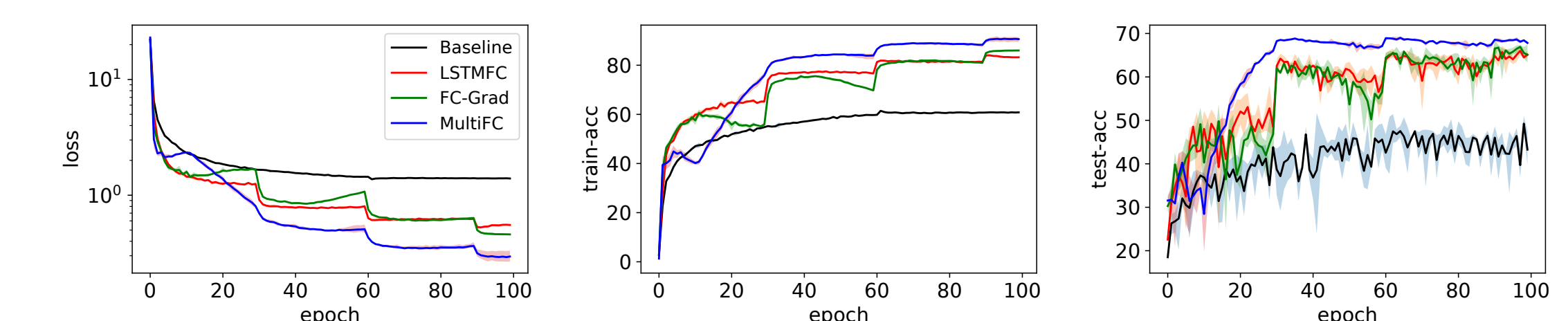


Fig. 6: Overall training of MetaQuant using ResNet110, CIFAR100, dorefa, SGD.