



## Motivation

Most prevailing low bits compression methods relied on a heavy training process with large amount of labeled data.

## Contribution

We propose a new layer-wise quantization method for limited training data scenario:

1. Layer-wise parameters are quantized using closed-form solution.
2. Preserves performance with theoretical guarantee.
3. A small portion of training data (1% in experiments) is required in the whole process.

## Cascade Layer-Wise Quantization

### Notations:

- Quantized weights:  $\hat{\Theta}$
- Un-Quantized weights:  $\bar{\Theta}$ ,  $\bar{\Theta}^{new}$
- Inputs to the  $l$ -th quantized layer:  $\hat{Y}^{l-1} = f(Y^0; \hat{\Theta}_{[1,...,l-1]})$
- Origin inputs:  $Y^{l-1}$

### Goal:

Divergence of final layer output before and after quantization is minimized:

$$\min_{\hat{\Theta}_{[l,...,L]}} \|f(\hat{Y}^{l-1}; \hat{\Theta}_{[l,...,L]}) - f(Y^{l-1}; \bar{\Theta}_{[l,...,L]})\|_F^2,$$

$$\text{s.t. } \hat{\Theta}_{[l,...,L]} \in \Omega_{[l,...,L]}.$$

### Approximation:

$$\|f(\hat{Y}^{l-1}; \hat{\Theta}_{[l,...,L]}) - f(Y^{l-1}; \bar{\Theta}_{[l,...,L]})\|_F^2$$

$$\leq \underbrace{\|f(\hat{Y}^{l-1}; \hat{\Theta}_{[l,...,L]}) - f(\hat{Y}^{l-1}; \bar{\Theta}_{[l,...,L]}^{new})\|_F^2}_{\text{Quantization}}$$

$$+ \underbrace{\|f(\hat{Y}^{l-1}; \bar{\Theta}_{[l,...,L]}^{new}) - f(Y^{l-1}; \bar{\Theta}_{[l,...,L]})\|_F^2}_{\text{Weights Update}}$$

## Quantization

- Layer output with quantized weights:  $\hat{Z}^l = \hat{Y}^{l-1} \cdot \hat{\Theta}_l$
- Layer output with origin weights:  $Z_*^l = \hat{Y}^{l-1} \cdot \bar{\Theta}_l^{new}$
- $\delta\Theta_l = \hat{\Theta}_l - \bar{\Theta}_l^{new}$

### Error Function:

$$E^l = E(\hat{Z}^l) = \frac{1}{n} \|\hat{Z}^l - Z_*^l\|_F^2$$

$$= \left( \frac{\partial E^l}{\partial \Theta_l} \right)^\top \delta\Theta_l + \frac{1}{2} \delta\Theta_l^\top H_l \delta\Theta_l + O(\|\delta\Theta_l\|_2^3)$$

### Optimization Objective:

$$\min_{\hat{\Theta}_l} f(\hat{\Theta}_l) = \frac{1}{2} (\hat{\Theta}_l - \bar{\Theta}_l^{new})^\top H_l (\hat{\Theta}_l - \bar{\Theta}_l^{new}),$$

$$\text{s.t. } \hat{\Theta}_l \in \Omega_l,$$

## Quantization with ADMM

By introducing an auxiliary parameter  $G$  to relax  $\hat{\Theta}$  to continuous:

$$L_\rho(\hat{\Theta}, G, \lambda)$$

$$= f(\hat{\Theta}) + I_\Omega(G) + \frac{\rho}{2} \|\hat{\Theta} - G + \lambda\|_2^2 - \frac{\rho}{2} \|\lambda\|_2^2.$$

- **Proximal Step** optimize  $\hat{\Theta}$ :

$$(\mathbf{H} + \text{diag}(\rho)) \hat{\Theta}^{k+1} = \mathbf{H} \bar{\Theta}^{new} + \text{diag}(\rho)(G^k - \lambda^k).$$

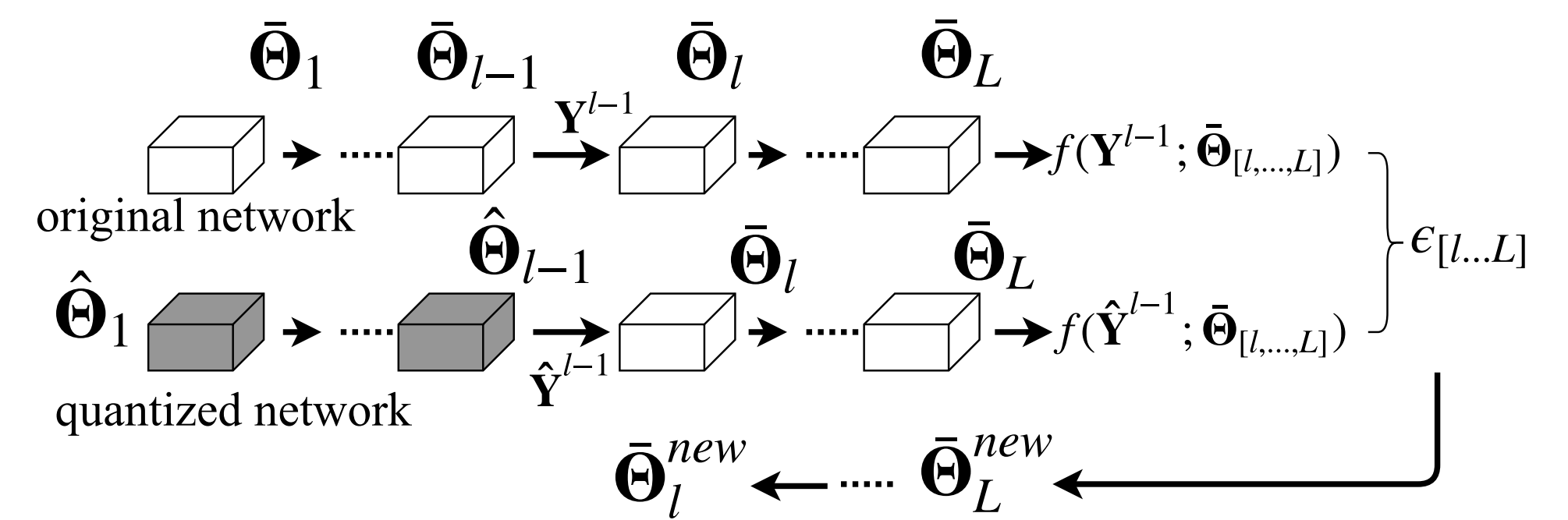
- **Projection Step** optimize  $G$ :

$$\min_G \|\hat{\Theta}^{k+1} - G + \lambda^k\|_2^2, \text{ s.t. } G \in \Omega.$$

- **Dual Update Step** optimize  $\lambda$ :

$$\lambda^{k+1} = \lambda^k + \hat{\Theta}^{k+1} - G^{k+1}.$$

## Remaining Non-quantized Weights Update



## Experiments

- Using only 1% of CIFAR10 and ImageNet dataset:

Dataset	Network	Method	bits	Improve(%)	Full-Precision
CIFAR10	ResNet20	TTQ	3	-77.25	91.77
		INQ	15	-48.48	90.02
		ExNN	3	-11.15	91.5
		VQ	3	-11.27	
		DQ	3	-19.92	
		L-DNQ	3	<b>-4.30</b>	
ImageNet	ResNet18	TTQ	3	-69.48/-88.49	69.6/89.2
		INQ	15	-61.27/-64.22	68.27/88.69
		ExNN	3	-43.53/-37.82	69.76/89.02
		VQ	3	-35.69/-29.08	
		DQ	3	-61.22/-65.64	
		L-DNQ	3	<b>-16.43/-10.67</b>	
		DQ	8	-56.78/-58.92	
			16	-13.81/-8.68	
			32	-2.82/-1.51	
		L-DNQ	9	<b>-2.73/-0.90</b>	

- Performance as number of instances climbs:

