# Leveraging Inter-Layer Dependency for Post-Training Quantization

Changbao Wang    Dandan Zheng    Yuanliu Liu    Liang Li
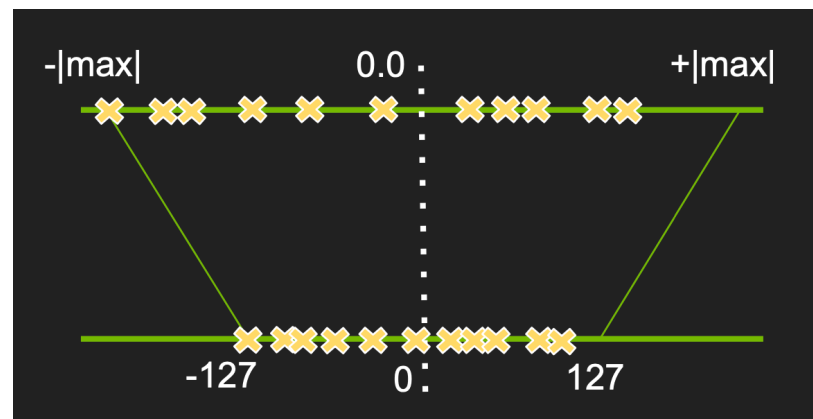
Ant Technology Group Co., Ltd.

蚂蚁集团 ANT GROUP | 支付宝

# Model Quantization

$$x_q = clip(\lfloor\frac{x}{s}\rceil, q_-, q_+), x \in \mathbb{R}^D, s \in \mathbb{R}, q_-, q_+ \in \mathbb{Z}$$

$$err_q = |x_q * s - x|$$

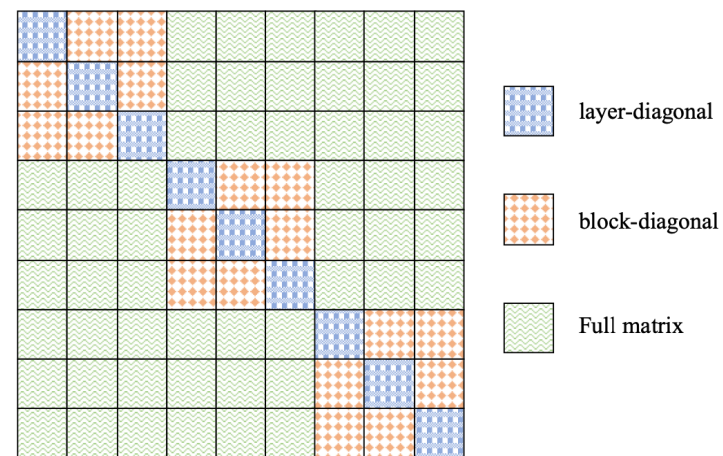# Previous Post-Training Quantization Approaches

$$w_q = clip(\lfloor \frac{w}{s} \rceil, q_-, q_+)$$

ceil?     floor?

$$\underset{\mathbf{V}}{\arg\min} \quad \left\| \mathbf{W}\mathbf{x} - \widetilde{\mathbf{W}}\mathbf{x} \right\|_F^2 + \lambda f_{reg}(\mathbf{V})$$

[AdaRound by Nagel et al. ICML 2020]
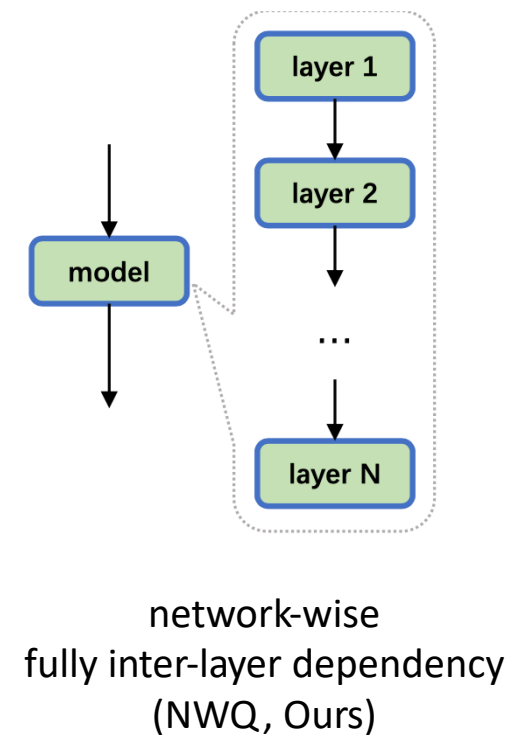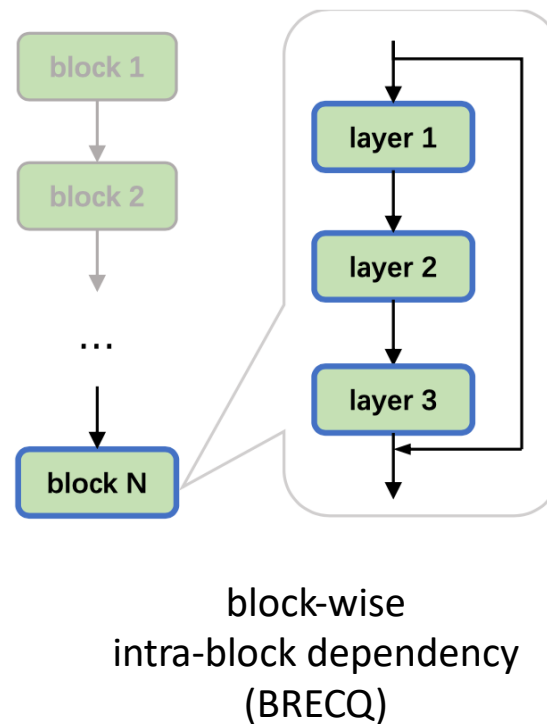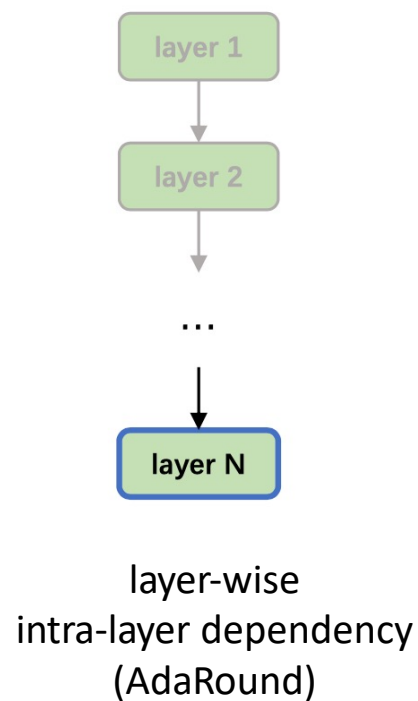


layer-diagonal

block-diagonal

Full matrix

$$\underset{\hat{\mathbf{w}}}{\min} \mathbb{E}\left[ \Delta\mathbf{z}^{(\ell),\mathsf{T}} \mathbf{H}^{(\mathbf{z}^{(\ell)})} \Delta\mathbf{z}^{(\ell)} \right] = \underset{\hat{\mathbf{w}}}{\min} \mathbb{E}\left[ \Delta\mathbf{z}^{(\ell),\mathsf{T}} \mathrm{diag}\left( (\frac{\partial L}{\partial\mathbf{z}_1^{(\ell)}})^2, \ldots, (\frac{\partial L}{\partial\mathbf{z}_a^{(\ell)}})^2 \right) \Delta\mathbf{z}^{(\ell)} \right]$$
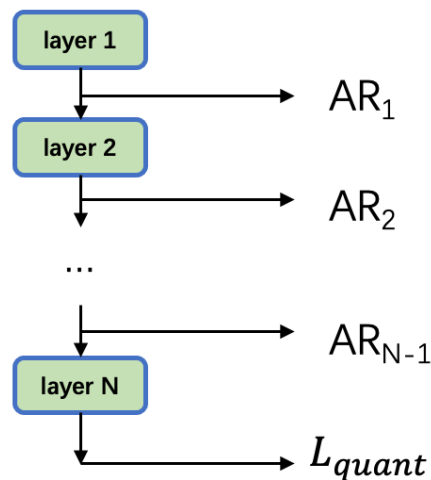
[BRECQ by Li et al. ICLR 2021]

# Network-Wise Quantization



layer-wise
intra-layer dependency
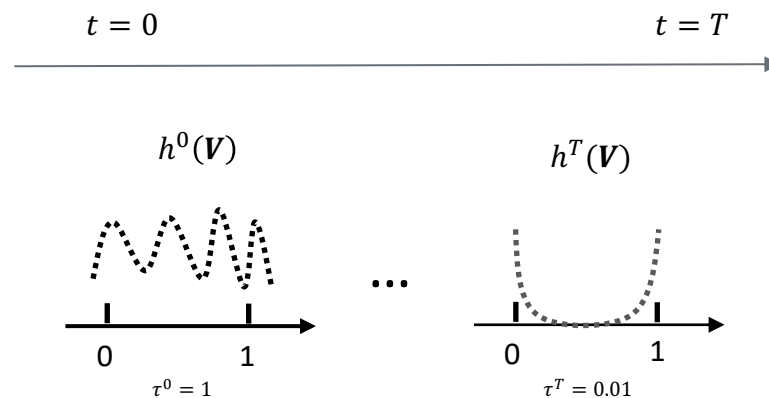(AdaRound)

block-wise
intra-block dependency
(BRECQ)

network-wise
fully inter-layer dependency
(NWQ, Ours)

Challenges of Naïve NWQ
- Higher risk of overfitting
- More difficult of discrete optimization

# Our Approaches

layer 1 → $AR_1$

layer 2 → $AR_2$

... → $AR_{N-1}$

layer N → $L_{quant}$

$t = 0$          $t = T$

$h^0(\boldsymbol{V})$          $h^T(\boldsymbol{V})$

0          1          0          1

$\tau^0 = 1$          $\tau^T = 0.01$

$t = 0$          $t = T$

$p(0) = 0.5$          $p(T) = 0$

$$loss = L_{quant} + \sum_{i=1}^{N-1} AR_i(x_i, \widehat{x_i})$$

$$x_q = clip\left(\left\lfloor\left|\frac{w}{s}\right|\right\rceil\right) + h^t(\boldsymbol{V}), q_-, q_+),$$

$$h^t(\boldsymbol{V}) = softmax(\frac{\boldsymbol{V}}{\tau^t})$$

$$\hat{x} = where(randn() < p(t), x, x_q)$$

Activation Regularization          Annealing Softmax          Annealing Mixup

# Overview

# Comparing with previous works



Accuracy of W2A2