

---

# Toward Extremely Low Bit and Lossless Accuracy in DNNs with Progressive ADMM

---

Sheng Lin<sup>1</sup> Xiaolong Ma<sup>1</sup> Shaokai Ye<sup>2</sup> Geng Yuan<sup>1</sup> Kaisheng Ma<sup>2</sup> Yanzhi Wang<sup>1</sup>

## Abstract

Weight quantization is one of the most important techniques of Deep Neural Networks (DNNs) model compression method. A recent work using systematic framework of DNN weight quantization with the advanced optimization algorithm ADMM (Alternating Direction Methods of Multipliers) achieves one of state-of-art results in weight quantization. In this work, we first extend such ADMM-based framework to guarantee solution feasibility and we have further developed a multi-step, progressive DNN weight quantization framework, with dual benefits of (i) achieving further weight quantization thanks to the special property of ADMM regularization, and (ii) reducing the search space within each step. Extensive experimental results demonstrate the superior performance compared with prior work. Some highlights: we derive the first lossless and fully binarized (for all layers) LeNet-5 for MNIST; And we derive the first fully binarized (for all layers) VGG-16 for CIFAR-10 and ResNet for ImageNet with reasonable accuracy loss. Our models and sample codes are released in anonymous link <http://bit.ly/2YYqzJv>.

## 1. Introduction

With the development of machine learning technologies, Deep Neural Networks (DNNs) have shown their extraordinary performance for their high accuracy and excellent scalability (Krizhevsky et al., 2012). However, DNNs are suffering from both intensive computation and huge stor-

age. A number of prior work have focused on developing *model compression* techniques for DNNs. These techniques, which are applied during the training phase of the DNN, aim to simultaneously reduce the model size and accelerate the computation for inference – all these to be achieved with non-negligible classification accuracy loss. Indeed the accuracy of a DNN inference engine after model compression is typically higher than that of a shallow neural network with no compression (Han et al., 2015; Wen et al., 2016). One of the most important categories of DNN model compression techniques is *weight quantization*.

We have investigated weight quantization of DNNs in many recent work (Leng et al., 2017; Park et al., 2017; Zhou et al., 2017; Lin et al., 2016; Wu et al., 2016; Rastegari et al., 2016; Hubara et al., 2016; Courbariaux et al., 2015). In these work, both storage and computational requirements of DNNs have been greatly reduced with tolerable accuracy loss. We know that multiplication operations are costly and it can be eliminated when applying binary, ternary, or power-of-2 weight quantizations (Rastegari et al., 2016; Hubara et al., 2016; Courbariaux et al., 2015).

To overcome the limitation of the highly heuristic nature in prior work, a recent work (Leng et al., 2017) developed a systematic framework of DNN weight quantization using the advanced optimization technique ADMM (Boyd et al., 2011; Hong et al., 2016). Through the adoption of ADMM, the original weight quantization problem is decomposed into two sub-problems, one effectively solved using stochastic gradient descent as original DNN training, while the other solved optimally and analytically via Euclidean projection. This method achieves one of state-of-art in weight quantization results. However, the direct application of ADMM technique lacks the guarantee on solution feasibility due to the non-convex nature of objective function (loss function), while there is also margin of improvement for solution quality.

In this work, we first make the following extensions on the ADMM-based weight compression (Zhang et al., 2018): (i) develop an integrated framework of dynamic ADMM regularization and quantized weight projection, thereby guaranteeing solution feasibility and providing high solution quality; (ii) incorporate the multi- $\rho$  updating technique for

<sup>1</sup>Northeastern University, Boston <sup>2</sup>Tsinghua University, Beijing.  
Correspondence to: Sheng Lin <lin.sheng@husky.neu.edu>, Xiaolong Ma <ma.xiaol@husky.neu.edu>, Shaokai Ye <shaokaiyeah@gmail.com>, Geng Yuan <yuan.geng@husky.neu.edu>, Kaisheng Ma <kaisheng@mail.tsinghua.edu.cn>, Yanzhi Wang <yanz.wang@northeastern.edu>.

faster and better ADMM convergence.

Extensive experimental results demonstrate that the proposed progressive framework consistently outperforms prior work. Some highlights: we derive the first lossless, fully binarized (for all layers) LeNet-5 for MNIST; and we derive fully binarized (for all layers) VGG-16 model for CIFAR-10 and ResNet model for ImageNet with reasonable accuracy loss.

## 2. DNN Model Compression

In this section, we give a detailed description to achieve a good quantization result for DNNs with progressive ADMM.

### 2.1. Framework Design

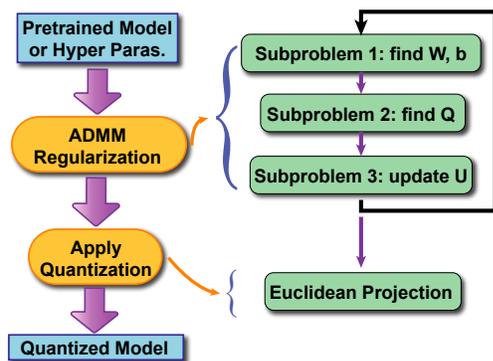


Figure 1. Illustration of one step progressive DNN model quantization.

The ADMM-based weight quantization is performed multiple times, each as a step in the progressive framework. Figure 1 illustrates one step of proposed progressive DNN weight quantization framework. The quantization result from the previous step is evaluated with current quantization result, and serve as intermediate results and starting point for the subsequent step if it is better than current result. The reason to develop a progressive model compression framework is that the multi-step procedure reduces the search space for weight quantization within each step.

Through extensive investigations, we conclude that the progressive comparison will be in general sufficient for weight quantization, in which each step requires approximately the same number of training epochs. And during the process, we adjust the penalty factor of ADMM to speed up the convergence and achieve better quantized result. The specific ADMM optimization process is introduced in following subsection.

### 2.2. ADMM-based Weight Quantization

ADMM(Boyd et al., 2011) is an advanced optimization technique which decompose an original problem into sub-problems that can be solved separately and iteratively. By adopting ADMM regularized optimization, the framework can provide high solution quality and with no obvious accuracy degradation.

First, the *progressive DNN weight quantization* starts from a pre-trained full size DNN model without compression. Consider an  $N$ -layer DNNs, sets of weights of the  $i$ -th (CONV or FC) layer are denoted by  $\mathbf{W}_i$ , respectively. And the *loss function* associated with the DNN is denoted by  $f(\{\mathbf{W}_i\}_{i=1}^N)$ . In this paper,  $\{\mathbf{W}_i\}_{i=1}^N$  characterize the set of weights from layer 1 to layer  $N$ . The overall weight quantization problem is defined by

$$\begin{aligned} & \underset{\{\mathbf{W}_i\}}{\text{minimize}} && f(\{\mathbf{W}_i\}_{i=1}^N), \\ & \text{subject to} && \mathbf{W}_i \in \mathcal{Q}_i, i = 1, \dots, N. \end{aligned} \quad (1)$$

For weight quantization, elements in  $\mathcal{Q}_i$  are the solutions of  $\mathbf{W}_i$ . Assume set of  $q_{i,1}, q_{i,2}, \dots, q_{i,M_i}$  is the available quantized values, where  $M_i$  denotes the number of available quantization level in layer  $i$ . Suppose  $q_{i,j}$  indicates the  $j$ -th quantization level in layer  $i$ , which gives

$$q_{i,j} \in \{-\alpha_i, \alpha_i\} \text{ or } \{-\alpha_i, 0, \alpha_i\}. \quad (2)$$

And  $\alpha_i$  is the scaling factor, which is initialized by the average of weight values in layer  $i$ .

Then the original problem (1) can be equivalently rewritten as

$$\begin{aligned} & \underset{\{\mathbf{W}_i\}}{\text{minimize}} && f(\{\mathbf{W}_i\}_{i=1}^N) + \sum_{i=1}^N h_i(\mathbf{Q}_i), \\ & \text{subject to} && \mathbf{W}_i = \mathbf{Q}_i, i = 1, \dots, N. \end{aligned} \quad (3)$$

We incorporate auxiliary variables  $\mathbf{Q}_i$ , dual variables  $\mathbf{U}_i$ , then apply ADMM to decompose problem (3) into simpler subproblems. Then solve these subproblems iteratively until convergence. The augmented Lagrangian formation of problem (3) is

$$\underset{\{\mathbf{W}_i\}}{\text{minimize}} f(\{\mathbf{W}_i\}_{i=1}^N) + \sum_{i=1}^N \frac{\rho_i}{2} \|\mathbf{W}_i - \mathbf{Q}_i + \mathbf{U}_i\|_F^2 \quad (4)$$

The first term in problem (4) is the differentiable loss function of the DNN, and the second term is a quadratic regularization term of the  $\mathbf{W}_i$ , which is differentiable and convex. As a result, subproblem (4) can be solved by stochastic gradient descent algorithm (Kingma & Ba, 2014) as the original DNN training.

The standard ADMM algorithm (Boyd et al., 2011) steps proceed by repeating, for  $k = 0, 1, \dots$ , the following sub-problems iterations:

$$\mathbf{W}_i^{k+1} := \arg \min_{\mathbf{W}_i} L_P(\{\mathbf{W}_i\}, \{\mathbf{Q}_i^k\}, \{\mathbf{U}_i^k\}), \quad (5)$$

$$\mathbf{Q}_i^{k+1} := \arg \min_{\mathbf{Q}_i} L_P(\{\mathbf{W}_i^{k+1}\}, \{\mathbf{Q}_i\}, \{\mathbf{U}_i^k\}), \quad (6)$$

$$\mathbf{U}_i^{k+1} := \mathbf{U}_i^k + \mathbf{W}_i^{k+1} - \mathbf{Q}_i^{k+1}. \quad (7)$$

which (5) is the proximal step, (6) is projection step and (7) is dual variables update.

### 3. Experimental Results

**Binary Weight Quantization Results on LeNet-5:** To the extent of authors’ knowledge, we achieve the first lossless, fully binarized LeNet-5 model in which weights all layers are binarized. The accuracy is still 99.21%, lossless compared with baseline. For example, recent work (Cheng et al., 2018) results in 2.3% accuracy degradation on MNIST for full binarization, with baseline accuracy 98.66%.

Table 1. Comparisons of fully binary weight quantization results on LeNet-5 for MNIST dataset.

Method	Accuracy	Num. of bits
Baseline (Cheng et al., 2018)	98.66%	32
Binary (Cheng et al., 2018)	96.34%	1
<b>Our binary</b>	99.21%	1

**Weight Quantization on CIFAR-10:** We also achieve fully binarized VGG-16 for CIFAR-10 with negligible loss in accuracy, in which weights all layers (including the first and the last) are binarized. The accuracy is 93.58%. We would like to point out that fully ternarized quantization results in 94.02% accuracy. Table 2 shows our results and comparisons.

Table 2. Comparisons of fully binary (ternary) weight quantization results on VGG-16 for CIFAR-10 dataset.

Method	Accuracy	Num. of bits
Baseline (Cheng et al., 2018)	84.80%	32
8-bit (Cheng et al., 2018)	84.07%	8
Binary (Cheng et al., 2018)	81.56%	1
<b>Our baseline</b>	94.70 %	32
<b>Our ternary</b>	94.02%	2
<b>Our binary</b>	93.58%	1

**Binary Weight Quantization Results on ResNet for ImageNet Dataset:** The binarization of ResNet models on ImageNet data set is widely acknowledged as a very challenging task. As a result, there are very limited prior work (e.g., the

one-shot ADMM (Leng et al., 2017)) with binarization results on ResNet models. As (Leng et al., 2017) targets ResNet-18 (which is even more challenging than ResNet-50 or larger ones), we make a fair comparison on the same model. Table 3 demonstrates the comparison results (Top-5 accuracy loss). In prior work, it is by default that the first and last layers are not quantified (or quantized to 8 bits) as these layers have a significant effect on overall accuracy. When leaving the first and last layers unquantized, our framework is not progressive, but an extended one-shot ADMM-based framework. We can observe the higher accuracy compared with the prior method under this circumstance (first and last layers unquantized while the rest of layers binarized). The Top-1 accuracy has similar result: 3.8% degradation in our extended one-shot and 4.3% in (Leng et al., 2017).

Table 3. Comparisons of weight quantization results on ResNet-18 for ImageNet dataset.

Method	Relative Top-5 acc. loss	Num. of bits
Uncompressed	0.0%	32
One-shot ADMM quantization (Leng et al., 2017)	2.9%	1 (32 for the first and last)
<b>Our method</b>	2.5%	1 (32 for the first and last)
<b>Our method</b>	5.8%	1

Using the progressive framework, we can derive a fully binarized ResNet-18, in which weights in all layers are binarized. The accuracy degradation is 5.8%, which is noticeable and shows that the full binarization of ResNet is a challenging task even under the progressive framework. We did not find prior work for comparison on this result.

### 4. Conclusion and On-going Work

In this work, we extended the prior ADMM-based framework and developed a multi-step, progressive DNN weight quantization framework, in which we achieve further weight quantization results and provide better convergence rate.

Considering the good performance of our method, we plan to test more different networks for ImageNet dataset. And we are working on testing our method for different applications and datasets.

### References

- Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J., et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- Cheng, H.-P., Huang, Y., Guo, X., Huang, Y., Yan, F., Li,

- H., and Chen, Y. Differentiable fine-grained quantization for deep neural network compression. *arXiv preprint arXiv:1810.10351*, 2018.
- Courbariaux, M., Bengio, Y., and David, J.-P. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in neural information processing systems*, pp. 3123–3131, 2015.
- Han, S., Pool, J., Tran, J., and Dally, W. Learning both weights and connections for efficient neural network. In *Advances in neural information processing systems*, pp. 1135–1143, 2015.
- Hong, M., Luo, Z.-Q., and Razaviyayn, M. Convergence analysis of alternating direction method of multipliers for a family of nonconvex problems. *SIAM Journal on Optimization*, 26(1):337–364, 2016.
- Hubara, I., Courbariaux, M., Soudry, D., El-Yaniv, R., and Bengio, Y. Binarized neural networks. In *Advances in neural information processing systems*, pp. 4107–4115, 2016.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- Leng, C., Li, H., Zhu, S., and Jin, R. Extremely low bit neural network: Squeeze the last bit out with admm. *arXiv preprint arXiv:1707.09870*, 2017.
- Lin, D., Talathi, S., and Annapureddy, S. Fixed point quantization of deep convolutional networks. In *International Conference on Machine Learning*, pp. 2849–2858, 2016.
- Park, E., Ahn, J., and Yoo, S. Weighted-entropy-based quantization for deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- Rastegari, M., Ordonez, V., Redmon, J., and Farhadi, A. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European Conference on Computer Vision*, pp. 525–542. Springer, 2016.
- Wen, W., Wu, C., Wang, Y., Chen, Y., and Li, H. Learning structured sparsity in deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 2074–2082, 2016.
- Wu, J., Leng, C., Wang, Y., Hu, Q., and Cheng, J. Quantized convolutional neural networks for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4820–4828, 2016.
- Zhang, T., Ye, S., Zhang, K., Tang, J., Wen, W., Fardad, M., and Wang, Y. A systematic dnn weight pruning framework using alternating direction method of multipliers. *European Conference on Computer Vision (ECCV)*, 2018.
- Zhou, A., Yao, A., Guo, Y., Xu, L., and Chen, Y. Incremental network quantization: Towards lossless cnns with low-precision weights. *arXiv preprint arXiv:1702.03044*, 2017.