



- inference frameworks.
- get further accuracy improvement.



$$m_i = l + (i + 0.5)\Delta \text{ and } s = \frac{1}{\tanh(0.5k\Delta)}.$$
(3)

$$\varphi(x) = s \tanh\left(k(x - m_i)\right), \quad \text{if } x \in \mathcal{P}_i, \tag{4}$$

$$Q_S(x) = \begin{cases} l, & x < l, \\ u, & x > u, \\ l + \Delta \left(i + \frac{\varphi(x) + 1}{2} \right), & x \in \mathcal{P}_i \end{cases}$$
(5)

Differentiable Soft Quantization: Bridging Full-Precision and Low-Bit Neural Networks

Ruihao Gong^{1,2}, Xianglong Liu^{1*}, Shenghu Jiang², Tianxiang Li^{2,3}, Peng Hu², Jiazhen Lin², Fengwei Yu², Junjie Yan²

¹State Key Laboratory of Software Development Environment, Beihang University, ²SenseTime Group Limited, ³Beijing Institute of Technology

Table 1: Ratio of time for calling MLA (8-bit register) and SADDW

1			I
b	2	3	4
MLA/SADDW	31/1	7/1	1/1

	DSQ 2-bit	NCNN 8-bit [6]
time (ms)	551.22	935.51
* NCNN was to	stad with commit	t d263 ad 5 an 2010 3 15

x > u, (9)

Experiments

(11)

Table 3: At	plation study or	n 2-bit uniform	n quantization.
N	lethod	Top-1 (%)	Top-5 (%)
Standard Quantization		86.63	99.35
Fixed α		86.95	99.50
Learnt α		87.25	99.49
Learnt α, l, u		88.44	99.5 0
able 4: Compa Model	arison of 1-bit of Method	quantized mod Bit-Width (W/A)	lels on CIFAR-
VGG-Small	FP	32/32	91.65
	BNN [3]	1/1	89.90
	XNOR [7]	1/1	89.80
	Ours	1/1	91.72
	FP	32/32	90.78
	DoReFa [11]	1/1	79.30
DocNat 20	DoReFa [11] Ours	1/1 1/1	79.30 84.11
ResNet-20	DoReFa [11] Ours DoReFa [11]	1/1 1/1 1/32	79.30 84.11 90.00
ResNet-20	DoReFa [11] Ours DoReFa [11] LQ-Net [10]	1/1 1/1 1/32 1/32	79.30 84.11 90.00 90.10

Table 5: Comparison of different quantized models on ImageNet.

Model	Method	Bit-Width	Accuracy (%)	
		(W/A)		
	FP	32/32	69.90	
	BWN [7]	1/32	60.80	
	HWGQ [1]	1/32	61.30	
	TWN [4]	2/32	61.80	
	Ours	1/32	63.71	
ResNet-18	PACT [2]	2/2	64.40	
	LQ-Net [10]	2/2	64.90	
	Ours	2/2	65.17	
	ABC-Net [5]	3/3	61.00	
	PACT [2]	3/3	68.10	
	LQ-Net [10]	3/3	68.20	
	Ours	3/3	68.66	
	BCGD [9]	4/4	67.36†	
	Ours	4/4	69.56 [†]	
ResNet-34	FP	32/32	73.80	
	LQ-Net [10]	2/2	69.80	
	Ours	2/2	70.02	
	ABC-Net [5]	3/3	66.70	
	LQ-Net [10]	3/3	71.90	
	Ours	3/3	72.54	
	BCGD [9]	4/4	70.81 [†]	
	Ours	4/4	72.76 [†]	
Mobile- NetV2	FP	32/32	71.87	
	PACT [2, 8]	4/4	61.40	
	Ours	4/4	64.80	

* The † represents the results of full quantization for activations and weights across all convolution layers.

References

- [2] J. Choi, Z. Wang, S. Venkataramani, P. I.-J. Chuang, V. Srinivasan, and K. Gopalakrishnan. Pact: Parameterized clipping activation for quantized neural networks. arXiv preprint arXiv:1805.06085, 2018.
- [4] F. Li, B. Zhang, and B. Liu. Ternary weight networks. arXiv preprint arXiv:1605.04711, 2016.
- [6] nihui et al. Ncnn. https://github.com/Tencent/ncnn, 2017.
- [8] K. Wang, Z. Liu, Y. Lin, J. Lin, and S. Han. Haq: Hardware-aware automated quantization. arXiv preprint arXiv:1811.08886, 2018.

- [11] S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen, and Y. Zou. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. arXiv preprint arXiv:1606.06160, 2016.

[1] Z. Cai, X. He, J. Sun, and N. Vasconcelos. Deep learning with low precision by half-wave gaussian quantization. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Jul 2017. [3] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio. Binarized neural networks. In Advances in Neural Information Processing Systems 29, pages 4107–4115. Curran Associates, Inc., 2016.

[5] X. Lin, C. Zhao, and W. Pan. Towards accurate binary convolutional neural network. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems 30, pages 345–353. Curran Associates, Inc., 2017.

[7] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. *Lecture Notes in Computer Science*, page 525–542, 2016.

[9] P. Yin, S. Zhang, J. Lyu, S. Osher, Y. Qi, and J. Xin. Blended coarse gradient descent for full quantization of deep neural networks. *Research in the Mathematical Sciences*, 6(1):14, 2019.

[10] D. Zhang, J. Yang, D. Ye, and G. Hua. Lq-nets: Learned quantization for highly accurate and compact deep neural networks. In *The European Conference on Computer Vision (ECCV)*, September 2018.