MEMORIZATION CAPACITY OF DEEP NEURAL NETWORKS UNDER PARAMETER QUANTIZATION

Yoonho Boo, Sungho Shin and Wonyong Sung

Department of Electrical and Computer Engineering Seoul National University Neural Processing Research Center Gwanak-gu, Seoul, Korea {dnsgh337, ssh9919, wysung}@snu.ac.kr

ABSTRACT

Most deep neural networks (DNNs) require complex models to achieve high performance. Parameter quantization is widely used for reducing the implementation complexities. Previous studies on quantization were mostly based on extensive simulation using training data on a specific model. We choose a different approach and attempt to measure the per-parameter capacity of DNN models and interpret the results to obtain insights on optimum quantization of parameters. This research uses artificially generated data and generic forms of fully connected DNNs, convolutional neural networks, and recurrent neural networks. We conduct memorization and classification tests to study the effects of the number and precision of the parameters on the performance. The model and the per-parameter capacities are assessed by measuring the mutual information between the input and the classified output. To get insight for parameter quantization when performing real tasks, the training and test performances are compared.

Index Terms— Deep neural network, parameter quantization, network capacity

1. INTRODUCTION

Deep neural networks (DNNs) have achieved impressive performance on various machine learning tasks. Several DNN architectures are known, and the most famous ones are fully connected DNNs (FCDNNs), convolutional neural networks (CNNs), and recurrent neural networks (RNNs). It is known that neural networks do not need full floating-point precision for inference [1, 2, 3]. A 32-bit floating-point parameter can be reduced to 8-bit, 4-bit, 2-bit, or 1-bit, but this can incur performance degradation. Therefore, precision should be optimized, which is primarily conducted by extensive computer simulations using training data. This not only takes much time for optimization but also can incorrectly predict the performance in real environments when the characteristics of input data are different from the training data.

In this study, we attempt to measure the capacity of DNNs, including FCDNN, CNN, and RNN, using a memorization and classification task that applies random binary input data. The perparameter capacities of various models are estimated by measuring the mutual information between the input data and the classification output. Then, the fixed-point performances of the models are measured to determine the relation between the quantization sensitivity and the per-parameter capacity. The memorization capacity analysis results are extended to real models for performing image classification, by which the parameter quantization sensitivity is compared between memorization and generalization tasks.

The contributions of this paper are as follows.

- We experimentally measure the memorization capacity of DNNs and estimate the per-parameter capacity. The capacity per parameter is between 2.3 bits to 3.7 bits, according to the network structure, which is FCDNN, CNN, or RNN.
- We show that the performance of the quantized networks is closely related to the capacity per parameter, and FCDNNs show the most resilient quantization performance while RNNs suffer most from parameter quantization. The network size hardly effects the quantization performance when DNN models are trained to use full capacity.
- We suggest the sufficient number of bits for representing weights of neural networks, which are approximately 6 bits, 8 bits, and 10 bits for FCDNNs, CNNs, and RNNs, respectively. This estimate of the number of bits for implementing neural networks is very important considering that many accelerators are designed without any specific training data or applications.
- The study with real-models shows that neural networks are more resilient to quantization when performing generalization tasks than conducting memorization.

2. RELATED WORKS AND BACKGROUNDS

2.1. Neural network capacity

The capacity of neural networks has been studied since the early days of DNN research. Although the capacity can be defined in many ways, it is related to the learnability of networks. The capacity of networks is shown as the number of uncorrelated random samples that can be memorized [4]. A single-layer perceptron with n parameters can memorize at least 2n random samples [5]. In other words, the network can always construct a hyperplane with n parameters that divides 2n samples. Additionally, the capacity of a three-layer perceptron is proportional to the number of parameters [6]. Recently, RNNs were trained with random data to measure the capacity per parameter [7]. Our study is strongly motivated by this research, and

This work was supported by Samsung Advanced Institute of Technology through Neural Processing Research Center (NPRC) in Seoul National University. This work was also supported in part by the Brain Korea 21 Plus Project and the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIP) (No. 2018R1A2A1A05079504).

extends it to the quantization performance interpretation of generic DNN models, including FCDNN, CNN, and RNN. Recent studies have showed that neural networks have a generalization ability even if the expressive capacity of the model is sufficiently large [8, 9]. In this paper, we also discuss the effect of network quantization when performing generalization tasks.

2.2. Fixed-point deep neural networks

Early works on neural network quantization usually employed 16-bit parameters obtained by directly quantizing the floating-point numbers [1]. Recently, a retraining technique was developed to improve the performance of quantized networks [2, 3]. Retrainingbased quantization was applied to CNN and RNN models, showing superior performance compared to directly quantized ones [10, 11]. Many studies attempting extreme quantization have been published, such as 2-bit ternary [2, 12, 13], 1-bit binary weight quantization, and XNOR networks [14, 15]. Some aggressive model compression techniques also employed vector quantization or table look-up [16, 17]. However, not all CNNs show the same quantization performance. For example, AlexNet [18] shows almost the same performance with only 1-bit quantized parameters. However, the same quantization technique incurs a very severe performance loss when applied to ResNet [15]. A previous study shows that large sized networks are more resilient to severe quantization than smaller ones [19]. Theoretical works and many practical implementation optimization techniques have been studied [20, 21, 22, 23, 24]. Recent work increases the number of network parameters to preserve the performance under low-precision quantization [25]. Our works are not targeted to a specific data or model, but introduce the general understanding of parameter quantization.

3. MEMORIZATION CAPACITY MEASUREMENTS OF DNNS

3.1. Capacity measurements on a memorization task

We assess the network capacity of DNN models using a random data memorization and classification task [7]. In this task, N random binary vectors, X, are generated and each is randomly and uniformly assigned to the output label Y. The size of the binary vector depends on the DNN model. For FCDNN, the input X is a one dimensional vector whose size is determined by the hidden layer dimension. In CNN, the input needs to be a 2-D or 3-D tensor. Input samples of CNNs are generated by concatenating and reshaping random binary vectors. During the training process, the DNN is trained to correctly predict the label, which is 0 or 1, of the random input X. As the number of input data size, N, increases, the classification accuracy drops because of the limited memorization capacity. Note that the accuracy for the memorization task refers to the training performance after convergence because there is no proper test dataset for random training samples.

The capacity is measured using the mutual information, defined as a measure of the amount of information that one random variable contains about another random variable [26]. The mutual information of a trained network with N input samples is calculated as follows:

$$I(Y; \hat{Y}_{\theta}|X) = H(Y|X) - H(Y|\hat{Y}_{\theta}, X)$$

= $N\left(1 - (p \log_2 \frac{1}{p} + (1-p) \log_2 \frac{1}{(1-p)})\right),$ (1)

where p is the mean classification accuracy for all samples under trained parameter θ . If the training accuracy is 1, the model memorizes all random samples and the $I(Y; \hat{Y}_{\theta}|X)$ becomes the number of samples N. If the training accuracy is 0.5, $I(Y; \hat{Y}_{\theta}|X)$ goes to 0.

The network capacity is defined as

$$C = \max_{\alpha} I(Y; \hat{Y}_{\theta} | X).$$
⁽²⁾

The accuracy, p, may vary depending on the training method of the model. We find N and p that maximize the mutual information of the networks by iteratively training the models. This optimization employs both grid search- and Bayesian optimization-based hyperparameter tuning [27]. The optimization procedure consists of three stages. First, we try to find the largest input data size whose accuracy is slightly lower than 1. Second, we perform a grid search to determine the boundary values of the hyper-parameters. The searched hyper-parameters can include initialization, optimizer, initial learning rate, learning rate decay factor, batch size, and optimizer variables. Finally, we conduct hyper-parameter tuning within the search space using Scikit-learn library [28]. We add the number of training samples N as a hyper-parameter and use the mutual information of Eq. (1) as the metric for the optimization.

3.2. Network quantization and parameter capacity

Quantization of model parameters perturbs the trained network, therefore, fixed-point training or retraining with full-precision backpropagation is usually needed [2, 12, 14, 29]. However, the performance of the quantized networks does not always meet that of the floating-point models, even after retraining. This suggests that model capacity is reduced by quantization, especially when the number of bits used is very small.

In this research, we observe the memorization capacity degradation caused by quantization in generic FCDNN, CNN, and RNN models. The uniform quantization is used for the sake of convenient arithmetic, and the same step size is assigned to a layer or feature map in the FCDNN and CNN. We assigned two step sizes in a LSTM layer. The biases are not quantized, because they have a large dynamic range. It is important to note that the weights connected to the output are not quantized, because their optimum bit-widths depend on the number of labels in the output. Quantization is performed from floating-point to 8, 6, 5, 4, 3, and 2-bit precision, in sequence. Retraining is performed after every quantization, but requires only a small number of epochs, because only fine-tuning is needed [2].

4. EXPERIMENTAL RESULTS ON MEMORIZATION CAPACITY OF FLOATING-POINT DNNS

The capacities of FCDNNs, CNNs, and RNNs are measured via the memorization task explained in Section 3.1. The models used for the test employ floating-point parameters.

4.1. Capacity of FCDNNs

The training data for FCDNNs is a 1-D vector of size n_{in} . N input data are used as for the training data. The output, Y, is the randomly assigned label, either 0 or 1, for each input. Thus, inputs, X and Y, are represented as $X \in \{0, 1\}^{N \times n_{in}}$ and $Y \in \{0, 1\}^N$, respectively. The input data dimension, n_{in} , should be larger than $\log_2 N$ so that no overlapped data is contained among N input data. In the experiments for FCDNNs, the input vector dimension, n_{in} , is chosen to be equal to the number of units in the hidden layer.



0.9 Accuracy 0.75-bit 6-bit 10-bit 2-bit 3-bi 4-bit 8-bit float (a) Full capacity FCDNN(32)×2 0. Accuracy FCDNN(128)×3 CNN model A CNN model B LSTM(32)×1 0 LSTM(64)×1 10-bit 2-bit 3-bit 4-bit 5-bit 6-bit 8-bit float

Fig. 2. Quantization effect when networks use the (a) full capacity and (b) half capacity.

(b) Half capacity

Fig. 1. (a) Mutual information according to the number of inputs N. (b) The relationship between the number of parameters and the capacity of networks in FCDNNs and CNNs.

We conduct experiments for FCDNNs with hidden layer dimensions of 32, 64, 128, and 256, and with hidden layer depths of 1, 2, 3, and 4. For each model, experiments are conducted to measure the accuracy of memorization while increasing the size of the input data, N. Note that only the training error is measured in this memorization task, because there is no unseen data. Experimental results are based upon the best accuracy obtained when attempted with different hyper parameters.

Fig. 1(a) shows the memorization accuracy and the mutual information obtained using Eq. (1) on the FCDNN. The model is composed of three layers and the hidden layer of size 64. Here, we find that the amount of mutual information steadily increases as the input data size grows. However, it begins to drop as the input size grows farther, and the memorization accuracy drops. By analyzing the accuracy trend of the model, it is possible to distinguish the input data size into three regions: the over-parameterized, the maximum performance, and the under-parameterized sections, as shown in Fig. 1(a). For example, if the model is trained to memorize only 10,000 data, it can be regarded as over-parameterized. The number of data that can be memorized by maximally utilizing all the parameters is between 30,000 and 40,000. In over-parameterized regions, performance can be maintained, even if the capacity of the networks is reduced.

The per-parameter capacity of FCDNNs is shown in Fig. 1(b). Regardless of the width or depth, one parameter has a capacity of 1.7 to 2.5 bits, and FCDNNs have an average of 2.3-bit capacity per parameter. Note that per parameter capacity is constant regardless of the width or depth of layers. The total capacity of the model may be interpreted as an optimal storage that can store a maximum of random binary samples [5, 30].

4.2. Capacity of CNNs

The capacity of CNNs is also measured via a similar memorization task. CNNs can have a variety of structures according to the number of channels, the size of the kernels, and the number of layers. The kernel size of CNNs in this test are either (3×3) or (5×5) , which are the same for all layers and the number of convolution layers is from 3 to 9. The dimension of the inputs is $(32 \times 32 \times 1)$ for all experiments. Three max-pooling operations are applied to reduce the number of parameters in the fully connected layer. As shown in Fig. 1(b), the convolution layers have the per-parameter capacity of between 2.86 and 3.09, which is higher than that of FCDNNs. The average capacity per parameter of the tested models is 3.0 bits.

Results show that the per-parameter memorization capacity of CNNs is higher than that of FCDNNs, even when CNNs memorize uncorrelated data. Note that one parameter of FCDNNs is used only once for each inference. However, the parameter of CNNs is used multiple times. This parameter-sharing nature of CNNs seems to increase the amount of information that one parameter can store.

4.3. Capacity of RNNs

It has been shown that the various structures of RNNs all have similar capacity per parameter of 3 to 6 bits [7]. We train RNNs using a dataset with no sequence correlation to show the capacity of the parameters. The random input dataset is composed of inputs, $X \in \{0,1\}^{N \times n_{seq} \times n_{in}}$ and labels $Y \in \{0,1\}^N$, which are uniformly set to 0 or 1. The training loss is calculated using the crossentropy of the label at the output of the last step.

We train RNNs with a single LSTM layer of 32-D. The input dimension, n_{in} , is also 32-D and the amount of unrolling sequence, n_{seq} , is five-step. We apply 5 input random vectors, X_0 , X_1 , X_2 , X_3 , and X_4 , each with 32-D, and assign one label to this 160-D vector at the last time step. The error propagates from the last step only, and the outputs at intermediate time-steps are ignored. The number of parameters in the network is 8,386. In this case, the maximum mutual information is obtained when the number of samples is 32K, and the memorization accuracy is 99.52 %. Therefore, the



Fig. 3. Performance degradation according to the quantization precision on CNN trained with CIFAR-10 dataset. Solid and dashed lines denote the accuracy on the training and test dataset, respectively.

per-parameter capacity of the model is 3.7 bits. The RNN shows higher per-parameter capacity than FCDNNs and CNNs.

5. EXPERIMENTAL RESULTS OF PARAMETER QUANTIZATION

5.1. Capacity under parameter quantization

We have shown that FCDNNs, CNNs, and RNNs have different per-parameter capacities. According to the parameter-data ratio, a trained DNN can be an over-parameterized, max-capacity, or underparameterized model. Thus, we can assume that the DNN performance under quantization would depend on not only the network structure, such as FCDNN, CNN, or RNN, but also the parameterdata ratio. The experiments are divided into two cases. The first is to measure performance degradation via quantization precision when each model is in the maximum capacity region. The second analyzes performance when the models are in the over-parameterized region.

When the FCDNN, CNN, and RNN are trained to have the maximum memorization capacity, the performances with parameter quantization are shown in Fig. 2(a). The fixed-point performances of two FCDNNs, two CNNs, and two RNNs are illustrated. CNN model A and B have (5×5) and (3×3) kernels, respectively. Both models consist of 6 convolutional layers and an output fully connected layer. With 6-bit parameter quantization, the FCDNN shows no accuracy drop. However, those for CNNs and RNNs are 5 % and 18 %, respectively. Because the RNN contains the largest amount of information at each parameter, the loss caused by parameter quantization seems to be the most severe. We also find that there is no decline in performance until the parameter precision is lowered to 6-bit for FCDNNs, 8-bit for CNNs, and 10-bit for RNNs, even when all models use full capacity.

Next, we show the fixed-point performance of DNNs when they are trained to be in the over-parameterized region. Note that the per-parameter capacity is lowered in the over-parameterized region. We conducted simulation with half size of the maximum number of data that can be memorized. For example, an FCDNN used for the measurement has 3 hidden layers with a hidden-layer dimension of 128; the capacity of the corresponding model is about 2^{17} bits. The network is over-parameterized when the number of memorized samples is 2^{16} . Fig. 2(b) shows that the FCDNN model memorizes all samples even with 4-bit parameter quantization when the model uses half of the capacity. We can find that over-parmeterized models are less sensitive to bit-precision on CNNs and RNNs.

5.2. Quantization experiments on real tasks

We have assessed the required precision of networks for performing memorization tasks. The memorization test only uses the training data that are artificially generated. However, most neural networks should conduct more than memorization because the test data are not seen during the training. In this section, we analyze the effects of network quantization for performing real tasks. We train two different sized CNN models with CIFAR-10 data. The structures of the two models are as follows:

$$2 \times kC - MP2 - 2 \times 2kC - MP2 - 2 \times 4kC - MP2 \quad (3)$$
$$-2 \times 8kFC - 10 output.$$

The size of kernels is (3×3) , kC represents a convolution layer with k channels and 8kFC means a fully connected layer with 8kneurons. We experimented 3 types of CNNs and they are denoted as *Large*, *Small*, and *Tiny*, where the number of channels, k, is 32, 16, and 8, respectively.

We analyze the impact of network quantization on the test performance. Fig. 3 shows the training and test data based performance of fixed-point CNNs on real data. When applying the training data that may have been memorized during the training phase, the Large model shows indifferent performance surface regardless of the parameter precision. But, for the Small model, the 2-bit quantization results in quite degraded performance when compared to the floating-point network. However, the test accuracy of the Small 2bit model is not much lowered. The Tiny model has insufficient memorization capacity so that training accuracy is 0.87 even at the full precision. In this case, the training and test accuracy does not decrease until 6-bit precision. The training accuracy drops sharply as the precision decreases to 5 bits, but the test accuracy degradation is very small. This observation suggests that the quantized networks are more resilient when performing generalization tasks. Thus, the required precision of the network obtained with the memorization task can be considered a conservative estimate.

6. CONCLUDING REMARKS

Quantization of parameters is a straightforward way of reducing the complexity of DNN implementations, especially when VLSI or special-purpose neural processing engines are used. Our study employed simulations on various DNN models. Memorization tests using random binary input data were conducted to determine the capacity by measuring the mutual information. Our simulation results show that the per-parameter capacity is not sensitive to the model size, but is dependent on the structure of the networks, such as FCDNN, CNN, and RNN. The maximum per-parameter memorization capacities of FCDNNs, CNNs, and RNNs are approximately 2.3 bits, 3.0 bits, and 3.7 bits per parameter. Thus, RNNs have the tendency of demanding more bits when compared to FCDNNs. We quantized DNNs under various capacity-utilization regions and showed that the memorization capacity of parameters are preserved up to 6 bits, 8 bits, and 10 bits on FCDNNs, CNN, and RNNs, respectively. The performance of the quantized networks was also tested with the image classification task. The results show that networks need more parameter precision when conducting memorization tasks, rather than inferencing with unseen data. Thus, the precision obtained through the memorization test can be considered a conservative estimate in implementing neural networks for solving real problems. This research not only gives valuable insights on the memorization capacity of neural networks but also provides practical strategies for training and optimizing fixed-point DNNs.

7. REFERENCES

- Gunhan Dundar and Kenneth Rose, "The effects of quantization on multilayer neural networks," *IEEE Transactions on Neural Networks*, vol. 6, no. 6, pp. 1446–1451, 1995.
- [2] Kyuyeon Hwang and Wonyong Sung, "Fixed-point feedforward deep neural network design using weights+ 1, 0, and- 1," in *Signal Processing Systems (SiPS), 2014 IEEE Workshop on*. IEEE, 2014, pp. 1–6.
- [3] Darryl Lin, Sachin Talathi, and Sreekanth Annapureddy, "Fixed point quantization of deep convolutional networks," in *International Conference on Machine Learning*, 2016, pp. 2849–2858.
- [4] Thomas M Cover, "Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition," *IEEE transactions on electronic computers*, , no. 3, pp. 326–334, 1965.
- [5] E Gardner and B Derrida, "Optimal storage properties of neural network models," *Journal of Physics A: Mathematical and general*, vol. 21, no. 1, pp. 271, 1988.
- [6] S Akaho and S Amari, "On the capacity of three-layer networks," in *Neural Networks, 1990.*, 1990 IJCNN International Joint Conference on. IEEE, 1990, pp. 1–6.
- [7] Jasmine Collins, Jascha Sohl-Dickstein, and David Sussillo, "Capacity and trainability in recurrent neural networks," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [8] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, "Understanding deep learning requires rethinking generalization," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [9] Devansh Arpit, Stanislaw Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al., "A closer look at memorization in deep networks," in *International Conference on Machine Learning*, 2017, pp. 233–242.
- [10] Sajid Anwar, Kyuyeon Hwang, and Wonyong Sung, "Fixed point optimization of deep convolutional neural networks for object recognition," in *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on.* IEEE, 2015, pp. 1131–1135.
- [11] Sungho Shin, Kyuyeon Hwang, and Wonyong Sung, "Fixedpoint performance analysis of recurrent neural networks," in Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on. IEEE, 2016, pp. 976–980.
- [12] Fengfu Li, Bo Zhang, and Bin Liu, "Ternary weight networks," arXiv preprint arXiv:1605.04711, 2016.
- [13] Chenzhuo Zhu, Song Han, Huizi Mao, and William J Dally, "Trained ternary quantization," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [14] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David, "Binaryconnect: Training deep neural networks with binary weights during propagations," in Advances in neural information processing systems, 2015, pp. 3123–3131.
- [15] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," in *European Conference on Computer Vision*. Springer, 2016, pp. 525–542.

- [16] Song Han, Huizi Mao, and William J Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [17] Yoonho Boo and Wonyong Sung, "Structured sparse ternary weight coding of deep neural networks for efficient hardware implementations," in *Signal Processing Systems (SiPS), 2017 IEEE International Workshop on.* IEEE, 2017.
- [18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," in Advances in neural information processing systems, 2012, pp. 1097–1105.
- [19] Wonyong Sung, Sungho Shin, and Kyuyeon Hwang, "Resiliency of deep neural networks under quantization," arXiv preprint arXiv:1511.06488, 2015.
- [20] Song Han, Jeff Pool, John Tran, and William Dally, "Learning both weights and connections for efficient neural network," in *Advances in neural information processing systems*, 2015, pp. 1135–1143.
- [21] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and; 0.5 mb model size," arXiv preprint arXiv:1602.07360, 2016.
- [22] Minje Kim and Paris Smaragdis, "Bitwise neural networks," arXiv preprint arXiv:1601.06071, 2016.
- [23] Charbel Sakr, Yongjune Kim, and Naresh Shanbhag, "Analytical guarantees on numerical precision of deep neural networks," in *Proceedings of the 34th International Conference on Machine Learning*, Doina Precup and Yee Whye Teh, Eds., International Convention Centre, Sydney, Australia, 06– 11 Aug 2017, vol. 70 of *Proceedings of Machine Learning Research*, pp. 3007–3016, PMLR.
- [24] Christos Louizos, Karen Ullrich, and Max Welling, "Bayesian compression for deep learning," in Advances in Neural Information Processing Systems, 2017, pp. 3290–3300.
- [25] Asit Mishra, Eriko Nurvitadhi, Jeffrey J Cook, and Debbie Marr, "Wrpn: Wide reduced-precision networks," arXiv preprint arXiv:1709.01134, 2017.
- [26] Thomas M Cover and Joy A Thomas, *Elements of information theory*, John Wiley & Sons, 2012.
- [27] Eric Brochu, Vlad M Cora, and Nando De Freitas, "A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning," arXiv preprint arXiv:1012.2599, 2010.
- [28] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al., "Scikit-learn: Machine learning in python," *Journal of machine learning research*, vol. 12, no. Oct, pp. 2825–2830, 2011.
- [29] Aojun Zhou, Anbang Yao, Yiwen Guo, Lin Xu, and Yurong Chen, "Incremental network quantization: Towards lossless cnns with low-precision weights," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [30] D Bollé, P Dupont, and J Van Mourik, "The optimal storage capacity for a neural network with multi-state neurons," *EPL* (*Europhysics Letters*), vol. 15, no. 8, pp. 893, 1991.