

Deep Model Compression Via Two-Stage Deep Reinforcement Learning

Huixin Zhan¹[0000-0001-8926-1941], Wei-Ming Lin²[0000-0002-9350-6646], and
Yongcan Cao²[0000-0003-3383-0185]

¹ Texas Tech University, Lubbock TX 79415, USA
huixin.zhan@ttu.edu

² The University of Texas at San Antonio, San Antonio TX 78249, USA
{weiming.lin, yongcan.cao}@utsa.edu

Abstract. Besides accuracy, the model size of convolutional neural networks (CNN) models is another important factor considering limited hardware resources in practical applications. For example, employing deep neural networks on mobile systems requires the design of accurate yet fast CNN for low latency in classification and object detection. To fulfill the need, we aim at obtaining CNN models with both high testing accuracy and small size to address resource constraints in many embedded devices. In particular, this paper focuses on proposing a generic reinforcement learning-based model compression approach in a two-stage compression pipeline: pruning and quantization. The first stage of compression, i.e., pruning, is achieved via exploiting deep reinforcement learning (DRL) to co-learn the accuracy and the FLOPs updated after layer-wise channel pruning and element-wise variational pruning via information dropout. The second stage, i.e., quantization, is achieved via a similar DRL approach but focuses on obtaining the optimal bits representation for individual layers. We further conduct experimental results on CIFAR-10 and ImageNet datasets. For the CIFAR-10 dataset, the proposed method can reduce the size of VGGNet by $9\times$ from 20.04MB to 2.2MB with a slight accuracy increase. For the ImageNet dataset, the proposed method can reduce the size of VGG-16 by $33\times$ from 138MB to 4.14MB with no accuracy loss.

Keywords: Compression · Computer vision · Deep reinforcement learning.

1 Introduction

CNN has shown advantages in producing highly accurate classification in various computer vision tasks evidenced by the development of numerous techniques, e.g., VGG [26], ResNet [9], DenseNet [15], and numerous automatic neural architecture search approaches [29, 33]. Albeit promising, the complex structure and large number of weights in these neural networks often lead to explosive computation complexity. Real world tasks often aim at obtaining high accuracy under limited computational resources. This motivates a series of works towards a light-weight architecture design and better speed-up ratio-accuracy trade-off, including Xception [5], MobileNet/MobileNet-V2 [13], ShuffleNet [34], and CondenseNet [14], where group and deep convolutions are crucial.

In addition to the development of the aforementioned efficient CNN models for fast inference, many results have been reported on the compression of large scale models, e.g., reducing the size of large-scale CNN models with little or no impact on their accuracies. Examples of the developed methods include low-rank approximation [7, 20], network quantization [23, 30], knowledge distillation [12], and weight pruning [8, 36, 11, 16, 21], which focus on identifying unimportant channels that can be pruned. However, one key limitation in these methods is the lack of automatic learning of the pruning policies or quantization strategies for reduced models.

Instead of identifying insignificant channels and then conducting compression during training, another potential approach is to use reinforcement learning (RL) based policies to determine the compression policy automatically. There are limited results on RL based model compression [10, 31]. In particular, [10, 31] proposed a deep deterministic policy gradient (DDPG) approach that uses reinforcement learning to efficiently sample the designed space for the improvement of model compression quality. While DDPG can provide good performance in some cases, it often suffers from performance volatility with respect to the hyper-parameter setup and other tuning methods. Besides, these RL-based methods don't directly deal with leveraging the sparse features of CNN, i.e., pruning the small weight connections.

Recently, RL based search strategies have been developed to formulate neural architecture search. For example, [35, 37] considered the generation of a neural architecture via considering agent's action space as the search space in order to model neural architecture search as a RL problem. Different RL approaches were developed to emphasize different representations of the agent's policies along with the optimization methods. In particular, [37] used a recurrent neural network based policy to sequentially sample a string that in turn encodes the neural architecture. Both REINFORCE policy gradient algorithm [28] and Proximal Policy Optimization (PPO) [25] were used to train the network. Differently, [3] used Q-learning to train a policy that sequentially chooses the type of each layer and its corresponding hyper-parameters. Note that [35, 37] focuses on generating CNN models with efficient architectures, while not on the compression of large scale CNN models.

In this paper, we propose to develop a novel two-stage DRL framework for deep model compression. In particular, the proposed framework integrates layer-wise pruning rate learning based on testing accuracy and FLOPs, element-wise variational pruning, and per-layer bits representation learning. In the pruning stage, we first conduct channel pruning that will prune the input channel dimension (i.e., C dimension) with minimized accumulated error in feature maps with the obtained per-layer pruning rate. Then fine-tuning with element-wise pruning via information dropout is conducted to prune the weights in the kernel (i.e., from H and W dimensions).

Briefly, this paper has three main contributions:

1. We propose a novel DRL algorithm that can obtain stabilized policy and address Q-value overestimation in DDPG by introducing four improvements: (1) computational constrained PPO: Instead of collecting T timesteps of action advantages in each of M parallel actors and updating the gradient in each iteration based on MT action advantages in one iteration of the typical PPO, we propose to collect Q-values in each timestep of M parallel actors and update the gradient each timestep based on the M

sampled Q-values; (2) PPO-Clip Objective: We propose to modify the expected return of the policy by clipping subject to policy change penalization. (3) smoothed policy update: Our algorithm first enables multiple agents to collect one minibatch of Q-values based on the prior policy and updates the policy while penalizing policy change. The target networks are then updated by slowly tracking the learned policy network and critic network; and (4) target policy regularization: We propose to smooth Q-functions along regularized actions via adding noise to the target action. The four improvements altogether can substantially improve performance of DRL to yield more stabilized layer-wise prune ratio and bit representations for deep compression, hence outperforming the traditional DDPG. We experimentally show the volatility of DDPG-based compression method in order to backup some common failure mode of policy exploitation in DDPG-based method as shown in Figure 2.

2. Pruning: We propose a new **ppo** with variational pruning compression structure with element-wise variational pruning that can prune three dimensions of CNN. We further learn the *Pareto front* of a set of models with two-dimensional outputs, namely, model size and accuracy, such that at least one output is better than, or at least as good as, all other models by constraining the actions. More compressed models can be obtained with little or no accuracy loss.

3. Quantization: We propose a new quantization method that uses the same DRL-supported compression structure, where the optimal bit allocation strategy (layer-wise bits representation) is obtained in each iteration via learning the updated accuracy. Fine-tuning is further executed after each rollout.

2 A Deep Reinforcement Learning Compression Framework

In this section, we focus on presenting the proposed new generic reinforcement learning based model compression approach in a two-stage compression pipeline: pruning and quantization. Figure 1 shows the overall structure. The two-stage compression pipeline

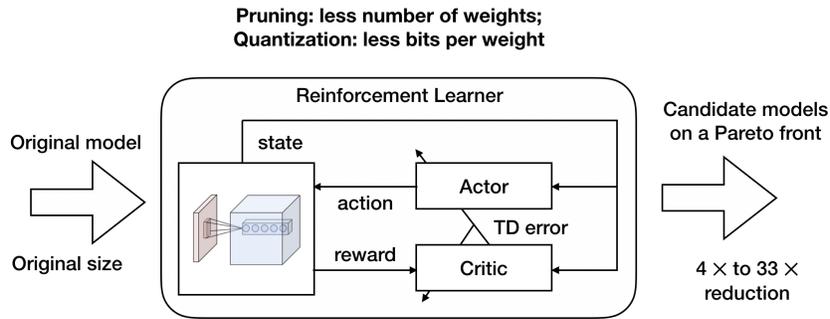


Fig. 1: The proposed deep reinforcement learning compression framework.

includes pruning and quantization. Adopting the pipeline can achieve a typical model

compression rate between $4\times$ and $33\times$. Investigating the Pareto front of candidate compression models shows little or no accuracy loss.

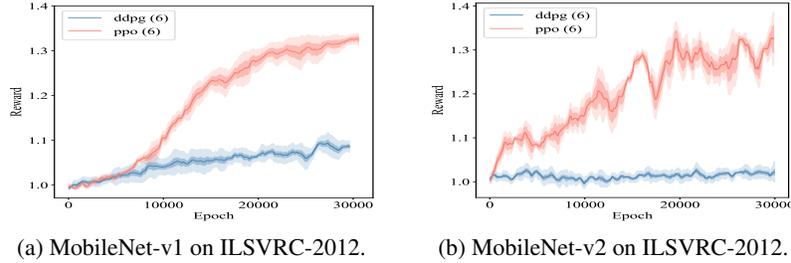


Fig. 2: Comparison of RL-based pruning methods, e.g., **ppo (modified)** and **ddpg**, for MobileNet-v1 and MobileNet-v2 on ILSVRC-2012 for 6 runs.

2.1 State

In both pruning and quantization, in order to discriminate each layer in the neural network, we use a 8-dimension vector space to model a continuous state space:

$$s_t = [N_{L_r}, N, C, H, W, Stride, A_H^t, FLOPs], \quad (1)$$

where N_{L_r} is the index of the layer, N and C are the dimension of, respectively, output channels and input channels, H is the kernel height, W is the kernel width, $Stride$ is the number of pixels shifts over the input matrix, A_H^t is the maximum pruning rate in pruning (respectively, the maximum and minimum bits representation in quantization) with respect to layer t , and $FLOPs$ is the number of floating point operations in each layer.

2.2 Action

In pruning, determining the compression policy is challenging because the pruning rate of each layer in CNN is related in an unknown way to the accuracy of the post-compression model. Since our goal is to simultaneously prune the C , H , and W dimensions. As the dimension of channels increases or the network goes deeper, the computation complexity increases exponentially. Instead of searching over a discrete space, a continuous reinforcement learning control strategy is needed to get a more stabilized scalar continuous action space, which can be represented as $a_t = \{pr_t | pr_t \in [pr_h, pr_l]\}$, where pr_l and pr_h are the lowest and highest and pruning rates, respectively. The compression rate in each layer is taken as a replacement of high-dimensional discrete masks at each weight of the kernels. Similarly, in quantization, the action is also modeled in a scalar continuous action space, which can be represented as $a_t = \{b_t | b_t \in \mathbb{N}^+\}$, where b_t is the number of bits representation in layer t .

2.3 Reward

To evaluate the performance of the proposed two-stage compression pipeline, we propose to construct two reward structures, labeled $r1$ and $r2$. $r1$ is a synthetic reward system as the normalization of current accuracy and FLOPs. $r2$ is an accuracy-concentrated reward system. In pruning, the reward for each layer can be chosen from $r_t \in \{r1, r2\}$. In quantization, we use $r2$ as our selected reward structure. In particular, $r1 = 1 - \frac{FLOP_{s_t} - FLOP_{s_{low}}}{FLOP_{s_{high}} - FLOP_{s_{low}}} + p_{ac}$ and $r2 = p_{ac}$, where p_{ac} is the current accuracy, $FLOP_{s_{high}}$ and $FLOP_{s_{low}}$ are the highest and lowest FLOPs in observation.

2.4 The Proposed DRL Compression Structure

In the proposed model compression method, we learn the Pareto front of a set of models with two-dimensional outputs (model size and accuracy) such that at least one output is better than (or at least as good as) all other outputs. We adopt a popular asynchronous actor critic [22] RL framework to compress a pre-trained network in each layer sequentially. At time step t , we denote the observed state by s_t , which corresponds to the per-layer features. The action set is denoted by \mathcal{A} of size 1. An action, $a_t \in \mathcal{A}$, is drawn from a policy function distribution: $a_t \sim \mu(s_t|\theta^\mu) + \mathcal{N}_t \in \mathbb{R}^1$, referred to as an actor, where θ^μ is the current policy network parameter and the noise $\mathcal{N}_t \in \mathcal{N}(0, \epsilon)$. The actor receives the state s_t , and outputs an action a_t . After this layer is compressed with pruning rate or bits representation a_t , the environment then returns a reward r_t according to the reward function structure $r1$ or $r2$. The updated state s_{t+1} at next time step $t + 1$ is observed by a known state transition function $s_{t+1} = f(s_t; a_t)$, governed by the next layer. In this way, we can observe a random minibatch of transitions consisting of a sequence of tuples $\mathbb{B} = \{(s_t; a_t; r_t; s_{t+1})\}$. In typical PPO, the surrogate objective is represented by $\hat{\mathbb{E}}_t[\frac{\pi^{\theta^\mu}(a|s_t)}{\pi^{\theta^{\mu^-}}(a|s_t)} \hat{A}_t]$, where the expectation $\hat{\mathbb{E}}_t[\cdot]$ is the empirical average over a finite batch of samples and θ^{μ^-} is the prior policy network parameter. If we compute the action advantage \hat{A}_t in each layer, T -step time difference rewards are needed, which is computationally intensive. In resource constrained PPO, we propose to replace the action advantages by Q-functions given by $Q(s_t, a_t) = \mathbb{E}[\sum_{i=t}^{t+T} \gamma^{i-t} r_i | s_t, a_t]$, referred to as critic.

The policy network parameterized by θ^μ and the value function parameterized by θ^Q are then jointly modeled by two neural networks. Let $a = \mu(s_i|\theta^\mu)$, we can learn θ^Q via Q-function regression, namely, Equation (2), and learn θ^μ over the tuples \mathbb{B} with PPO-Clip objective stochastic policy gradient, namely, Equation (3) as

$$\theta^Q = \arg \min_{\theta} \frac{1}{|\mathbb{B}|} \sum_{(s_i, a_i, r_i, s_{i+1}) \in \mathbb{B}} (y_i - Q(s_i, a_i|\theta))^2, \quad (2)$$

$$\theta^\mu = \arg \max_{\theta} \hat{\mathbb{E}}_{(s_i, a_i, \dots) \in \mathbb{B}} \min \left\{ \frac{\pi^\theta(a|s_i)}{\pi^{\theta^{\mu^-}}(a|s_i)} Q(s_i, a|\theta^Q), \right. \\ \left. clip\left(\frac{\pi^\theta(a|s_i)}{\pi^{\theta^{\mu^-}}(a|s_i)}, 1 - c, 1 + c\right) \times Q(s_i, a|\theta^Q) \right\}, \quad (3)$$

where c is the probability ratio of the clipping. A pseudocode of DRL compression structure is shown in Algorithm 1.

Algorithm 1: The proposed DRL compression structure in pruning.

Data: Randomly initialize critic network $Q(s, a|\theta^Q)$ and actor $\mu(s|\theta^\mu)$ with weights θ^Q and θ^μ . Initialize target network Q' and μ' with weights $\theta^{Q'} \leftarrow \theta^Q$, $\theta^{\mu'} \leftarrow \theta^\mu$, the learning rate of the target network ρ , A_H^t , and empty replay buffer \mathcal{D}

Result: Weights θ^Q and θ^μ .

```

1 initialization;
2 while Episode < M do
3   Initialize a random process  $\mathcal{N}$  for action exploration;
4   Receive initial observation state  $s_1$ ;
5    $M \leftarrow M + 1$ ;
6   for  $t = 1, \dots, T$  do
7     Select action  $a_t = clip(\mu(s_t|\theta^\mu) + \mathcal{N}_t, A_H^t)$  according to the current policy
      and exploration noise;
8     Execute  $a_t$ ;
9     Store  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $\mathcal{D}$ ;
10    for  $t = 1, \dots, \mathbb{B}$  do
11      Sample a random minibatch of  $\mathbb{B}$  trajectories from  $\mathcal{D}$ ;
12      Set  $y_t = r_t + \gamma Q(s_{t+1}, \mu(s_{t+1})|\theta^{Q'})$ ;
13      Update the policy by maximizing the ‘‘surrogate’’ objective via stochastic
      gradient ascent with Adam in Equation 3;
14      Pruning the  $C$  dimension of  $t$ -th layer with pruning rate  $a_t$ ;
15      Executing the element-wise variational pruning in Algorithm 2;
16      Update the critic by minimizing the combinatorial loss via stochastic
      gradient descent in Equation 2;
17      Update the target networks via  $\theta^{Q'} \leftarrow \rho\theta^{Q'} + (1 - \rho)\theta^Q$ ,
       $\theta^{\mu'} \leftarrow \rho\theta^{\mu'} + (1 - \rho)\theta^\mu$ ;
18    end
19  end
20 end
```

3 Pruning

In this section, we present two schemes to compress CNN with little or no loss in accuracy by employing reinforcement learning to co-learn the layer-wise pruning rate and the element-wise variational pruning via information dropout. Similar to the aforementioned a3c framework, the layer-wise pruning rate is computed by the actor. After obtaining pruning rate a_t , layer s_t is pruned by a typical channel pruning method [11], whose detail will be given below, to select the most representative channels and reduce the accumulated error of feature maps. In other words, after we get the pruning rate, channel pruning can be used to determine which specific channels are less important or we can simply prune based on the weight magnitude. In each iteration, the CNN layer is further compressed by variational pruning. In particular, we start by learning the connectivity via normal network training. Then, we prune the small-weight connections: all connections with weights that create a representation of the data that is minimal sufficient for the task of reconstruction are remained. Finally, we retrain the network to learn the final weights for the remaining sparse connections.

In pruning, β is a vector whose dimension matches the 4D tensor with shape $N \times C \times W \times H$ in each layer. We also define β^i , the i -th entry of β , as a binary mask for each weight in the kernel. Figure 3 shows the pruning flow in our two schemes. In DRL compression framework, the scalar mask of the j -th weight w_j with mask β^j is set to zero if the weights are pruned based on LASSO regression, discussed in subsection. If pruned, these weights are moved to \bar{s}^j , defined as a set of pruned weights. Otherwise, if the weights are pruned based on information dropout, discussed in subsection, the scalar mask of the j -th weight w_j with mask β^j will be moved to \bar{s}^j with probability $p_{\alpha_j}(\xi^{(j)})$. The weights that play more important role in reducing the classification error are less likely to be pruned.

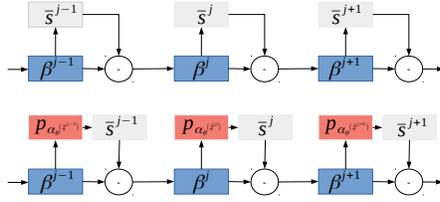


Fig. 3: Comparison of the pruning on two aforementioned schemes, i.e., channel pruning and variational pruning via information dropout.

3.1 Pruning from C Dimension: Channel Pruning

The C -dimension channel pruning can be formulated as:

$$\begin{aligned} \arg \min_{\beta, W} \frac{1}{2N} \left\| Y - \sum_{i=1}^C \beta X_i W_i^T \right\|_F^2 + \lambda \|\beta\|_1 \quad (4) \\ \text{subject to } \|\beta\|_0 \leq pr \times C \\ \|W_i\|_F = 1, \forall i, \end{aligned}$$

where pr is the pruning rate, X_i and Y are the input volume and the output volume in each layer, W_i is the weights, β is the coefficient vector of length C for channel selection, and λ is a positive weight to be selected by users. Then we assign $\beta_i \leftarrow \beta_i \|W_i\|_F$ and $W_i \leftarrow \frac{W_i}{\|W_i\|_F}$.

3.2 Pruning from H and W Dimensions: Variational Pruning

In information dropout, we propose a solution to: (1) efficiently approximate posterior inference of the latent variable \mathbf{z} given an observed value \mathbf{x} based on parameter θ , where \mathbf{z} is a representation of \mathbf{x} and defined as some (possibly nondeterministic) function of

\mathbf{x} that has some desirable properties in some coding tasks \mathbf{y} and (2) efficiently approximate marginal inference of the variable \mathbf{x} to allow for various inference tasks where a prior over \mathbf{x} is required.

Without loss of generality, let us consider Bayesian analysis of some dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$ consisting of N i.i.d samples of some discrete variable \mathbf{x} . We assume that the data are generated by some random process, involving an unobserved continuous random variable \mathbf{z} . Bayesian inference in such a scenario consists of (1) updating some initial belief over parameters \mathbf{z} in the form of a prior distribution $p_{\theta^*}(\mathbf{z})$, and (2) a belief update over these parameters in the form of (an approximation to) the posterior distribution $p_{\theta}(\mathbf{z}|\mathbf{x})$ after observing data \mathbf{x} . In variational inference [17], inference is considered as an optimization problem where we optimize the parameters θ of some parameterized model $p_{\theta}(\mathbf{z})$ such that $p_{\theta}(\mathbf{z})$ is a close approximation of $p_{\theta}(\mathbf{z}|\mathbf{x})$ as measured by the KL divergence $D_{KL}(p_{\theta}(\mathbf{z}|\mathbf{x})|p_{\theta}(\mathbf{z}))$. The divergence between $p_{\theta}(\mathbf{z}|\mathbf{x})$ and the true posterior is minimized by minimizing the negative variational lower bound $\mathcal{L}(\theta)$ of the marginal likelihood of the data, namely,

$$\begin{aligned} \mathcal{L}(\theta, \theta^*; \mathbf{x}^{(i)}) = & - \sum_{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in \mathcal{D}} \mathbb{E}_{p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p(\mathbf{y}^{(i)}|\mathbf{z})] \\ & + \alpha D_{KL}(p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})|p_{\theta^*}(\mathbf{z})). \end{aligned} \quad (5)$$

As shown in [18], the neural network weight parameters θ are less likely to overfit the training data if adding input noise during training. We propose to represent \mathbf{z} by computing a deterministic map of activations $f(\mathbf{x})$, and then multiply the result in an element-wise manner by a random noise ξ , drawn from a parametric distribution p_a with the variance that depends on the input data \mathbf{x} , as

$$\begin{aligned} \mathbf{z} &= (\mathbf{x} \circ \xi)\theta, \\ \xi_{i,j} &\sim p_{a_{\theta}}(\mathbf{x})(\xi_{i,j}), \end{aligned} \quad (6)$$

where \circ denotes the element product operation of two vectors. A choice for the distribution $p_{a_{\theta}}(\mathbf{x})(\xi_{i,j})$ is the log-normal distribution $\log(p_{a_{\theta}}(\mathbf{x})(\xi_{i,j})) = \mathcal{N}(0, a_{\theta}^2(\mathbf{x}))$ [2] that makes the normally fixed dropout rates p_a adaptive to the input data, namely,

$$\begin{aligned} \log(p_{a_{\theta}}(\mathbf{x})(\xi_{i,j})) &\sim \mathcal{N}(\mathbf{z}; 0, a_{\theta}^2(\mathbf{x})I), \\ \log(p_{\theta^*}) &\sim \mathcal{N}(\mathbf{z}; \mu, \sigma^2 I), \end{aligned} \quad (7)$$

where $a_{\theta}(\mathbf{x})$ is an unspecified function of \mathbf{x} . The resulting estimator becomes

$$\begin{aligned} \mathcal{L}(\theta; \mathbf{x}^{(i)}) &\sim \frac{1}{N} \sum_{j=1}^N [-\log p(\mathbf{y}^{(i)}|\mathbf{z}^{(i,j)})] \\ &+ \alpha \left[\frac{1}{2\sigma^2} (a_{\theta}^2(\mathbf{x}^{(i)}) + \mu^2) - \log \frac{a_{\theta}^2(\mathbf{x}^{(i)})}{\sigma} - \frac{1}{2} \right], \end{aligned} \quad (8)$$

where $\mathbf{z}^{(i,j)} \sim (\mathbf{x}^{(i)} \circ \xi^{(i,j)})\theta$ and $\xi^{(i,j)} \sim p_{a_{\theta}}(\xi) = \log \mathcal{N}(0, a_{\theta}^2(\mathbf{x}))$. This loss can be optimized using stochastic gradient descent. A pseudocode of this variational pruning is shown in Algorithm 2 and an illustrative experiment is given in subsection 5.4.

Algorithm 2: Variational pruning.

Data: Pruned model parameters at this iteration θ , the number of fine-tuning iterations \mathcal{Z} , learning rate γ and decay of learning rate τ .

Result: Further compressed and tuned model parameters θ .

```

1 for  $\mathcal{Z}$  iterations do
2   Randomly choose a mini-batch of samples from the training set;
3   Compute gradient of  $\mathcal{L}(\theta; \mathbf{x}^{(i)})$  by  $\frac{\partial \mathcal{L}(\theta; \mathbf{x}^{(i)})}{\partial \theta}$ , where  $\mathcal{L}(\theta; \mathbf{x}^{(i)})$  is computed
   by Equation (8);
4   Update  $\theta$  using  $\theta \leftarrow \theta - \gamma \frac{\partial \mathcal{L}(\theta; \mathbf{x}^{(i)})}{\partial \theta}$ ;
5    $\gamma \leftarrow \tau \gamma$ 
6 end

```

4 Quantization

In the proposed DRL-based quantization-aware training, the RL agent automatically searches for the optimal bit allocation representation strategy for each layer. The modeling of quantization state, action, and rewards are defined in Section 2. The DRL structure is the same as the one for pruning in subsection 2.4. In the fine-tuning step of the quantized CNN, we apply Straight-Through Estimator (STE) [4]. The idea of this estimator of the expected gradient through stochastic neurons is simply to back-propagate through the hard threshold function, e.g., sigmoidal non-linearity function [32]. The gradient is 1 if the argument is positive and 0 otherwise.

5 Experiments

5.1 Settings

In all experiments, the MNIST and CIFAR-10 dataset are both divided by 50000 samples for training, 5000 samples for validation, and 5000 samples for evaluation. The ILSVRC-12 dataset is divided by 1281167 samples for training, 10000 samples for validation, and 50000 samples for evaluation. We adopt a neural network policy with one hidden layer of size 64 and one fully-connected layer using sigmoid as the activation function. We use the proximal policy optimization clipping algorithm with $c = 0.2$ as the optimizer. The critic also has one hidden layer of size 64. The discounting factor is selected as $\gamma = 0.99$. The learning rate of the actor and the critic is set as 1×10^{-3} . In CIFAR-10, the per GPU batch size for training is 128 and the batch size for evaluation is 100. In ILSVRC-12, the per GPU batch size for training is 64 and the batch size for evaluation is 100. The fine tuning steps for each layer are selected as 2000 in the quantization. The parameters are optimized using the SGD with momentum algorithm [27]. For MNIST and CIFAR-10, the initial learning rate is set as 0.1 for LeNet, ResNet, and VGGNet. For ILSVRC-12, the initial learning rate is set as 1 and divided by 10 at rollouts 30, 60, 80, and 90. The decay of learning rate is set to 0.99. All experiments were performed using TensorFlow, allowing for automatic differentiation through the gradient updates [1], on 8 NVIDIA Tesla K80 GPUs.

5.2 MNIST and CIFAR-10

The MNIST [6] and CIFAR-10 dataset [19] consists of images with a 32×32 resolution. Table 1 shows the performance of the proposed method. It can be observed that the proposed method can not only reduce model size but also improve the accuracy (i.e., reduce error rate). In MNIST, comparing with the most recent Drop-Pruning method, our method for LeNet obtains $10\times$ model

compression with a slightly accuracy increase (0.68%). In CIFAR-10, comparing with the baseline model, our method for the VGGNet achieves $9\times$ model compression with a slightly accuracy increase (0.34%). In addition, we compare our algorithm with the commonly adopted weight magnitude channel selection strategy and channel pruning strategy to demonstrate the importance of variational pruning. Please refer to the subsection 5.3 for more details.

Table 1: Results on MNIST and CIFAR-10 dataset.

LeNet	Error (%)	Para.	Pruned Para. (%)	FLOPs (%)
LeNet-5 (DropPruning)	0.73	60K	87.0	-
LeNet-5	0.34	5.94K	90.1	16.4
CIFAR-10	Error (%)	Para.	Pruned Para. (%)	FLOPs (%)
VGGNet(Baseline)	6.54	20.04M	0	100
VGGNet	6.33	2.20M	89.0	48.7
VGGNet	6.20	2.29M	88.6	49.1
ResNet-152 (Baseline)	Error (%)	Para.	Pruned Para. (%)	FLOPs (%)
ResNet-152	5.37	1.70M	0	100
ResNet-152	5.19	1.30M	23.5	71.2
ResNet-152	5.33	1.02M	40.0	55.1

5.3 ImageNet

Table 2: Results on ImageNet dataset for ResNet-18 and ResNet-50 with different speed-ups.

Model	Top-1/Top-5 Error (%)	Pruned Para. (%)	FLOPs (%)		Speed-up \times
			Pruning	Pruning + Quantization	
ResNet-18 (Baseline)	29.36/10.02	0	100		1
ResNet-18	30.29/10.43	30.2	71.4		11.4
ResNet-18	30.65/11.93	51.0	44.2		16.0
ResNet-18	33.40/13.37	76.7	29.5		28.2
ResNet-50 (Baseline)	24.87/6.95	0	100		1
ResNet-50	23.42/6.93	31.2	66.7		12.0
ResNet-50	24.21/7.65	52.1	47.6		16.0
ResNet-50	28.73/8.37	75.3	27.0		29.6

To evaluate the effect of different pruning rates A_H^t , we select 30%, 50%, and 70% for ResNet-18 and ResNet-50 and then evaluate the model pruning on ImageNet ILSVRC-2012 dataset [24]. Experimental results are shown in Table 2 while the per-layer weight bits policy for the quantization is shown in Figure 4. From Table 2, it can be seen that the error increases as the pruning rate increases. However, our pruned ResNet-50 with 30% pruning rate outperforms the pre-trained baseline model in the top-1 accuracy and our pruned ResNet-50 with 30% and 50% pruning rate outperforms the pre-trained baseline model in the top-1 accuracy. In Figure 4, the 8-bit uniform

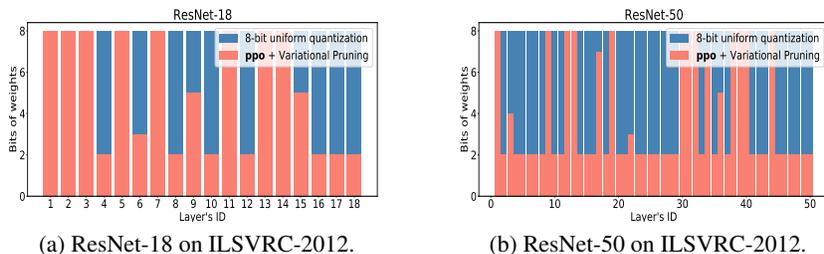


Fig. 4: ResNet-18 and ResNet-50 with different bit allocation strategies.

quantization strategy is shown in blue bar, and the **ppo** with variational pruning policy is shown in red bar. The DRL-supported policy generates a more compressed model with a faster inference speed. By observing the DRL-supported policy, the 3×3 layer is more important than the 1×1 layers because the 1×1 layers are represented by less bits naturally.

To show the importance of our DRL-supported compression structure with variational pruning, we compare RL with channel pruning and RL with variational pruning. Table 3 shows that **ppo** with channel pruning can find the optimal layer-wise pruning rates while **ppo** with variational pruning can further decrease the testing error of the compressed model. Another observation is with the same compression scope, e.g., 50% FLOPs, our model’s accuracy outperforms **ddpg** based algorithms. A comparison of the reward r_1 for AMC [10] (DDPG-based pruning) and our proposed method (PPO-based pruning) is also shown in Figure 2 for $A_H^t = 50\%$ in 6 runs. A typical failure mode of **ddpg** training is the Q-value overestimation, which leads to a lower reward in r_1 . We also report the results for ILSVRC-12 on MobileNet-v2 on Table 3. Although for 70% FLOPs, our model’s performance is competitive to the **ddpg** approach, we achieve lower accuracy decrease on smaller models such as 30% and 60% FLOPs.

Table 3: MobileNet-v1 and MobileNet-v2 on ILSVRC-12.

Model	FLOPs (%)	$\Delta acc(\%)$
MobileNet-v1 (Baseline)	100	0
MobileNet-v1 (ppo + Channel Pruning)	50	-0.2
MobileNet-v1 (ppo + Channel Pruning)	40	-1.1
MobileNet-v1 (ppo + Variational Pruning)	47	+0.1
MobileNet-v1 (ppo + Variational Pruning)	40	-0.8
MobileNet-v1 (ddpg) [10]	50	-0.4
MobileNet-v1 (ddpg) [10]	40	-1.7
0.75 MobileNet-v1 (Uniform) [13]	56	-2.5
0.75 MobileNet-v1 (Uniform) [13]	41	-3.7
Model/Pruning	FLOPs (%)	$\Delta acc(\%)$
MobileNet-v2 (Baseline)	100	0
MobileNet-v2 (ppo + Variational Pruning)	21	-2.4
MobileNet-v2 (ppo + Variational Pruning)	59	-1.0
MobileNet-v2 (ppo + Variational Pruning)	70	-0.8
MobileNet-v2 (ddpg) [10]	30	-3.1
MobileNet-v2 (ddpg) [10]	60	-2.1
MobileNet-v2 (ddpg) [10]	70	-1.0
0.75 MobileNet-v2 (Uniform) [13]	70	-2.0
Model/Quantization	Model Size	$\Delta acc(\%)$
MobileNet-v2 (ppo + Variational Pruning)	0.95M	-2.9
MobileNet-v2 (ppo + Variational Pruning)	0.89M	-3.2
MobileNet-v2 (ppo + Variational Pruning)	0.81M	-3.6
MobileNet-v2 (HAQ) [30]	0.95M	-3.3
MobileNet-v2 (Deep Compression) [8]	0.96M	-11.9

We further examine the results when applying both pruning and quantization on ILSVRC-12. We use the VGG-16 model with 138 million parameters as the reference model. Table 4 shows that VGG-16 can be compressed to 3.0% of its original size (i.e., $33\times$ speed-up) when weights in the convolution layers are represented with 8 bits, and fully-connected layers with 5 bits. Again, the compressed model outperforms the baseline model in both the top-1 and top-5 errors.

Table 4: Comparison with another non-RL two-stage compression method (VGG-16 on ILSVRC-12).

Model (MobileNet-v1)	Layer	Parameters	Pruned Para. (%)	Weight bits		Speed-up \times
				Pruning + Quantization	Pruning + Quantization	
ppo + Variational Pruning	conv1.1 / conv1.2	2K / 37K	42 / 89	8 / 8		2.5 / 10.2
	conv2.1 / conv2.2	74K / 148K	72 / 69	8 / 8		7.0 / 6.8
	conv3.1 / conv3.2 / conv3.3	295K / 590K / 590K	50 / 76 / 58	8 / 8 / 8		4.6 / 10.3 / 5.9
	conv4.1 / conv4.2 / conv4.3	1M / 2M / 2M	68 / 88 / 76	8 / 8 / 8		7.6 / 9.1 / 7.2
	conv5.1 / conv5.2 / conv5.3	1M / 2M / 2M	70 / 76 / 69	8 / 8 / 8		7.1 / 8.5 / 7.1
	fc.6 / fc.7 / fc.8	103M / 17M / 4M	96 / 96 / 77	5 / 5 / 5		62.5 / 66.7 / 14.1
	Total	138M	93.1		5	
Deep Compression [8]	conv1.1 / conv1.2	2K / 37K	42 / 78	5 / 5		2.5 / 10.2
	conv2.1 / conv2.2	74K / 148K	66 / 64	5 / 5		6.9 / 6.8
	conv3.1 / conv3.2 / conv3.3	295K / 590K / 590K	47 / 76 / 58	5 / 5 / 5		4.5 / 10.2 / 5.9
	conv4.1 / conv4.2 / conv4.3	1M / 2M / 2M	68 / 73 / 66	5 / 5 / 5		7.6 / 9.2 / 7.1
	conv5.1 / conv5.2 / conv5.3	1M / 2M / 2M	65 / 71 / 64	5 / 5 / 5		7.0 / 8.6 / 6.8
	fc.6 / fc.7 / fc.8	103M / 17M / 4M	96 / 96 / 77	5 / 5 / 5		62.5 / 66.7 / 14.0
	Total	138M	92.5		5	

5.4 Variational Pruning via Information Dropout

The goal of this illustrative experiment is to validate the approach in subsection 3.2 and show that our regularized loss function $\mathcal{L}(\theta; \mathbf{x}^{(i)})$ shown in Equation (8) can automatically adapt to the data and can better exploit architectures for further compression. The random noise ξ is drawn from a distribution $p_{a_\theta(\mathbf{x})}(\xi)$ with a unit mean $u = 1$ and a variance $a_\theta(\mathbf{x})$ that depends on the input data \mathbf{x} . The variance $a_\theta(\mathbf{x})$ is parameterized by θ . To determine the best allocation of parameter θ to minimize the KL-divergence term $D_{KL}(p_\theta(\mathbf{z}|\mathbf{x})|p_{\theta^*}(\mathbf{z}))$, we still need to have a prior distribution $p_{\theta^*}(\mathbf{z})$. The prior distribution is identical to the expected distribution of the activation function $f(\mathbf{x})$, which represents how much data \mathbf{x} lets flow to the next layer. For a network that is implemented using the softplus activation function, a log-normal distribution is a good fit for the prior distribution (Achille and Soatto 2018). After we fix this prior distribution as $\log(p_{\theta^*}(\mathbf{z})) \sim \mathcal{N}(0, 1)$, the loss can be computed using stochastic gradient descent to back-propagate the gradient through the sampling of \mathbf{z} to obtain the optimized parameter θ . Even if $\log(p_{\theta^*}(\mathbf{z})) \sim \mathcal{N}(0, 1)$, the actual value of u is not necessarily equal to 1 during the runtime. Hence, the mean u and the variance $a_\theta(\mathbf{x})$ of the random noise ξ can be computed via solving the following two equations

$$E(\xi) = e^{u + \frac{a_\theta^2(\mathbf{x})}{2}}, \quad (9)$$

$$D(\xi) = (e^{a_\theta^2(\mathbf{x})} - 1)e^{a_\theta^2(\mathbf{x}) + 2u}, \quad (10)$$

where $E(\xi)$ is the mean of sampled ξ and $D(\xi)$ is the variance of sampled ξ . We add a constraint, $a_\theta(\mathbf{x}) \leq 0.8$, to avoid a large noise variance. Figure 5d shown the probability density function (PDF) of the noise parameter by experiment, which matches a log-normal distribution. The result shows that we optimize the parameters θ of the parameterized model $p_{\theta^*}(\mathbf{z})$ such that $p_{\theta^*}(\mathbf{z})$ is a close approximation of $p_\theta(\mathbf{z}|\mathbf{x})$ as measured by the KL divergence $D_{KL}(p_\theta(\mathbf{z}|\mathbf{x})|p_{\theta^*}(\mathbf{z}))$. After the noise distribution is known, the distribution of $p_\theta(\mathbf{z}|\mathbf{x})$ in Equation (5) can be obtained. In order to show how much information from images that information dropout is transmitting to the second layer, Figure 5b shows the latent variable \mathbf{z} while Figure 5c shows the weights. As shown in Figure 5b, the network lets through the input data (Figure 5a).

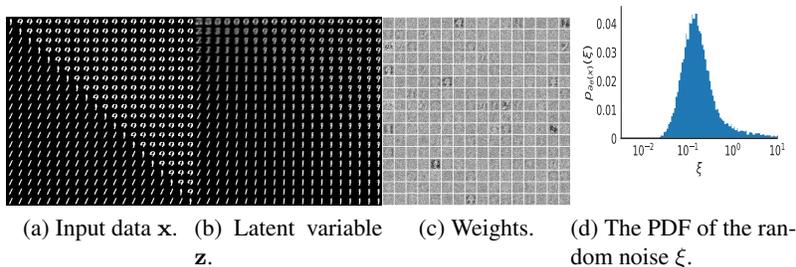


Fig. 5: An illustrative information dropout experiment. Figure 5a shows the input data \mathbf{x} . Figure 5b shows the plot of the latent variable \mathbf{z} at a choice of parameter θ at each spatial location in the third information dropout layer of LeNet trained on MNIST with $\alpha = 1$. The resulting representation \mathbf{z} is robust to nuisances, and provides good performance. Figure 5c shows the weights. Figure 5d shows the PDF of the noise parameter ξ .

5.5 Single Layer Acceleration Performance

In order to further show the importance of variational pruning after obtaining the optimized pruning rate based on the **ppo** algorithm, we test a simple 4-layer convolutional neural network, including 2 convolution (conv) layers and 2 fully connected (fc) layers, for image classification on the CIFAR-10 dataset. We evaluate single layer acceleration performance using the proposed **ppo** with variational pruning algorithm in Section 3 and compare it with the channel pruning strategy. A third typical weight magnitude pruning method is also tested for further comparison, i.e., pruning channels based on the weights’ magnitude (**ppo** + Weight Magnitude Pruning).

Figure 6 shows the performance comparison measured by the error increase after a certain layer is pruned. By analyzing this figure, we can observe that our method (**ppo** + Variational Pruning) earns the best performance in all layers. Since, **ppo** + Channel Pruning applies a LASSO regression based channel selection to minimize the reconstruction error, it achieves a better performance than the weight magnitude method.

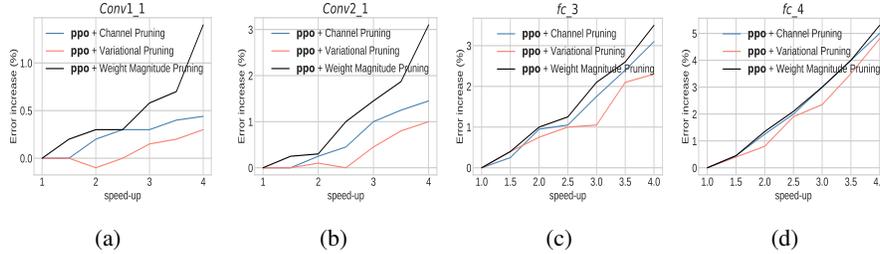


Fig. 6: Single layer error increase under different compression rates. To verify the importance of variational pruning, we considered two baselines: (1) **ppo** + Channel Pruning, and (2) **ppo** + Weight Magnitude Pruning.

Furthermore, the proposed policy considers the fully-connected layers more important than the convolutional layers because the error increase for fully-connected layers is typically larger under the same compression rate.

5.6 Time Complexity

A single convolutional layer with N kernels requires evaluating a total number of NC of the $2D$ kernels $W_n^c * z^c : F = \{W_n^c \in \mathbb{R}^{d \times d} | n = 1, \dots, N; c = 1, \dots, C\}$. Note that there are N kernels $F = \{W_n^c | n = 1, \dots, N\}$ operations on each input channel z^c with cost $O(NCd^2HW)$. The variational pruning via information dropout involves computing a total number of NC' of the $2D$ kernels $W_n^c * z^c$ with cost $O(NC'd^2HW)$, indicating that efficiency inference requires that $C' \ll C$. In subsection 3.2, we consider ameliorating the inference efficiency by information dropout. In the kernels $s^c = \{s_m^c | m = 1, \dots, M\}$, the cost can be reduced to $O(NC'd^2HW)$.

6 Conclusion

Using hand-crafted features to get compressed models requires domain experts to explore a large design space and the trade-off among model size, speed-up, and accuracy, which is often suboptimal and time-consuming. This paper proposed a deep model compression method that uses reinforcement learning to automatically search the action space, improve the model compression quality, and use the FLOPs obtained from fine-tuning with information dropout pruning for the further adjustment of the policy to balance the trade-off among model size, speed-up, and accuracy. Experimental results were conducted on CIFAR-10 and ImageNet to achieve $4\times$ - $33\times$ model compression with limited or no accuracy loss, proving the effectiveness of the proposed method.

Acknowledgment

This work was supported in part by the Army Research Office under Grant W911NF-21-1-0103.

References

1. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al.: Tensorflow: A system for large-scale machine learning. In: 12th USENIX symposium on operating systems design and implementation. pp.265–283 (2016)
2. Achille, A., Soatto, S.: Information dropout: Learning optimal representations through noisy computation. *IEEE transactions on pattern analysis and machine intelligence* 40(12), 2897–2905 (2018)
3. Baker, B., Gupta, O., Naik, N., Raskar, R.: Designing neural network architectures using reinforcement learning. *arXiv preprint arXiv:1611.02167* (2016)
4. Bengio, Y., Léonard, N., & Courville, A. (2013). Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*.
5. Chollet, F.: Xception: Deep learning with depthwise separable convolutions. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 1251–1258 (2017)
6. Cohen, G., Afshar, S., Tapson, J., Van Schaik, A.: Emnist: Extending mnist to handwritten letters. In: *2017 International Joint Conference on Neural Networks*. pp. 2921–2926. (2017)
7. Denton, E.L., Zaremba, W., Bruna, J., LeCun, Y., Fergus, R.: Exploiting linear structure within convolutional networks for efficient evaluation. In: *Advances in neural information processing systems*. pp. 1269–1277 (2014)
8. Han, S., Mao, H., Dally, W.J.: Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149* (2015)
9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778(2016)
10. He, Y., Lin, J., Liu, Z., Wang, H., Li, L.J., Han, S.: Amc: Automl for model compression and acceleration on mobile devices. In: *Proceedings of the european conference on computer vision*. pp. 784–800 (2018)
11. He, Y., Zhang, X., Sun, J.: Channel pruning for accelerating very deep neural networks. In: *Proceedings of the IEEE international conference on computer vision*. pp. 1389–1397(2017)
12. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015)
13. Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861* (2017)
14. Huang, G., Liu, S., Van der Maaten, L., Weinberger, K.Q.: Condensenet: An efficient densenet using learned group convolutions. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 2752–2761 (2018)
15. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4700–4708 (2017)
16. Jia, H., Xiang, X., Fan, D., Huang, M., Sun, C., Meng, Q., He, Y., Chen, C.: Droppruning for model compression. *arXiv preprint arXiv:1812.02035* (2018)
17. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013)
18. Kingma, D.P., Salimans, T., Welling, M.: Variational dropout and the local reparameterization trick. In: *Advances in neural information processing systems*. pp. 2575–2583 (2015)
19. Krizhevsky, A., Nair, V., Hinton, G.: The cifar-10 dataset. online: <http://www.cs.toronto.edu/kriz/cifar.html> 55(2014)

20. Lebedev, V., Ganin, Y., Rakhuba, M., Oseledets, I., Lempitsky, V.: Speeding-up convolutional neural networks using fine-tuned cp-decomposition. arXiv preprint arXiv:1412.6553(2014)
21. Liu, Z., Li, J., Shen, Z., Huang, G., Yan, S., Zhang, C.: Learning efficient convolutional networks through network slimming. In: Proceedings of the IEEE international conference on computer vision. pp. 2736–2744 (2017)
22. Mnih, V., Badia, A.P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., Kavukcuoglu, K.: Asynchronous methods for deep reinforcement learning. In: Proceedings of the international conference on machine learning. pp. 1928–1937 (2016)
23. Rastegari, M., Ordonez, V., Redmon, J., Farhadi, A.: Xnor-net: Imagenet classification using binary convolutional neural networks. In: Proceedings of the european conference on computer vision. pp. 525–542. Springer (2016)
24. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. International journal of computer vision 115(3), 211–252 (2015)
25. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347 (2017)
26. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
27. Sutskever, I., Martens, J., Dahl, G., Hinton, G.: On the importance of initialization and momentum in deep learning. In: Proceedings of the international conference on machine learning. pp. 1139–1147(2013)
28. Sutton, R.S., McAllester, D.A., Singh, S.P., Mansour, Y.: Policy gradient methods for reinforcement learning with function approximation. In: Advances in neural information processing systems. pp. 1057–1063 (2000)
29. Veit, A., Belongie, S.: Convolutional networks with adaptive inference graphs. In: Proceedings of the european conference on computer vision. pp. 3–18 (2018)
30. Wang, K., Liu, Z., Lin, Y., Lin, J., Han, S.: Haq: Hardware-aware automated quantization with mixed precision. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8612–8620 (2019)
31. Wu, J., Zhang, Y., Bai, H., Zhong, H., Hou, J., Liu, W., Huang, W., Huang, J.: Pocketflow: An automated framework for compressing and accelerating deep neural networks (2018)
32. Yin, X., Goudriaan, J., Lantinga, E.A., Vos, J., Spiertz, H.J.: A flexible sigmoid function of determinate growth. Annals of botany 91(3), 361–371 (2003)
33. Yu, X., Yu, Z., Ramalingam, S.: Learning strict identity mappings in deep residual networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4432–4440 (2018)
34. Zhang, X., Zhou, X., Lin, M., Sun, J.: Shufflenet: An extremely efficient convolutional neural network for mobile devices. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 6848–6856 (2018)
35. Zhong, Z., Yan, J., Wu, W., Shao, J., Liu, C.L.: Practical block-wise neural network architecture generation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2423–2432 (2018)
36. Zhuang, Z., Tan, M., Zhuang, B., Liu, J., Guo, Y., Wu, Q., Huang, J., Zhu, J.: Discrimination-aware channel pruning for deep neural networks. In: Advances in neural information processing systems. pp. 875–886 (2018)
37. Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V.: Learning transferable architectures for scalable image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8697–8710 (2018)