# ASCAI: Adaptive Sampling for acquiring Compact AI

**Mojan Javaheripi** and **Mohammad Samragh** and **Tara Javidi** and **Farinaz Koushanfar**

*Department of Electrical and Computer Engineering, UC San Diego, USA*

{*mojan, msamragh, tjavidi, farinaz*}@ucsd.edu

## Abstract

This paper introduces ASCAI, a novel adaptive sampling methodology that can *learn* how to effectively compress Deep Neural Networks (DNNs) for accelerated inference on resource-constrained platforms. Modern DNN compression techniques comprise various hyperparameters that require per-layer customization to ensure high accuracy. Choosing such hyperparameters is cumbersome as the pertinent search space grows exponentially with the number of model layers. To effectively traverse this large space, we devise an intelligent sampling mechanism that adapts the sampling strategy using customized operations inspired by genetic algorithms. As a special case, we consider the space of model compression as a vector space. The adaptively selected samples enable ASCAI to automatically *learn* how to tune per-layer compression hyperparameters to optimize the accuracy/model-size trade-off. Our extensive evaluations show that ASCAI outperforms rule-based and reinforcement learning methods in terms of compression rate and/or accuracy.

## 1. Introduction

With the growing range of applications for Deep Neural Networks (DNNs) on embedded platforms, various DNN compression techniques have been developed to enable execution of such models on resource-limited devices. Some examples of DNN compression methods include pruning [Lin et al. (2018)], quantization [Zhou et al. (2016b); Ghasemzadeh et al. (2018)], nonlinear encoding [Han et al. (2015); Samragh et al. (2019b, 2017)], and tensor decomposition [Kim et al. (2015); Samragh et al. (2019a)].

Compression techniques require tailoring for each DNN architecture, which is generally realized by tuning certain hyperparameters at each layer. Fig. 1 illustrates the importance of intelligent hyperparameter selection for the example of pruning. In general, finding optimal hyperparameters is quite challenging as the space of possibilities grows exponentially with the number of DNN layers. Such large search-space renders manual or computerized greedy hyperparameter tuning algorithms sub-optimal or infeasible. Heuristics that compress one layer at a time overlook the existing inter-dependencies among layers. As such, an intelligent policy that globally tunes the pertinent hyperparameters for all layers is highly desired.
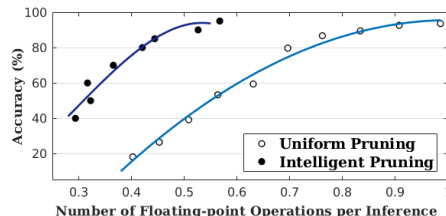


Figure 1: Accuracy/FLOPs Pareto frontiers for pruning a pre-trained VGG network on CIFAR-10. Our intelligent policy achieves a better Pareto curve compared to rule-based (uniform) pruning.

In this paper, we propose ASCAI, an adaptive sampling methodology that automates hyperparameter selection for DNN compression. Genetic algorithms are leveraged to iteratively adapt the sampling strategy. We devise a customized translator that encodes each compressed DNN as a fixed-length sample (vector) of the space, called an *individual*. We then adaptively sample from the pertinent vector space with the goal of finding individuals that render higher compression rate and inference accuracy. Our algorithm initializes a random *population* of individuals and iteratively *evolve*s them towards higher quality generations. To assess the quality of individuals, we develop a scoring mechanism that captures the trade-off between model accuracy and computational complexity. ASCAI evolution aims to encourage survival of individuals with high scores and elimination of the weak ones. Towards this goal, each evolution iteration is modeled by a set of consecutive genetic operations. First, individuals with high scores are selected to generate a new population (*evaluation* and *selection*). Next, the chosen individuals are combined and perturbed to produce children of similar quality (*crossover* and *mutation*). The same procedure continues until convergence.

## 2. Background and Related Work

**Preliminaries and Key Insights.** Recent research showcase the superiority of genetic algorithms to other exploration methods (e.g., reinforcement learning and random search) for large search-spaces [Xie and Yuille (2017)] due to their intriguing characteristics: (1) Genetic algorithms are inexpensive to implement as they do not rely on backpropagation. (2) They are highly parallelizable. (3) Genetic algorithms support a variety of scores and do not suffer in settings with sparse/discrete rewards [Such et al. (2017)]. (4) Finally, these methods can adopt multiple DNN compression tasks.

A genetic algorithm works on a *generation* of *individual*s. The core idea is to encourage creation of superior individuals and elimination of the inferior ones. To this end, an iterative process (Fig. 2) evolves the previous generation into a new, more competent population by performing a set of bio-inspired actions, i.e., *selection*, *crossover*, and *muta-*



Figure 2: Operations of genetic algorithms.

*tion*. During evaluation, individuals are assigned scores representing their quality (*fitness*). The selection step performs a sampling (with replacement) based on individuals' fitness scores. Crossover and mutation create new individuals from existing ones.
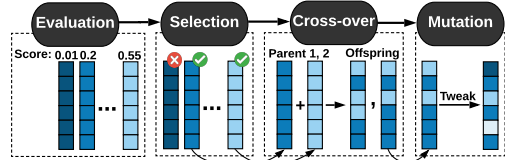
**Related Work.** In the context of deep learning, genetic algorithms have been applied to Neural Architecture Search (NAS) [Xie and Yuille (2017); Real et al. (2017); Huang and Wang (2018)], where the goal is to build a neural network architecture. The objective of NAS is generally achieving higher inference accuracy with no or little emphasis on the execution cost. Different from NAS, this paper focuses on learning hyperparameters for DNN customization, which simultaneously targets execution cost and inference accuracy. Network compression has been studied in contemporary research [He et al. (2018a, 2017); Wang et al. (2017); Jiang et al. (2018); Li et al. (2016); Lin et al. (2017, 2018); Luo et al. (2017)]. These algorithms aim at eliminating redundancies from pre-trained network architectures to reduce computational complexity while preserving inference accuracy. Reinforcement Learning (RL) is proposed as an automated tool that searches for improved model compression quality [He et al. (2018b)]. Although effective in finding near-optimal solutions, RL relies on gradient-based training, which can lead to a high computational burden and a slow convergence. [Hu et al. (2018)] develop a novel pruning scheme that selects the pruned filters using genetic algorithms rather than magnitude-based or gradient-based approaches. Our work is different in that we aim to learn compression hyperparameters rather than the compression technique itself. As a result, ASCAI can be applied to generalized compression techniques.

## 3. ASCAI Approach

An overview of ASCAI methodology is shown in Fig. 3. A translation scheme represents each customized DNN using a vector of per-layer hyperparameters (Sec. 3.1). The initial population is established using a random sampling scheme (Sec. 3.2). At each iteration, the evaluation, selection, crossover, and mutation operations (Sec. 3.3) are performed to update the population towards a new generation. By iteratively applying these operations, ASCAI finds near-optimal compression hyperparameters for a desired DNN. The rest of this section provides details on the search algorithm.
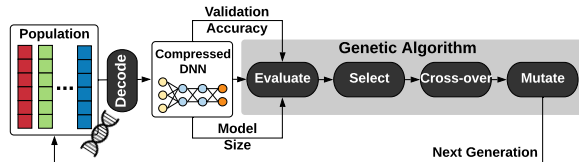


Figure 3: ASCAI solution for DNN customization.

### 3.1 Genetic Translation

Although our proposed approach is applicable to various DNN compression techniques, in this paper we direct our focus on four compression tasks: structured and non-structured Pruning [Li et al. (2016); He et al. (2017)], Singular Value Decomposition [Zhou et al. (2016a)], and Tucker-2 approximation [Kim et al. (2015)]. To construct the individual vectors for our genetic algorithm, we append per-layer hyperparameters as described in the following.

**Pruning.** We allocate a continuous value $p \in [0, 1]$ for each layer to represent the ratio of pruned values. For an $L$-layer DNN, each individual would be a vector $v \in \mathbb{R}^L$ with elements $v_i \in [0, 1]$.

**SVD.** We apply SVD on weight parameters ($W$) of fully-connected layers ($W \in \mathbb{R}^{m \times n}$) and point-wise convolutions ($W \in \mathbb{R}^{m \times n \times 1 \times 1}$). To represent the approximation rank in each layer, we discretize the possibilities into 64 values and encode them as follows:

$$rank \in \{1 \leftarrow \frac{R}{64},\ 2 \leftarrow \frac{2R}{64},\ \ldots,\ 64 \leftarrow R\},\ \ R = min\{m, n\}$$

**Tucker-2.** Tucker decomposition is a generalized form of low-rank approximation for arbitrary-shaped tensors. We apply this method to 4-way weights in convolutional layers, $W \in \mathbb{R}^{m \times n \times k \times k}$. We focus on Tucker-2 which only decomposes the tensor along $m$ and $n$ directions, i.e., output and input channels. For each layer, a tuple of approximation ranks $(rank_m, rank_n)$ should be provided. We quantize the space of decomposition ranks to 8 bins per-way as follows:

$$rank_m \in \{1 \leftarrow \frac{m}{8},\ 2 \leftarrow \frac{2m}{8},\ \ldots,\ 8 \leftarrow m\}, \quad rank_n \in \{1 \leftarrow \frac{n}{8},\ 2 \leftarrow \frac{2n}{8},\ \ldots,\ 8 \leftarrow n\}$$

When applying low-rank approximation, for a DNN that has a total of $L_1$ fully connected and point-wise convolutions ($1 \times 1$ filters) and $L_2$ regular convolution layers ($k \times k$ filters), the individual would be a vector of length $L_1 + 2L_2$ that represent the encoded ranks described above.

### 3.2 Warm Initialization

A naïve initialization of individuals can result in a slow and sub-optimal convergence. To address this, we utilize *warm* population initialization to reduce search time by eliminating unnecessary exploration of low-score regions, i.e., regions at which the inference accuracy is low. Let us denote an individual vector as $v \in \mathbb{R}^L$, the validation dataset as $D = \{(x_m, y_m)\}_{m=1}^M$, and the corresponding classification accuracy as $acc_{D|v}$. During initialization, we only accept randomly sampled (*i.i.d.*) individuals that satisfy an accuracy threshold: $acc_{D|v} > acc_{thr}$. To this end, we find a threshold vector $\theta$ that specifies the maximum per-layer compression when all other layers are uncompressed. Below we explain how to obtain $\theta$ for continuous and discrete hyperparameters.

**Pruning.** We obtain a threshold vector $\theta \in \mathbb{R}^L$ with the $i$-th element $\theta_i$ specifying the maximum pruning rate for the $i$-th layer such that the accuracy threshold $acc_{thr}$ is not violated:

$$\theta_i = max\{p\}\ \ s.t.\ \ v_j = \begin{cases} p & j = i \\ 0 & j \neq i \end{cases} \ \&\ \ acc_{D|v} > acc_{thr} \tag{1}$$

For each individual $v$, the $i$-th element $v_i$ is sampled from a Normal distribution $\mathcal{N}(\theta_i/2, \theta_i/2)$.

**Decomposition.** The threshold vector $\theta$ represents per-layer minimum ranks that satisfy $acc_{thr}$:

$$\theta_i = min\{rank\}\ \ s.t.\ \ v_j = \begin{cases} rank & j = i \\ rank_{max} & j \neq i \end{cases} \ \&\ \ acc_{D|v} > acc_{thr} \tag{2}$$

where $rank_{max}$ corresponds to the non-decomposed layer parameters (see Sec. 3.1). Once $\theta$ is obtained, we uniformly sample $v_i$ from integers $\{\theta_i, \theta_i + 1, \ldots\}$.

### 3.3 Genetic Operations

To enable efficient exploration of the underlying search-space, we devise customized genetic operations, namely, evaluation, selection, crossover, and mutation. Below we describe each step.

**Evaluation.** To assess a population of individuals, we first decode them to their corresponding compressed DNN architectures; this is done by applying compression to each layer of the original DNN based on the corresponding hyperparameter in the individual. We develop a customized scoring mechanism that reflects the acquired model's accuracy and compression rate. Given validation dataset, $D = \{(x_m, y_m)\}_{m=1}^{M}$, the score of individual $v$ evaluated on dataset $D$ is:

$$score(v, D) = \frac{\Delta FLOPS(v)}{PEN_{acc}(D|v)} \tag{3}$$

The numerator encourages reduction in model FLOPs while the denominator penalizes the decrease in model accuracy caused by compression. Here, $\Delta FLOPS(v)$ represents the difference in FLOPS between the uncompressed DNN and the compressed model after applying $v$. $PEN_{acc}(D|v)$ is a measure for accuracy degradation:

$$PEN_{acc}(D|v) = acc_{D|o} - acc_{D|v} + T(acc_{D|v} < acc_{thr}) \times e^{acc_{thr} - acc_{D|v}} \tag{4}$$

where $acc_{D|o}$ is the validation accuracy of the original (uncompressed) model, $acc_{D|v}$ is the accuracy after applying compression with $v$, and $T(.)$ returns the Boolean assessment of the provided inequality. To prevent undesirable drop of accuracy, we greatly diminish the score of individuals that cause lower accuracies than the set constraint, $acc_{thr}$; Having an accuracy constraint is crucial since the genetic algorithm will converge to a model size of zero otherwise. To ensure efficiency, we only use a small portion of the training samples as validation data.

**Selection.** The selection stage in genetic algorithms attempts to choose high-quality individuals to generate the next population. Let us denote the population at the beginning of $t$-th iteration by $\mathbb{P}_t = \{v_n^t\}_{n=1}^{N}$ with fitness scores $\{s_n\}_{n=1}^{N}$ obtained from Eq. 3. Following [Xie and Yuille (2017)], we normalize the scores as $s_n \leftarrow \frac{s_n - s_{min}}{\sum_{n=1}^{N}(s_n - s_{min})}$, where $s_{min}$ is the minimum score in current population. Subtraction of the minimum score ensures that the probability of selecting the weakest individual is zero. The new population, $\mathbb{P}_{t+1} = \{v_n^{t+1}\}_{n=1}^{N}$, is generated by non-uniform random sampling (with replacement) from the old population, $\mathbb{P}_t$. In this non-uniform sampling, the probability of selecting an individual, $v_n^t$, is proportional to its score. This method eliminates weak individuals and passes the high-quality ones to the next generation.

**Crossover.** Given a selected population $\mathbb{P}_{t+1}$, crossover generates two offsprings by operating on each pair of adjacent parent individuals $\{v_{2k-1}^{t+1}, v_{2k}^{t+1}\}_{k=1}^{\frac{N}{2}}$. We use two parameters to control the degree of crossover operation: $p_{cross}$ determines the probability of applying crossover between two individuals, and $p_{swap}$ is the per-element swapping probability. Crossover allows superior individuals to exchange their learned patterns and enables knowledge transfer across the population.

**Mutation.** Mutation randomly tweaks each individual in the population. Similar to crossover, we define two control parameters: $p_{mutate}$ is the probability that the individual gets mutated and $p_{tweak}$ determines the per-element tweaking probability. Mutation allows exploration of the neighborhood of candidate points in the search-space. Each element of a continuous-valued individual is mutated by adding a random value drawn from a zero-mean Normal distribution $\mathcal{N}(0, 0.2)$. Discrete-valued individuals are mutated by randomly incrementing or decrementing vector elements. The values are clipped to the valid ranges after mutation.

# 4. Experiments

We provide extensive evaluations on CIFAR-10 and ImageNet benchmarks. The baseline networks are trained from scratch using PyTorch library. We conduct experiments with non-structured pruning ($P_n$), structured pruning ($P_s$), decomposition ($D$), and combination of multiple compression methods ($D + P_s$). Tab. 1 summarizes the results of ASCAI compressed networks and compares them with prior work that utilize pruning as the compression technique. For brevity, we compare ASCAI with best existing works and exclude other related works.

Table 1: Comparison of ASCAI with state-of-the-art compression methods, namely CP [He et al. (2017)], AMC [He et al. (2018b)], SFP [He et al. (2018a)], FP [Li et al. (2016)], SSS [Wang et al. (2017)], GDP [Lin et al. (2018)], ThiNet [Luo et al. (2017)], and RNP [Lin et al. (2017)]. For ImageNet, we follow common practice in prior work and compare our top-5 accuracy with them.

**CIFAR-10**

| Model | Policy | Acc (%) | Params (%) |
|---|---|---|---|
| ResNet-50 (Non-structured) | Baseline | 93.7 | 100 |
| | AMC | 93.5 | 40 |
| | ASCAI ($P_n$) | 93.6 | 30 |

| Model | Policy | Acc (%) | FLOPs (%) |
|---|---|---|---|
| | Baseline | 93.6 | 100 |
| | CP | 91.8 | 50 |
| | AMC | 91.9 | 50 |
| ResNet-56 | SFP | 93.3 | 47.4 |
| | ASCAI ($P_s$) | 93.2 | 44 |
| | ASCAI ($D$) | 93.5 | 59 |
| | ASCAI ($D + P_s$) | 93.2 | 37 |
| | Baseline | 94.0 | 100 |
| | FP | 93.3 | 61.4 |
| ResNet-110 | SFP | 93.8 | 59.2 |
| | ASCAI ($P_s$) | 93.6 | 41.2 |
| | ASCAI ($D$) | 93.9 | 55 |
| | ASCAI ($D + P_s$) | 92.6 | 21 |

**ImageNet**

| Model | Policy | Acc (%) | FLOPs (%) |
|---|---|---|---|
| | Baseline | 90.0 | 100 |
| | GDP | 87.9 | 24.5 |
| | RNP | 86.3 | 20 |
| VGG-16 | SPP | 87.6 | 20 |
| | AMC | 88.2 | 20 |
| | ASCAI ($P_s$) | 88.1 | 20 |
| | ASCAI ($D$) | 90.1 | 31 |
| | ASCAI ($D + P_s$) | 88.5 | 14 |

**Non-structured Pruning** ($P_n$ **in Tab. 1**). We perform non-structured pruning on ResNet-50 trained on CIFAR-10. We prune the parameters with lowest absolute value as in AMC [Han et al. (2015)], the state-of-the-art that utilizes RL for automated DNN compression. Similar to AMC, we do not perform fine-tuning on the compressed model in this experiment. As shown, ASCAI achieves better accuracy with $1.33\times$ lower parameters.

**Structured Pruning** ($P_s$ **in Tab. 1**). We implement structured pruning by adding masks after $ReLU$ activation layers. Following [Molchanov et al. (2016)], we prune the activations in each layer based on the sum of absolute gradients at the ReLU output. We base our comparisons on the number of operations per inference, i.e., FLOPs, compared to the uncompressed baseline. On CIFAR-10 networks, ASCAI achieves on average $1.25\times$ lower FLOPs, while achieving similar classification accuracy compared to prior art. On VGG-16, ASCAI outperforms all heuristic methods and gives competing results with AMC [He et al. (2018b)].

**Decomposition and Pruning** ($D + P_s$ **in Tab. 1**). To unveil the full potential of our method, we allow ASCAI to learn and combine multiple compression techniques, namely, structures pruning, SVD, and Tucker decomposition. We also report the FLOPs reduction achieved by decomposition separately (shown by $D$ in Tab. 1). Combining multiple techniques allows ASCAI to push the limits of compression. On VGG-16, ASCAI pushes the state-of-the-art FLOPs reduction from $5\times$ to $7.2\times$ with $0.3\%$ higher accuracy.

## 4.1 Analysis and Discussion

To illustrate ASCAI methodology, we consider VGG architecture trained on ImageNet and compressed with structured (filter) pruning. The population is visualized at the initial step (Fig. 4-a) and after 50 iterations of genetic updating (Fig 4-b). Upon convergence, individuals strongly resemble one another and have similarly high scores. ASCAI successfully learns expert-designed rules: first and last rows in Fig. 4-b (first convolution and last fully-connected) are given high densities to maintain inference accuracy. ASCAI performs *whole-network* compression by capturing the state of all layers in each genetic individual. As such, our algorithm can learn which configuration of hyperparameters least affects model accuracy and most reduces the overall FLOPs. To show this capability, we present the per-layer FLOPs for VGG-16 network trained on ImageNet in Fig. 4-c. The bar for each layer shows the percentage of total FLOPs in the original model; the curve shows the percentage of pruned FLOPs in the compressed network. Different from prior art [Jiang et al. (2018)], ASCAI prunes the first convolutions more and relaxes pruning for late convolution and fully-connected layers as they have a minor role in FLOPs.
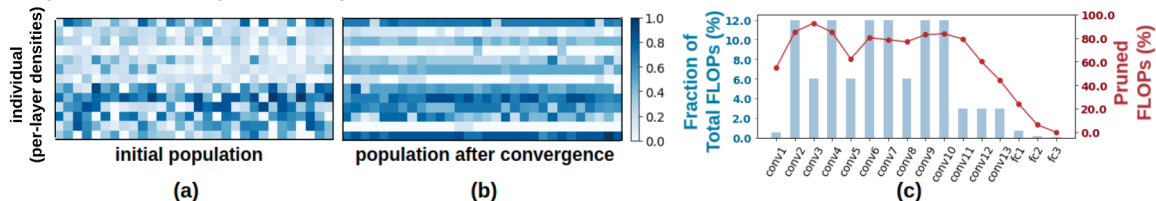


Figure 4: (a) Initialized population for structured pruning. (b) Population upon convergence. Here, each row corresponds to a DNN layer and each column denotes an individual in the population. (c) Per-layer FLOPs (bar charts) and percentage of pruned FLOPs (curve).

## 4.2 Ablation Study

In this section, we study the effect of ASCAI components on algorithm convergence and final FLOPs/accuracy. For brevity, we only focus on structured pruning for CIFAR-10. We show the trend lines as well as a fraction of individuals (black dots) across ASCAI iterations.

**Effect of Initialization.** Fig. 5-a shows the evolution of FLOPs ratio for two initialization policies, one with uniformly random samples and one with our proposed initialization scheme discussed in Sec. 3.2. As seen, naive initialization greatly harms the convergence rate and final FLOPs.

**Effect of Population Size.** Fig. 5-b presents the effect of population size on ASCAI convergence. A higher number of individuals results in a smoother convergence and lower final FLOPs. This effect saturates for a large enough population.



Figure 5: (a) Effect of initialization. (b) Effect of population size.

## 5. Conclusion

This paper introduces ASCAI, a method to automate DNN compression using adaptive sampling with genetic algorithms. Our algorithm learns how the compression hyperparameters should be set across layers to achieve a better performance than models designed by human experts. The core idea behind ASCAI is to translate compression hyperparameters into a vector of genes and explore the corresponding search-space using genetic operations. This approach allows ASCAI to be generic and applicable to any combination of post-processing DNN compression methods.
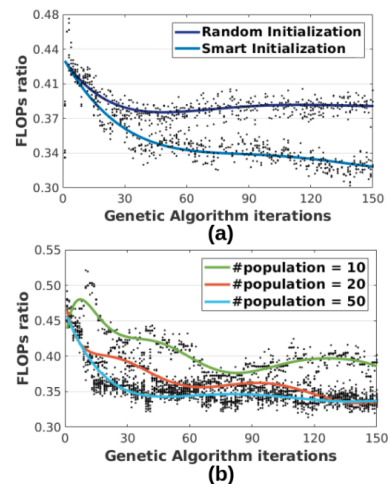
6

# References

Mohammad Ghasemzadeh, Mohammad Samragh, and Farinaz Koushanfar. Rebnet: Residual binarized neural network. In *2018 IEEE 26th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pages 57–64. IEEE, 2018.

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.

Yang He, Guoliang Kang, Xuanyi Dong, Yanwei Fu, and Yi Yang. Soft filter pruning for accelerating deep convolutional neural networks. *arXiv preprint arXiv:1808.06866*, 2018a.

Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *International Conference on Computer Vision (ICCV)*, volume 2, 2017.

Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. Amc: Automl for model compression and acceleration on mobile devices. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 784–800, 2018b.

Yiming Hu, Siyang Sun, Jianquan Li, Xingang Wang, and Qingyi Gu. A novel channel pruning method for deep neural network compression. *arXiv preprint arXiv:1805.11394*, 2018.

Zehao Huang and Naiyan Wang. Data-driven sparse structure selection for deep neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 304–320, 2018.

Chunhui Jiang, Guiying Li, Chao Qian, and Ke Tang. Efficient dnn neuron pruning by minimizing layer-wise nonlinear reconstruction error. In *IJCAI*, pages 2–2, 2018.

Yong-Deok Kim, Eunhyeok Park, Sungjoo Yoo, Taelim Choi, Lu Yang, and Dongjun Shin. Compression of deep convolutional neural networks for fast and low power mobile applications. *arXiv preprint arXiv:1511.06530*, 2015.

Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016.

Ji Lin, Yongming Rao, Jiwen Lu, and Jie Zhou. Runtime neural pruning. In *Advances in Neural Information Processing Systems*, pages 2181–2191, 2017.

Shaohui Lin, Rongrong Ji, Yuchao Li, Yongjian Wu, Feiyue Huang, and Baochang Zhang. Accelerating convolutional networks via global & dynamic filter pruning. In *IJCAI*, pages 2425–2432, 2018.

Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression. In *Proceedings of the IEEE international conference on computer vision*, pages 5058–5066, 2017.

Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient transfer learning. *arXiv preprint arXiv:1611.06440*, 3, 2016.

Esteban Real, Sherry Moore, Andrew Selle, Saurabh Saxena, Yutaka Leon Suematsu, Jie Tan, Quoc V Le, and Alexey Kurakin. Large-scale evolution of image classifiers. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2902–2911. JMLR. org, 2017.

Mohammad Samragh, Mohammad Ghasemzadeh, and Farinaz Koushanfar. Customizing neural networks for efficient fpga implementation. In *Field-Programmable Custom Computing Machines (FCCM), 2017 IEEE 25th Annual International Symposium on*, pages 85–92. IEEE, 2017.

Mohammad Samragh, Mojan Javaheripi, and Farinaz Koushanfar. Autorank: Automated rank selection for effective neural network customization. *ML-for-Systems workshop at the 46th International Symposium on Computer Architecture (ISCA)*, 2019a.

Mohammad Samragh, Mojan Javaheripi, and Farinaz Koushanfar. Codex: Bit-flexible encoding for streaming-based fpga acceleration of dnns. *arXiv preprint arXiv:1901.05582*, 2019b.

Felipe Petroski Such, Vashisht Madhavan, Edoardo Conti, Joel Lehman, Kenneth O Stanley, and Jeff Clune. Deep neuroevolution: Genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning. *arXiv preprint arXiv:1712.06567*, 2017.

Huan Wang, Qiming Zhang, Yuehai Wang, and Haoji Hu. Structured probabilistic pruning for convolutional neural network acceleration. *arXiv preprint arXiv:1709.06994*, 2017.

Lingxi Xie and Alan Yuille. Genetic cnn. *arXiv preprint arXiv:1703.01513*, 2017.

Hao Zhou, Jose M Alvarez, and Fatih Porikli. Less is more: Towards compact cnns. In *European Conference on Computer Vision*, pages 662–677. Springer, 2016a.

Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *arXiv preprint arXiv:1606.06160*, 2016b.