

Evolutionary Neural Architecture Search for Image Classification

Mid-term Report

Bowen Zheng, Shijie Chen, Shuxin Wang
Department of Computer Science and Engineering
Southern University of Science and Technology
Shenzhen, Guangdong, China

Abstract—In this project, we propose an elite-parent preserving evolutionary framework for Neural Architecture Search on image classification problems. We have finished the evolutionary algorithm framework and the construction of neural architecture. Currently we focus on a cell-based search space where a neural architecture is composed of multiple interconnected cells. Experiments shows the effectiveness of our algorithm in obtaining a neural architecture with good performance. In the future, we will try to add more evolutionary operations and try to enlarge the search space.

I. INTRODUCTION

The great leap of computing resources in the past few decades made it possible to fully utilize the potential of neural networks. In recent years neural networks outperformed traditional methods in many fields of research, especially image classification. However, the state-of-the-art architectures are carefully designed and tuned by researchers for a specific problem. Therefore, people start to think about automate the design of neural networks in the hope of finding the best-performing network architecture efficiently.

Neural architecture search is a research field focusing on automating the design of neural networks. Currently there are a few popular approaches, including reinforcement learning, bayesian optimization, tree-based searching and genetic-based evolutionary algorithms.

This project focuses on NAS for image classification problems. The reason is that this area is well explored and there exists many high-performance hand-crafted neural architectures. They provide a good guidance and target to our project. In addition, neural networks for image classification are mostly built upon basic cells including convolution, polling, normalization, and activation layers. This helps to shrink our search space.

II. RELATED WORKS

A lot of research works have been done in each of the three categories in NAS. Some proposed algorithms can design architectures that is on par of or even more capable than state-of-the-art hand-crafted networks.

A. Search Space

The search space of NAS determines the possible architectures a NAS algorithm can find.

The simplest search space is the simple multiple-layer structure, in which a neural network A is composed of multiple layers L_i connected to the neighboring layers. In this case, the search space can be described by (1) The maximum number of layers (2) The type and dimension of each layer and their hyper-parameters [1] [2].

In more recent studies, some researchers use a cell based search space in which the possible architecture of cells are explored. A cell is nothing but a smaller neural network. The entire neural network is constructed by connecting several pre-defined cells. A cell has less layers but allows more complex architectures like skip connections between any layers [3] [4]. The cell could be some hand-crafted neural networks that have already been proofed effective. The search space is therefore decreased to the possible arrangements of cells.

In contrast to the above direction, some researchers also tried to search for effective cells and connect them at last in a predefined manner [5] [3]. The search space is greatly decreased in that each cell is comparably small. This method can also be easily transferred to other datasets [5] since the structure of cells are not fixed.

Recently, some researchers managed to optimize the overall architecture as well as the cells at the same time and obtained state-of-the-art result [6].

B. Search Strategy

Many different search strategies can be applied to explore the search space discussed above. These methods includes bayesian optimization, evolutionary methods, reinforcement learning and gradient-based methods.

Evolutionary algorithms has been used to evolve neural networks since 1989 [7]. Earlier works use genetic algorithms to both optimize the structure of neural networks and train the networks [8]. However, with the birth of back-propagation (BP), recent works of neural-evolution use genetic algorithms only for optimizing neural architectures and use BP to train the networks [9]. In the context of NAS, the individuals in the genetic algorithm are neural network architectures and the genetic operations (crossover and mutation) are used to alter the architecture by adding/removing layers or change connectivity of nodes.

Genetic algorithms shows their diversity in how they sample parents, generate offspring and update population. Some work choose parents from a preto-optimal front [10] while others use tournament selection [11] [4] [9]. When generating offsprings, some algorithms randomly initialize weight of child networks. In comparison, Lamarckian inheritance is used in [10] so that child networks could inherit weight of its parent and the training cost is reduced. To update population, some algorithms abandon least capable individuals [9] while some delete the oldest individuals [4]. A more sophisticated policy is developed by [12] and [13] in which the age of the individuals are taken into account.

There are other methods that are used to implement NAS, including bayesian optimization, reinforcement learning ,tree-based search and gradient-based methods. We don't discuss them here since we use evolutionary algorithms in our project.

C. Performance Estimation Strategy

One important issue in neural architecture search is the estimation of neural network performance. This is critical in the population update policy of evolutionary algorithms.

The simplest way to estimate performance is to train every searched network from scratch and test the desired performance metric, e.g. accuracy on validation set. However, the training of neural network is very time and computation consuming.

An alternative is to estimate performance on lower fidelities. More specifically, to train the network for shorter period of time [5] or on some subset of the dataset [14]. However, the estimate must ensure the result ranking of different networks must be the same as that of complete training. That is to say, there exist a trade-off between computational load and estimation fidelity.

Another approach to estimate performance is based on learning curve extrapolation [15]. This method accelerate estimation by stop poor performance networks at the early state of training based on statistical patterns of learning curves. Other researchers propose ways to predict neural network performance based on architectural and cell properties [11].

III. METHODOLOGY

We will develop an evolutionary neural architecture search algorithm based on the *age* of individuals. By incorporating *age* as a part of gene, we can prolong the existence of good individuals [13] and kill average individuals when their age reach a certain limit [12].

Combining the advantage of the above works, we propose a population update police based on *age* and that preserves a parent whose child has good performance. An individual is dropped from the total population if its *age* exceeds a predefined *lifetime*. However, we prolong its life if its offspring performs well (e.g. within top $P\%$ in ranking). In this way, we hope to preserve good parent architectures in the population.

To test the effectiveness of our algorithm, we will experiment on image classification datasets including CIFAR-10, CIFAR-100 and will possibly extend to IMAGENET.

IV. SYSTEM DESIGN

A. Cell-based Search Space

We explore a cell-based search space with 6 hidden cells, as is illustrated in Fig.1. Each edge in Fig.1 represents a neural network unit. Each square represents a tensor in the neural network.

1) *Neural Network Units*: There are 7 possible neural network units that is suitable for image classification problems in our search space:

- 1) identity
- 2) 3×3 average polling
- 3) 3×3 max polling
- 4) 1×1 convolution
- 5) 3×3 depthwise-separable convolution
- 6) 3×3 dilated convolution
- 7) 3×3 convolution

The search space is the possible connections of the states using different neural network units.

2) *Cell Architecture Representation*: We represent the architecture of a cell by a vector R which is a vector of lists of tuple $R_i = (a_i, b_i)$ where a_i marks the input state of state i and b_i shows the unit used in the state transition connection.

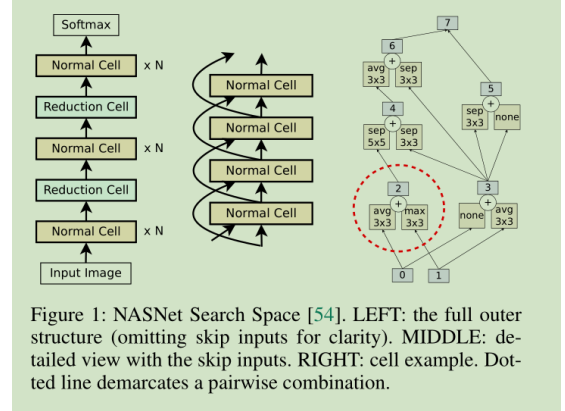


Figure 1: NASNet Search Space [54]. LEFT: the full outer structure (omitting skip inputs for clarity). MIDDLE: detailed view with the skip inputs. RIGHT: cell example. Dotted line demarcates a pairwise combination.

Fig. 1: The artitecture of a cell

B. Overall Neural Network Architecture

As is illustrated in Fig.2 overall network consists of input layer, 6 identical cells, 3 polling layers, 2 fully connected perceptron layers, a softmax layer and an output layer.

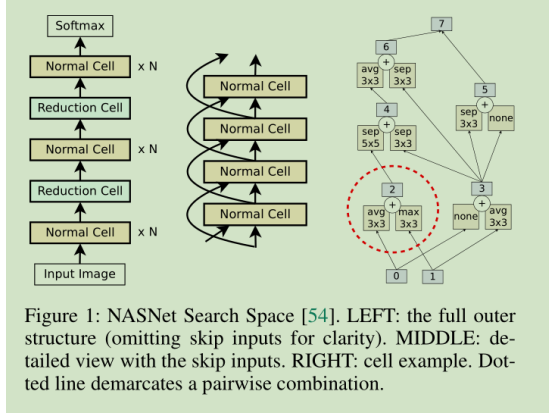


Fig. 2: The overall artitecture of the total network. LEFT: outer structure. RIGHT: detailed skip connection of cells.

For simplicity, all cells are the same.

C. Evolutionary Algorithm Framework

For now, we follow the design of a cell based evolutionary NAS algorithm described in [13] except that we propose a population update policy to preserve good parents. We will further examine initalization, crossover and mutation strategies to optimize our algorithm. We will also compare different performance estimation approaches to cut computation demand.

Algorithm 1 Elite Parent Preserving Evolution:

```

1:  $population \leftarrow \phi$ 
2:  $entireGen \leftarrow \phi$ 
3:  $generation \leftarrow 0$ 
4: while  $population\ size < N$  do
5:    $newNetwork \leftarrow networkInit()$ 
6:    $trainNetwork(newNetwork)$ 
7:   add  $newNetwork$  to  $popuption$  and  $entireGen$ 
8: end while
9: while  $genetation < G$  do
10:   $parent \leftarrow$  select a parent from population using
    tournament selection
11:   $child \leftarrow mutate(parent)$ 
12:   $trainNetwork(child)$ 
13:  add  $child$  to  $popuption$  and  $entireGen$ 
14:  if accuracy of  $child$  is better than  $P\%$  individuals in
    population then
15:    extend the  $lifetime$  of  $parent$  by  $t$ 
16:  end if
17:  for all  $individual$  in  $population$  do
18:    update the  $age$  of  $individual$ 
19:    remove current  $individual$  if its age reaches its
     $lifetime$ 
20:  end for
21: end while
22: return the network model with highest accuracy in
     $entireGen$ 

```

The *population* size is N and the algorithm evolve G generations in total. *entireGen* stores all network models that we generated.

In *networkInit()*, the initial network model is generated. Then its *age* is set to 1 and *lifetime* is set to the default value.

P and t are hyper-parameters controlling the actual lifespan of an individual.

V. -PROPOSED FRAMEWORK

Variable	Name
N_{total}	size of the total population

TABLE I: Notation

A. Fuzzy Rules

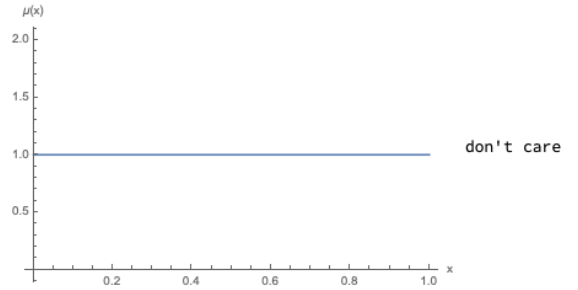


Fig. 3: The *don'tcare* function

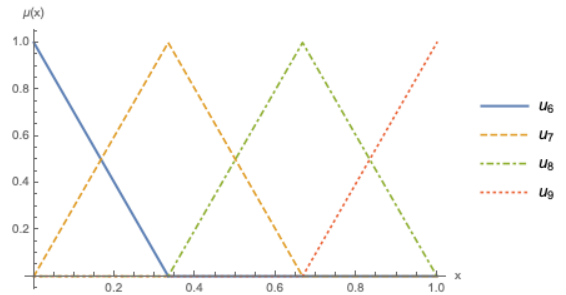


Fig. 4: Membership Functions with 4 intervals

The antecedent part of R_q is a set of antecedent fuzzy sets:

$$A_q = \{A_{qi} | i = 1, 2, \dots, n\} \quad (1)$$

VI. ASYNCHRONOUS PARALLEL DISTRIBUTED SYSTEM DESIGN

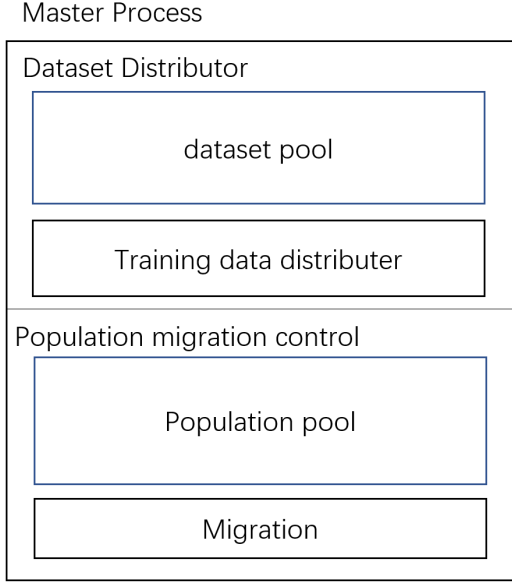


Fig. 5: The master process

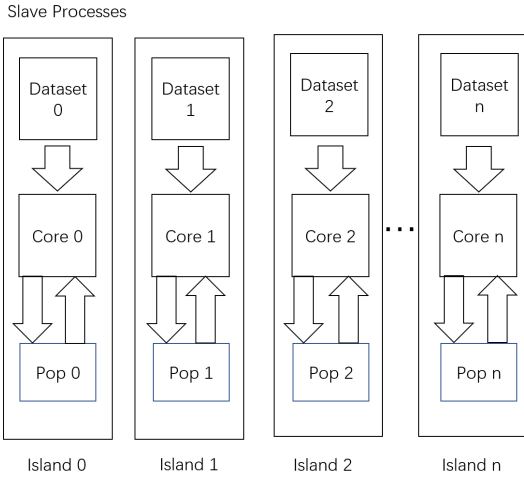


Fig. 6: The slave process

A. Master and Slave process

- 1) Master Process:
- 2) Slave Processes:

B. Dataset Distributor:

C. Population Migration Control:

D. Asynchronous Population Migration

VII. EXPERIMENTS

A. Pareto Front

- 1) With different cores: As is shown by Fig.7 and Fig.8 ($I_{update} = 100, N_{total} = 264$), our model is able to get

results that is similar to non-parallel models. Note that with the increase of the number of cores, the ability of the model to obtain classifiers with more rules and antecedent fuzzy sets is decreasing, due to the decrease of sub-population. It's acceptable since our objective is to minimize the two objectives.

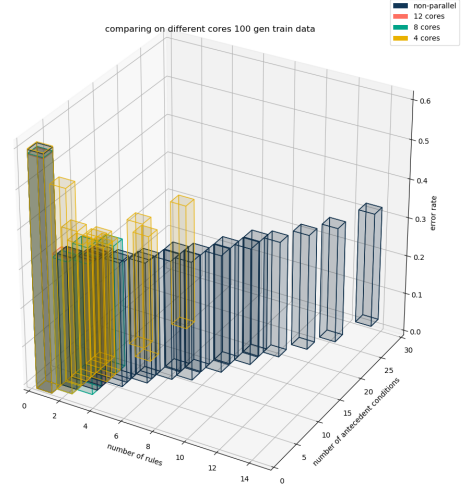


Fig. 7: Pareto Front on training data

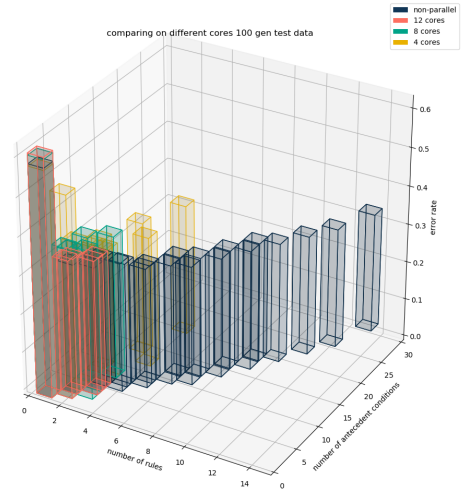


Fig. 8: Pareto Front on test data

- 2) With different exchange interval: We can see from Fig.9, Fig.10 and Fig.11 ($N_{total} = 264, N_{island} = 8$) that the choice of exchange interval also impacts the performance of our model. A longer exchange interval leads to less communication between master and slave processes, thus reduces running time.

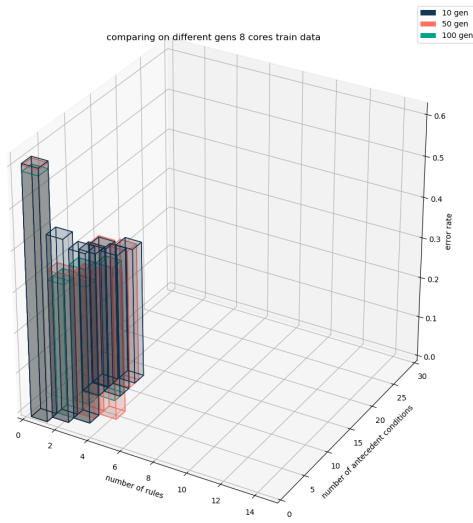


Fig. 9: Pareto Front on training data with different I_{update}

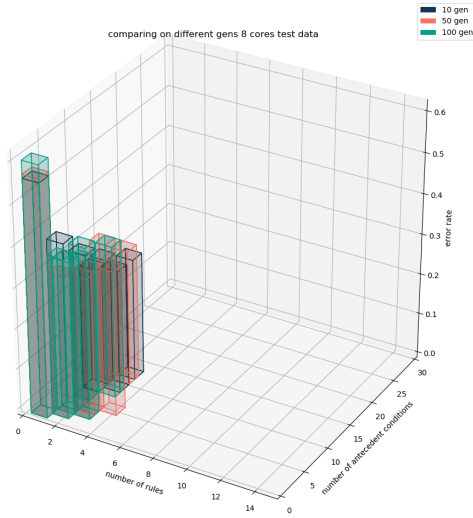


Fig. 10: Pareto Front on training data with different I_{update}

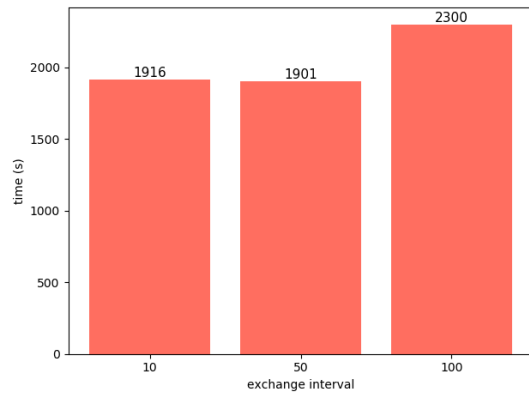


Fig. 11: Computational time with different I_{update}

B. Computation Time

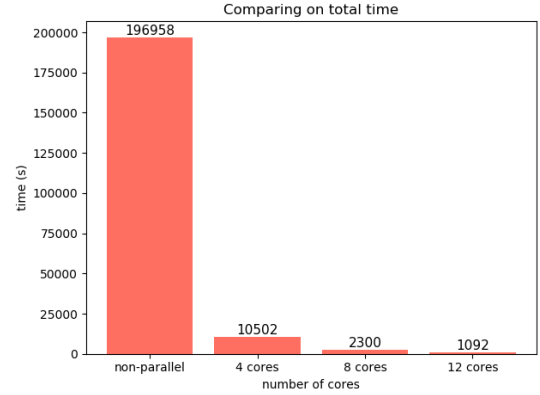


Fig. 12: Computation time of 3000 generations, $N_{pop} = 264$.

With n cores, we observe a speed up of up to n^2 times in our computational experiments.

C. Compare with other model

We compared our model with the synchronized model by Nojima et al. [?], as is shown in Fig.13, Fig.14 and Fig.15. The experiments are done using 8 cores.

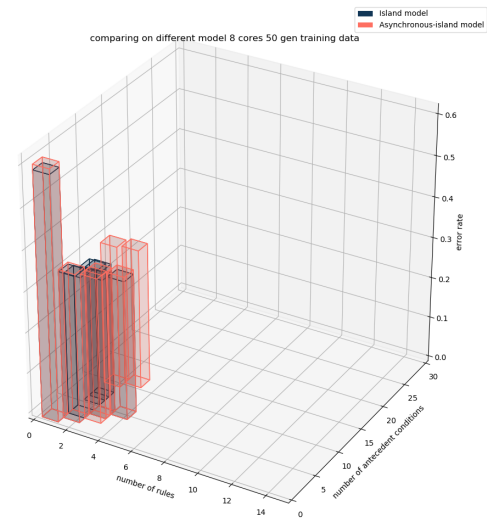


Fig. 13: Comparison on training data.

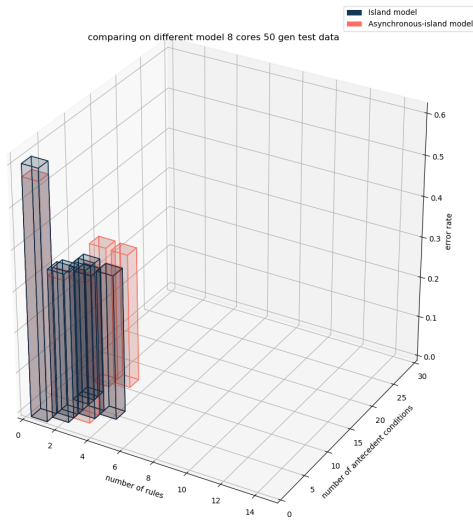


Fig. 14: Comparison on test data.

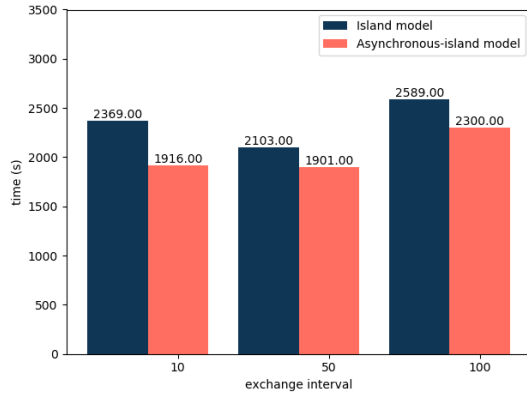


Fig. 15: Comparison on total time.

VIII. CONCLUSION

IX. CONTRIBUTION

REFERENCES

- [1] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1251–1258, 2017.
- [2] B. Baker, O. Gupta, N. Naik, and R. Raskar, "Designing neural network architectures using reinforcement learning," *arXiv preprint arXiv:1611.02167*, 2016.
- [3] H. Cai, J. Yang, W. Zhang, S. Han, and Y. Yu, "Path-level network transformation for efficient architecture search," *arXiv preprint arXiv:1806.02639*, 2018.
- [4] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, "Regularized evolution for image classifier architecture search," *arXiv preprint arXiv:1802.01548*, 2018.
- [5] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8697–8710, 2018.
- [6] C. Liu, L.-C. Chen, F. Schroff, H. Adam, W. Hua, A. Yuille, and L. Fei-Fei, "Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation," *arXiv preprint arXiv:1901.02985*, 2019.
- [7] G. F. Miller, P. M. Todd, and S. U. Hegde, "Designing neural networks using genetic algorithms," in *ICGA*, vol. 89, pp. 379–384, 1989.

- [8] K. O. Stanley and R. Miikkulainen, "Evolving neural networks through augmenting topologies," *Evolutionary computation*, vol. 10, no. 2, pp. 99–127, 2002.
- [9] E. Real, S. Moore, A. Selle, S. Saxena, Y. L. Suematsu, J. Tan, Q. V. Le, and A. Kurakin, "Large-scale evolution of image classifiers," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 2902–2911, JMLR. org, 2017.
- [10] T. Elsken, J. H. Metzen, and F. Hutter, "Efficient multi-objective neural architecture search via lamarckian evolution," 2018.
- [11] C. Liu, B. Zoph, M. Neumann, J. Shlens, W. Hua, L.-J. Li, L. Fei-Fei, A. Yuille, J. Huang, and K. Murphy, "Progressive neural architecture search," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 19–34, 2018.
- [12] G. S. Hornby, "Alps: The age-layered population structure for reducing the problem of premature convergence," in *Proceedings of the 8th Annual Conference on Genetic and Evolutionary Computation, GECCO '06*, (New York, NY, USA), pp. 815–822, ACM, 2006.
- [13] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, "Regularized evolution for image classifier architecture search," *CoRR*, vol. abs/1802.01548, 2018.
- [14] A. Klein, S. Falkner, S. Bartels, P. Hennig, and F. Hutter, "Fast bayesian optimization of machine learning hyperparameters on large datasets," *arXiv preprint arXiv:1605.07079*, 2016.
- [15] T. Domhan, J. T. Springenberg, and F. Hutter, "Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves," in *Twenty-Fourth International Joint Conference on Artificial Intelligence*, 2015.