# Automate Deep Learning in Image Processing

1  **Xinyu Ma**
2  Institute of Computing Technology
3  Chinese Academy of Science
4  Beijing, P.R.C
5  *xinyuma2016@gmail.com*

6  ## Abstract

7  Convolutional neural work(CNN) have achieved great success in image
8  processing. This kind of neural networks can recognize visual patterns
9  directly from pixel images with minimal preprocessing. But it is hard to
10  design the architecture and also restricted by the availability of the
11  computing resources. For example, GoogleNet has 22 hidden layers and
12  trained on 2 GPUS for a week; VGGNet consists of 16 convolutional layers
13  and trained on 4 GPUs for 2-3 weeks. I did a research about automated
14  exploring and designing a better architecture of neural networks and
15  combining some work I have done this semester about automated exploring
16  CNN architecture with reinforcement learning. Distributed neural network
17  training platform is also necessary for accelerating model training.
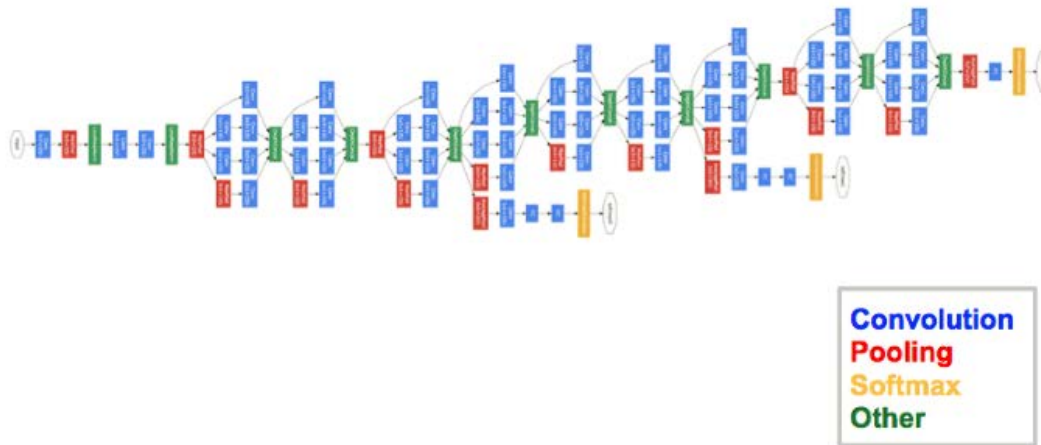
18

19  ## 1     Introduction

20  With the development of computing resource, like GPU and AI chip, machine learning
21  experts can design deeper and more complicated neural networks. The common problem in
22  deep learning area is that it's still very difficult to design neural network architecture and
23  model's hyper parameter tuning is too heady and too tedious. So I am very interested in how
24  to automatically design a good architecture of neural networks especially of convolutional
25  neural networks in a dataset without hyper parameter tuning and other works. It means that
26  you give a dataset, platform will give you a good neural network model.

27  Firstly, I will introduce some classic CNN architecture like GoogleNet, ResNet. Or maybe I
28  should introduce CNN first. A typical CNN architecture consists of several convolution,
29  pooling, and fully connected layers. While constructing a CNN, a network designer has to
30  make numerous design choices: the number of layers of each type, the ordering of layers, and
31  the hyper parameters for each type of layer, e.g., the receptive field size, stride, and number
32  of receptive fields for a convolution layer. The number of possible choices makes the design
33  space of CNN architectures extremely large and hence, infeasible for an exhaustive manual
34  search. While there has been some work (Pinto et al., 2009; Bergstra et al., 2013; Domhan et
35  al., 2015;Bowen et al., 2017;Hieu et al. 2018) on automated or computer-aided neural
36  network design, new CNN architectures or network design elements are still primarily
37  developed by researchers using new theoretical insights or intuition gained from
38  experimentation.

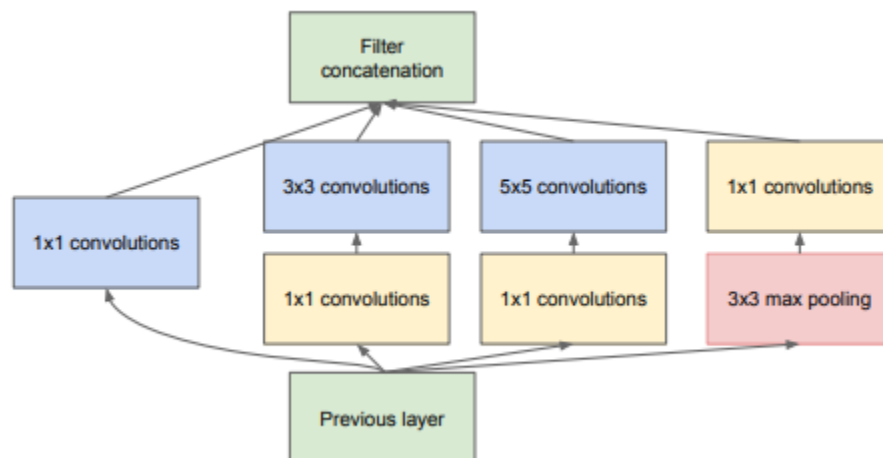39

40  ### 1.1     GoogleNet
41  GoogleNet won the ImageNet Large Scale Visual Recognition Challenge (ILSVRC 2014)
42  competition. It achieved a top-5 error rate of 6.67%. This was very close to human level
43  performance. The network used a CNN inspired by LeNet but implemented a novel element
44  which is dubbed an inception module. It used batch normalization, image distortions and

45  RMSprop. This module is based on several very small convolutions in order to drastically
46  reduce the number of parameters. Its architecture consisted of a 22 layer deep CNN but
47  reduced the number of parameters from 60 million (AlexNet) to 4 million.



**Convolution**
**Pooling**
**Softmax**
**Other**

48
49                                Figure 1. GoogleNet

50  We can see from figure 1 that GoogleNet's architecture is not only the simple combination of
51  convolution, pooling and fully connected layers' order, it has small modules called Inception
52  which is more complicated than all CNNs' architecture that we have ever seen.



53
54                                Figure 2.Inception

55  Inception module includes 1x1, 3x3 and 5x5 convolutions, 3x3 max pooling. They are
56  combined in a fixed order that embodies senior deep learning expert's wisdom.

57
58  **1.2     ResNet**

59  Residual networks(ResNet) won the ILSVRC 2015 classification task by achieving 3.57%
60  error on the ImageNet test set with an ensemble of these residual networks. This was almost
61  above human level performance. ResNet are easier to optimize and can gain accuracy from
62  considerably increased depth. It has 152 layers and includes 3x3 and 5x5 convolutions with
63  different filter size.

64  It's impossible to train such deeper neural network with current computing resource. So
65  ResNet proposed let layers fit a residual mapping not the original to avoid gradient exploding
66  and vanishing. Also ResNet proposed shortcut to skip some layers and connect layers. The
67  identity shortcuts can be directly used when the input and output are of the same dimensions.
68  When the dimensions increase, ResNet consider two options: (A) The shortcut still performs

69  identity mapping, with extra zero entries padded for increasing dimensions. This option
70  introduces no extra parameter; (B) The projection shortcut is used to match dimensions (done
71  by 1x1 convolutions). For both options, when the shortcuts go across feature maps of two
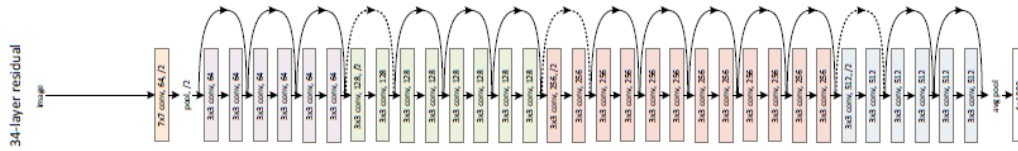72  sizes, they are performed with a stride of 2.

73



74  Figure 3. ResNet

75  As we can see, it's too complicate to design for general deep learning experts. So, here
76  comes the question: Can we design CNN automatly?

77

78  ## 2    Related Work

79  Sometimes designing neural network architectures and reinforcement learning can do
80  independently. Research on automating neural network design goes back to the 1980s when
81  genetic algorithm-based approaches were proposed to find both architectures and weights
82  (Schaffer et al., 1992). However, to the best of our knowledge, networks designed with
83  genetic algorithms, such as those generated with the NEAT algorithm (Stanley &
84  Miikkulainen, 2002), have been unable to match the performance of hand-crafted networks
85  on standard benchmarks (Verbancsics & Harguess, 2013). Other biologically inspired ideas
86  have also been explored; motivated by screening methods in genetics, Pinto et al. (2009)
87  proposed a high-throughput network selection approach where they randomly sample
88  thousands of architectures and choose promising ones for further training. In recent work,
89  Saxena & Verbeek (2016) propose to sidestep the architecture selection process through
90  densely connected networks of layers, which come closer to the performance of hand-crafted
91  networks. Google also proposed Efficient Neural Architecture Search(ENAS), a fast and
92  inexpensive approach for automatic model design. In ENAS, a controller discovers neural
93  network architectures by searching for an optimal subgraph within a large computational
94  graph. The controller is trained with policy gradient to select a subgraph that maximizes the
95  expected reward on a validation set. Details we will talk about later.

96  Recently there has been much work at the intersection of reinforcement learning and deep
97  learning. For instance, methods using CNNs to approximate the Q-learning utility function
98  (Watkins, 1989) have been successful in game-playing agents (Mnih et al., 2015; Silver et al.,
99  2016) and robotic control (Lillicrap et al., 2015; Levine et al., 2016). These methods rely on
100 phases of exploration, where the agent tries to learn about its environment through sampling,
101 and exploitation, where the agent uses what it learned about the environment to find better
102 paths. In traditional reinforcement learning settings, over-exploration can lead to slow
103 convergence times, yet over-exploitation can lead to convergence to local minima (Kaelbling
104 et al., 1996). However, in the case of large or continuous state spaces, the -greedy strategy of
105 learning has been empirically shown to converge (Vermorel & Mohri, 2005). Finally, when
106 the state space is large or exploration is costly, the experience replay technique (Lin, 1993)
107 has proved useful in experimental settings (Adam et al., 2012; Mnih et al., 2015). Bowen
108 incorporate these techniques—Q-learning, the $\epsilon$-greedy strategy and experience replay—in
109 algorithm design.

110

111 ## 3    AutoML

112 Automated Machine Learning(AutoML) is developed by Google and it mainly provides
113 methods and processes to make Machine Learning available for non-Machine Learning
114 experts, to improve efficiency of Machine Learning and to accelerate research on Machine
115 Learning. It's amazing if people do not need to select an appropriate model family and
116 optimize model hyper parameters. AutoML can also help you do these works: preprocess
117 clean the data, select and construct appropriate model family, post process machine learning

118　models and critically analyze the results obtained. Most machine learning experts will be
119　substituted by AutoML or other automated tools for machine learning. So I decide to do some
120　research for not be eliminated in the future. Hahaha…

121　AutoML did a lot of work in neural architecture search. Their team published many papers,
122　see the following website : https://www.ml4aad.org/automl/literature-on-neural-architecture
123　-search/.

124　They present an approach to automate the process of discovering optimization methods, with
125　a focus on deep learning architectures. They train a Recurrent Neural Network controller to
126　generate a string in a domain specific language that describes a mathematical update
127　equation based on a list of primitive functions, such as the gradient, running average of the
128　gradient, etc. The controller is trained with Reinforcement Learning to maximize the
129　performance of a model after a few epochs.

130
131　**3.1　Model Description**

132　CNN consists of many layers which can be described by string sequences, like "convolution,
133　RELU, max pooling, batch normalization …". This is very useful for large-scale model
134　generating by neural network. GoogleNet can be described as follows:

```
"op":"CONV:7X7:2X2:SAME",
"op":"MAX_POOLING:3X3:2X2:SAME",
"op":"LocalRespNorm",
"op":"CONV:3X3:1X1:SAME",
"op":"LocalRespNorm",
"op":"MAX_POOLING:3X3:2X2:SAME",
"op":[["CONV:1X1:1X1:SAME",["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SA
ME"],["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],["MAX_POOLING:3X3
:1X1:SAME","CONV:1X1:1X1:SAME"]],"DepthContact"],
"op":[["CONV:1X1:1X1:SAME",["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SA
ME"],["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],["MAX_POOLING:3X3
:1X1:SAME","CONV:1X1:1X1:SAME"]],"DepthContact"],
"op":"MAX_POOLING:3X3:2X2:SAME",
"op":[["CONV:1X1:1X1:SAME",["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SA
ME"],["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],["MAX_POOLING:3X3
:1X1:SAME","CONV:1X1:1X1:SAME"]],"DepthContact"],
"op":[["CONV:1X1:1X1:SAME",["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SA
ME"],["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],["MAX_POOLING:3X3
:1X1:SAME","CONV:1X1:1X1:SAME"]],"DepthContact"],
"op":[["CONV:1X1:1X1:SAME",["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SA
ME"],["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],["MAX_POOLING:3X
3:1X1:SAME","CONV:1X1:1X1:SAME"]],"DepthContact"],
"op":[["CONV:1X1:1X1:SAME",["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SA
ME"],["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],["MAX_POOLING:3X3
:1X1:SAME","CONV:1X1:1X1:SAME"]],"DepthContact"],
"op":[["CONV:1X1:1X1:SAME",["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SA
ME"],["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],["MAX_POOLING:3X3
:1X1:SAME","CONV:1X1:1X1:SAME"]],"DepthContact"],
"op":"MAX_POOLING:3X3:2X2:SAME",
"op":[["CONV:1X1:1X1:SAME",["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SA
ME"],["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],["MAX_POOLING:3X3
:1X1:SAME","CONV:1X1:1X1:SAME"]],"DepthContact"],
"op":[["CONV:1X1:1X1:SAME",["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SA
ME"],["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],["MAX_POOLING:3X3
:1X1:SAME","CONV:1X1:1X1:SAME"]],"DepthContact"],
"op":"AVERAGE_POOLING:7X7:1X1:VALID",
"op":"FC: RELU",
"op":"SOFTMAX",
"regularization":"Dropout",
"coefficient":"0.4"
```

135
136　Figure 4.CNN to String Sequence

137  Figure 4. is an example about translating GoogleNet architecture to string sequences.
138  GoogleNet has another name called as Inception(Figure 2.) which one inception module can
139  be described with:

140  [["CONV:1X1:1X1:SAME",

141  ["CONV:1X1:1X1:SAME","CONV:3X3:1X1:SAME"],

142  ["CONV:1X1:1X1:SAME","CONV:5X5:1X1:SAME"],

143  ["MAX_POOLING:3X3:1X1:SAME","CONV:1X1:1X1:SAME"]

144  ],

145  "DepthContact"].

146  It's interesting that CNN model can be described by such simple style. But we just have layer
147  operations without input, output and other operations like skip operation. Skip operation is in
148  ResNet which is very important for its architecture. Here we can use some special file like
149  JSON and XML to define model extending string sequence's shortage.

150
151  ## 3.2   Generate Model Representation

152  Next we can generate different model by generating different string sequences with recurrent
153  neural network. RNN is a class of artificial neural network where connections between nodes
154  form a directed graph along a sequence. This allows it to exhibit dynamic temporal behavior
155  for a time sequence. Unlike feedforward neural networks, RNNs can use their internal state
156  (memory) to process sequences of inputs. This makes them applicable to tasks such as
157  unsegmented, connected handwriting recognition or speech recognition. So, CNN
158  architecture can generate by RNN family.

159  In AutoML, an RNN controller is trained in a loop: the controller first samples a candidate
160  architecture, i.e. a child model (String sequence), and then trains it to convergence to
161  measure its performance on the task of desire which we will talk in the next section. The
162  controller then uses the performance as a guiding signal to find more promising architectures.
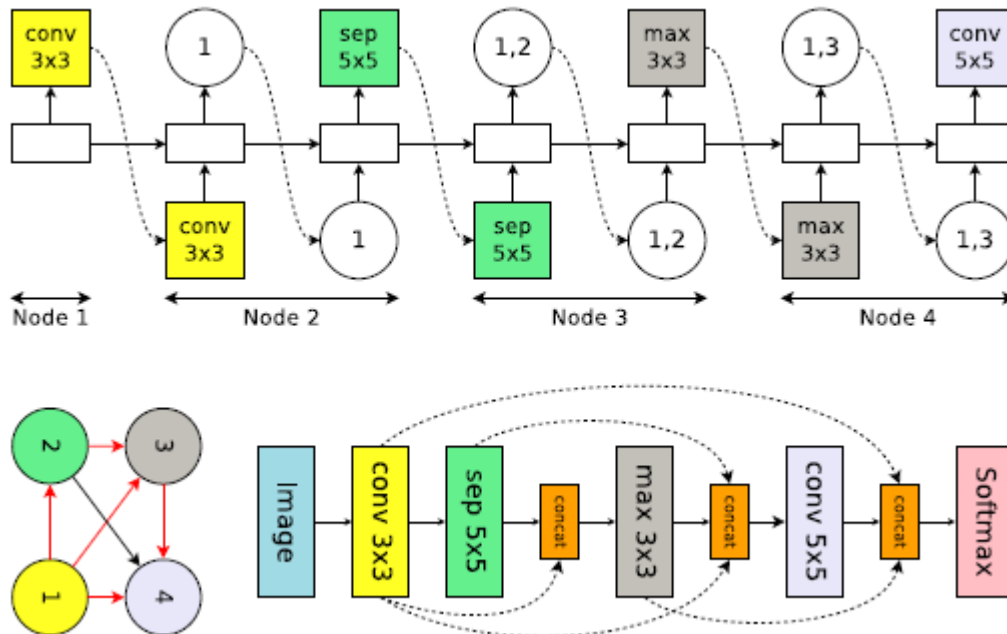


163

164                          Figure.5 Example

165  Figure 5. shows an example run of a recurrent cell in our search space with 4 computational
166  nodes, which represent 4 layers in a convolutional network. Top: The output of the controller
167  RNN. Bottom Left: The computational DAG corresponding to the network's architecture.

168    Red arrows denote the active computational paths. Bottom Right: The complete network.
169    Dotted arrows denote skip connections. The controller RNN samples two decisions at each
170    decision block: 1) what previous node to connect to and 2) what activation function to use. In
171    the search space for convolutional models, the controller RNN also samples two sets of
172    decisions at each decision block: 1) what previous nodes to connect to and 2) what
173    computation operation to use. These decisions construct a layer in the convolutional model.
174    We don't talk about details here.

175

176    **3.3    Neural Architecture Search**

177    There are two aspects to improve model performance, through increasing neural network
178    search space and generating better model. Increasing neural network model means that we
179    have more model representation to search. Generating better model means we can generate
180    model literately and better than previous version.

181    How can we increase neural network search space? Design convolutional cells and design
182    operation between layers. Rather than designing the entire convolutional network, one can
183    design smaller modules and then connect them together to form a network. Figure 6
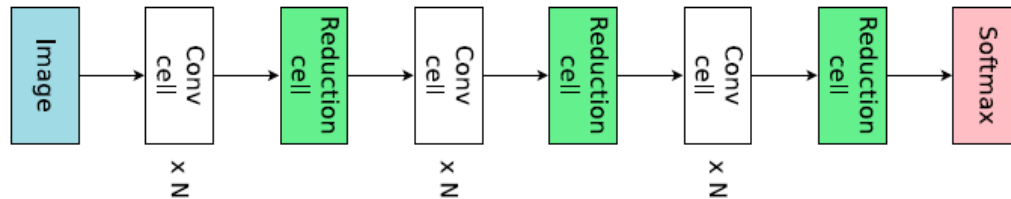184    illustrates this design, where the convolutional cell and reduction cell architectures are to be
185    designed.



186

187    Figure 6. Connecting 3 blocks, each with N convolution cells and 1 reduction cell, to make
188    the final network.

189    AutoML utilize the computational DAG with B nodes to represent the computations that
190    happen locally in a cell. In this DAG, node 1 and node 2 are treated as the cell's inputs,
191    which are the outputs of the two previous cells in the final network (see Figure 6). For each
192    of the remaining B − 2 nodes, they ask the controller RNN to make two sets of decisions: 1)
193    two previous nodes to be used as inputs to the current node and 2) two operations to apply to
194    the two sampled nodes. The 5 available operations are: identity, separable convolution with
195    kernel size 3×3 and 5×5, and average pooling and max pooling with kernel size 3×3. At each
196    node, after the previous nodes and their corresponding operations are sampled, the operations
197    are applied on the previous nodes, and their results are added.

198    How can we improve expected reward literately? Google tested many methods like grid
199    search, reinforcement learning, stochastic search, Bayesian search and so on. Here I just
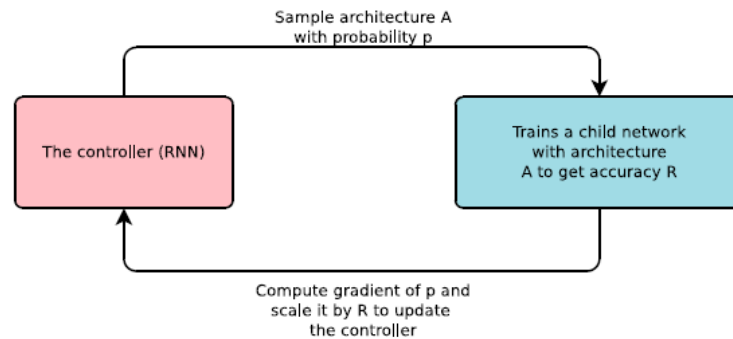200    introduce reinforcement learning.



201

202    Figure 7. An overview of neural architecture search with RL

The agent sequentially selects layers via the ε-greedy strategy until it reaches a termination state(See Figure 7). The CNN architecture defined by the agent's path is trained on the chosen learning problem, and the agent is given a reward equal to the validation accuracy. The validation accuracy and architecture description are stored in a replay memory, and experiences are sampled periodically from the replay memory to update Q-values. The agent follows an ε schedule which determines its shift from exploration to exploitation. Their method requires three main design choices: 1) reducing CNN layer definitions to simple state tuples, 2) defining a set of actions the agent may take, i.e., the set of layers the agent may pick next given its current state, and 3) balancing the size of the state-action space—and correspondingly, the model capacity—with the amount of exploration needed by the agent to converge. We do not talk details here.

## 4    Conclusions

Since I don't have much time to do detailed work, this report is a little simple and just summarizing the methods of automated deep learning. AutoML is the first automating machine learning system that did a lot of very instructive research. So I mainly introduce their work about how to automated search good convolutional neural network through two aspects: increasing neural architecture search space and improving model performance. First we can design unique CNN cells and create different block connections, second we can improve expected reward by reinforcement learning.

Everyone can do deep learning when deep learning can be automated done. It's so amazing that I am attracted by this topic, and also did some work that cannot write down here. This filed is still have a long way to go and I think also it will be very promising in the future.

### References

[1] Baker, Bowen, Gupta, Otkrist, Naik, Nikhil, and Raskar, Ramesh. Designing neural network architectures using reinforcement learning. In ICLR, 2017a

[2] Baker, Bowen, Otkrist, Gupta, Raskar, Ramesh, and Naik, Nikhil. Accelerating neural architecture search using performance prediction. *Arxiv, 1705.10823*, 2017b

[3] Bello, Irwan, Pham, Hieu, Le, Quoc V., Norouzi, Mohammad, and Bengio, Samy. Neural combinatorial optimization with reinforcement learning. In ICLR Workshop, 2017a.

[4] Brock, Andrew, Lim, Theodore, Ritchie, James M., and Weston, Nick. SMASH: one-shot model architecture search through hypernetworks. ICLR, 2018.

[5] Cai, Han, Chen, Tianyao, Zhang, Weinan, Yu, Yong., and Wang, Jun. Efficient architecture search by network transformation. In AAAI, 2018.

[6] Chollet, Francois. Xception: Deep learning with depthwise separable convolutions. In CVPR, 2017.

[7] Deng, Boyang, Yan, Junjie, and Lin, Dahua. Peephole: Predicting network performance before training. Arxiv, 1705.10823, 2017.

[8] DeVries, Terrance and Taylor, Graham W. Improved regularization of convolutional neural networks with cutout. Arxiv, 1708.04552, 2017.

[9] James Bergstra, Daniel Yamins, and David D Cox. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. ICML (1), 28:115–123, 2013.

[10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. arXiv preprint arXiv:1512.03385, 2015.

[11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In European Conference on Computer Vision, pp. 630–645. Springer, 2016.

[12] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. Nature, 521(7553):436–444, 2015.

[13] Ming Liang and Xiaolin Hu. Recurrent convolutional neural network for object recognition. CVPR, pp. 3367–3375, 2015.

[14] Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Learning to compose neural networks for question answering. In NAACL, 2016.

[15] Marcin Andrychowicz, Misha Denil, Sergio Gomez, MatthewWHoffman, David Pfau, Tom Schaul, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. arXiv preprint arXiv:1606.04474, 2016.

[16] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In ICLR, 2015.