

---

# From One to Zero: Causal Zero-Shot Neural Architecture Search by Intrinsic One-Shot Interventional Information

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 “Zero-shot” neural architecture search (ZNAS) is key to achieving real-time neural  
2 architecture search. ZNAS comes from “one-shot” neural architecture search but  
3 searches in a weight-agnostic supernet and consequently largely reduce the search  
4 cost. However, the weight parameters are agnostic in the zero-shot NAS and none  
5 of the previous methods try to explain it. We question whether there exists a  
6 way to unify the one-shot and zero-shot experiences for interpreting the agnostic  
7 weight messages. To answer this question, we propose a causal definition for “zero-  
8 shot NAS” and facilitate it with interventional data from “one-shot” knowledge.  
9 The experiments on the standard NAS-bench-201 and CIFAR-10 benchmarks  
10 demonstrate a breakthrough of search cost which requires merely **8 GPU seconds**  
11 **on CIFAR-10** while maintaining competitive precision.

## 12 1 Introduction

13 Neural architecture search has been an interesting topic in the AutoML community [27]. Traditional  
14 methods search by training the distinct neural architecture iteratively [31] whose training cost is  
15 huge. One-shot model cleverly use a supernet to merge all the singular neural architectures into  
16 one and consequently, the waste of search time is largely saved [16]. Further, the gradient-based  
17 one-shot method [12] is proposed which acquires robust results on NASNet [32]. Though the one-shot  
18 model largely reduces the search cost, it still suffers from a weight-sharing problem, and especially,  
19 gradient-based approaches cause degenerate architectures [29]. The work [25] gives theoretical proof  
20 for this and subtly uses a progressive tuning metric to discretize the one-shot supernet iteratively  
21 which gets awesome neural architectures. However, it still gets degenerate architectures with different  
22 training settings.

23 The brilliant work [5] from Google Brain gives a hint for searching neural networks without tuning the  
24 parameters. “To produce architectures that themselves encode solutions, the importance of weights  
25 must be minimized”. In this manner, a zero-shot neural architecture search (ZNAS) is born. The  
26 work [10] firsts propose the idea of ZNAS to be “it does not optimize network parameters during  
27 search”. From a one-shot perspective, the “zero-shot” is given credit by “one-shot” where single  
28 neural architectures are supposed to be selected from the weight-agnostic supernet [5]. Considering  
29 causal weight messages, the “zero-shot” select neural architecture with the minimum impact of  
30 any weight parameter [5]. Thus a causal definition is supposed to be that the weight messages are  
31 multi-environmentally distributed. Compared to one-shot NAS, zero-shot NAS gets imperfect weight  
32 messages due to random initialization and searching without training [10, 2].

33 A training-free approach is first proposed by the work [13]. Different from the previous zero-shot  
34 model [10], the work [13] samples well-trained architectures and get validation accuracy to train the

35 statistical proxy before it searches. The work [2] follows the way of the previous work [10] and uses  
 36 the DARTS search space to conduct zero-shot NAS on CIFAR-10 and ImageNet in a training-free  
 37 manner. However, the number of samples directly decides the belief of the final precision. The  
 38 “well-trained” architectures might not be “perfectly-trained” in different training settings.

39 Zero-shot NAS learns the representation of neural architectures to get the best one. Consistently  
 40 compared to one-shot NAS methods, zero-shot NAS methods ignore the weight information. By  
 41 merely measuring the architectural expressivity, they overlooked the impact of weights as a necessary  
 42 assessment element. From a one-shot NAS perspective, architectural information can be represented  
 43 by a list of neuron representations [25]. The message of training weights  $\omega$  supports the neuron’s  
 44 representation [15, 12, 25]. Because the structural dependencies of shared (mutual) messages across  
 45 neurons are all agnostic [5], in the zero-shot neural architecture search, the neuron’s representation is  
 46 harder to interpret due to the random messages. What is worse, the uninterpretability might result in  
 47 large bias and variances because the imprecise observational data might be misleading. Finally, it  
 48 will lead the search to get degenerate architectures through the process of accumulating errors.

49 We first propose to interpret the zero-shot NAS in a causal-representation-learning setting. According  
 50 to the weight-agnostic setting, we formulate the zero-shot NAS as a novel framework for imperfect-  
 51 information NAS. The structural information of zero-shot NAS is interpreted by impact with the  
 52 latent factors. As a consequence, intrinsic high-level interventional data acquired by one-shot NAS  
 53 is properly adopted to refine the imperfectness. Moreover, we reformulate the causality by game  
 54 theory and interpret the imperfect-information NAS as imperfect information game  $\mathcal{G}$ . Extensive  
 55 experiments on various benchmark datasets including CIFAR-10, NAS-Bench-201, and ImageNet  
 56 have shown the super search efficiency ( $10000\times$  faster than DARTS) of our methods. In this work,  
 57 our main contributions are as follows:

- 58 • We propose that the causal zero-shot NAS is to learn the neuron’s representation with latent  
 59 factors in observationally imperfect messages.
- 60 • We theoretically demonstrate the validation information of either a neuron or a neuron  
 61 ensemble obeys a Gaussian distribution given a Gaussian input.
- 62 • The proposed method uses high-level interventional data from one-shot NAS to facilitating  
 63 zero-shot NAS to solve the imperfectness.
- 64 • Our method sets the new state-of-the-art in zero-shot NAS of search cost (8 GPU seconds)  
 65 while maintaining comparable test accuracies.

## 66 2 Preliminaries and Related Work

67 In this section, we talk about the preliminaries and the previous works on one-shot NAS and zero-shot  
 68 NAS. We talk about the motivation to replace statistical proxy by introducing the basic knowledge on  
 69 causal interventional representation learning in causality [20, 1].

### 70 2.1 One-shot NAS

71 One-shot NAS methods [12, 16], that unify all the single-path neural architectures into one super-  
 72 network  $\mathcal{S}$  (supernet), select the single-path neural architecture as the best one by training the weights  
 73  $\omega$  in a weight-sharing manner and maximizing the validation accuracy ( $\mathcal{V}$ ) of architecture  $\mathcal{A}$  as  
 74 follows:

$$74 \quad \text{Max}_{\mathcal{A}}(\mathcal{V}(\mathcal{A}, \bar{\omega})) \quad \text{s.t.} \quad \bar{\omega} = \omega + \delta_{\mathcal{A}}\omega_{\mathcal{S}} \quad (1)$$

75 The iterative updating of  $\omega$  and selection of  $\mathcal{A}$  makes the one-shot NAS a bi-level optimization  
 76 problem that is NP-hard. Differentiable one-shot model also relies on the observational data from  
 77 unidely trained validation accuracies of differentiable subnets [12]. Wang et al. [25] propose a  
 78 selection-based approach to modify the output of differentiable one-shot NAS [12] to discretize a  
 79 single-path neural architecture that consists of operations (neurons) with strength. As a consequence,  
 80 the perturbation-based inductive bias is demonstrated to be helpful to solve the degeneration.

### 81 2.2 Statistical proxies in zero-shot NAS

82 We compare the various training-free and zero-shot NAS methods according to the usage of statistical  
 83 representation. Some training-free approaches use the statistic of validation accuracy to predict the

84 final architecture. NASWOT [13] samples a number ( $N$ ) of well-trained neural architectures from  
85 the NAS-Bench-201 dataset to learn the kernel. However, to get these representations, the training  
86 costs tremendously. The zero-shot methods directly use zero-cost statistical proxies to represent the  
87 expressivity without weights and validation accuracy. Zen-NAS [10] uses a Gaussian complexity to  
88 measure the network expressivity and evolve the architectures to maximize the expressivity. Other  
89 training-free approaches such as TE-NAS [2] and NASI [22] imitate the train of NAS by neural  
90 tangent kernel (NTK) which largely reduces the waste of train cost. TE-NAS [2] propose to maximize  
91 the number of linear region of activation patterns [14]. On the opposite, NASI [22] subtly optimize  
92 the trace of NTK by sampling.

93 Here raise the question that to what extent the validation accuracy outperforms the statistical proxy.  
94 Vice versa, we question if the statistical proxy is in substitute of the validation accuracy. Compared to  
95 the proxy-based methods with approximations, the validation-based method is more reproducible. The  
96 validation accuracy is an intrinsic robust and upper-bounded proxy to measure the neural architectures.  
97 Besides, previous arts of one-shot manner usually use the validation accuracy to be the objective to  
98 maximize. Despite these benefits, the zero-shot representation is imperfect due to the weight-agnostic  
99 messages.

## 100 2.3 Causal representation learning

101 The study [20] demonstrates that causality is a “subtle concept” which can not be fully described  
102 by Boolean or Probabilistic. It is more about reasoning. Reichenbach demonstrates a common  
103 cause principle to explain the causality by dependencies among variables [18]. Causal representation  
104 learning mainly deals with learning causally for representations. By observational data, we can hardly  
105 learn the real circumstances (environments), especially in complex scenes and high-dimensional data  
106 scenarios. Causal representation learning seeks to extract high-level information (dependencies) from  
107 low-level data. Interventions have taken a prominent role in representation learning literature on  
108 causation. The work [1] uses interventional data to facilitate the causal representations to get precise  
109 outcomes. Neural architecture search aims at learning the architectural representations automatically.  
110 The automatism of the previous arts of neural architecture search might not be causal especially in  
111 zero-shot setting.

## 112 3 Method

### 113 3.1 Imperfect information

114 Neural architecture search is a task aiming at interpreting the mechanism of architectural knowledge  
115 of neural networks given methods of evaluations. Activation patterns, statistical proxies, and naive  
116 validation accuracy are adopted to evaluate the score of a neural network. However, we can hardly  
117 understand any neural network and even hardly explain the weight distribution of any neural network  
118 without assumptions. Observational data are always imperfect due to the infinite environments (search  
119 spaces/training schemes/hardware/etc.) of all possible networks with finite observations and limited  
120 tools. Architecture information is not stand-alone.

121 In one-shot NAS, demonstrated in Equation 4, given a neural network, we first train the weights  $\omega$   
122 and the  $\omega$  combined with architecture  $\mathcal{A}$  can give a validation accuracy  $\mathcal{V}$ . After  $\mathcal{V}$  is given, we then  
123 update the  $\omega$  to get  $\bar{\omega}$  and a novel architecture  $\mathcal{A}$  until the validation accuracy  $\mathcal{V}$  is maximum. In  
124 the train, the architecture of a neural network is the key factor that impacts the other two factors  $\omega$   
125 and validation accuracy  $\mathcal{V}$ . The search is actually a reverse way of train to the aspect of the intrinsic  
126 dependency of accuracy  $\mathcal{V}$  on the weight  $\omega$  and architecture  $\mathcal{A}$ . However, we have overlooked a lot of  
127 factors like data distributions, batch sizes, rates of weight decay, and so on and on which we can not  
128 optimize as “one shot”. If the observational data alone can not interpret the phenomenon, it is a must  
129 to model the latent factors  $\mathcal{Z}$  that cause this uninterpretability. Figure 1 illustrates the dependencies  
130 of architecture  $\mathcal{A}$ , validation accuracy  $\mathcal{V}$ , and weights  $\omega$ . The dashed line reveals that  $\mathcal{Z}$  changes the  
131 dependencies of selected neurons (or searched architectures) on observational data of  $\omega$  and  $\mathcal{V}$  [23],  
132 which indeed implies strong causality [20]. In logical condition, the structural relationship between  
133  $\mathcal{V}$  and  $\omega$  can be almost broken<sup>1</sup>.

---

<sup>1</sup>See demonstration in Section 3.3, results in Section 4.

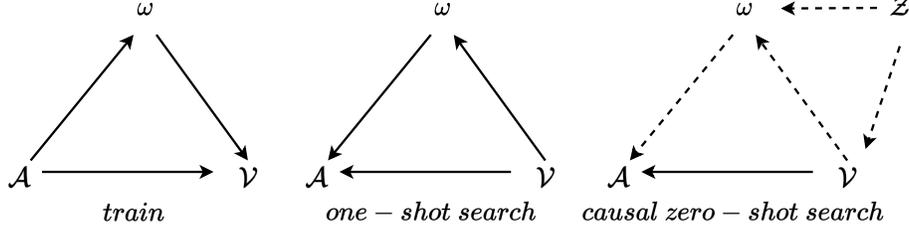


Figure 1: Illustrations of the dependencies of architecture  $\mathcal{A}$ , validation accuracy  $\mathcal{V}$ , and weights  $\omega$  with latent factor  $\mathcal{Z}$  on the train (left), one-shot neural architecture search (middle), and causal zero-shot neural architecture search (right).

134 We assume the validation accuracy  $\mathcal{V}$  of a set of neural architectures  $\{\mathcal{A}\}$  obeys a Gaussian distribution.

$$135 \mathcal{V} \sim \mathcal{N}(\mu, \sigma^2) \quad (2)$$

136 Due to the random weight information, artificial neural networks (ANN) themselves have architectural  
 137 information to deliver the neural networks' expressivity with large variances [5]. It is demonstrated  
 138 that the weight-agnostic neural network still preserves the 92% accuracy-level information for digit  
 139 classification by the work [5]. However, the weights are agnostic and consequently the validation  
 140 accuracies are imperfect. We assume the true validation accuracy is the difference of the observational  
 141  $\mathcal{V}^{obs}$  and latent impact of factor  $\mathcal{Z}$  demonstrated in Equation 3.

$$\mathcal{V} \sim \mathcal{N}(\mu_{obs} - \mu_{\mathcal{Z}}, \sigma_{obs}^2 - \sigma_{\mathcal{Z}}^2) \quad (3)$$

### 142 3.2 Problem formulation

143 In Zen-NAS, the adoption of statistical proxy on the feature map is impressive while it is constrained  
 144 to structural dependencies [10]. We question to what extent, when we search a neural network, the  
 145 statistical proxies can be replaced with the more robust functions such as validation accuracy causally  
 146 [20]. In some one-shot [16, 12] and training-free methods [13], the evaluation metrics are usually the  
 147 validation accuracy of the associated neural architectures.

148 Inspired by the previous work [25], we evaluate each neuron to select respectively in substitute.  
 149 Intuitively, we measure the importance of each neuron by a simple validation accuracy of a singular  
 150 associate neuron while resting other neurons on the same edge. DARTS+PT [25] the perturbation-  
 151 based approach mutes the irrelevant neurons to conduct an inference while saving the other paralleled  
 152 edges. For each paralleled edge (layer)  $\mathcal{E}$  that contains  $M$  neurons  $\mathcal{N}$ s, we mute the other neurons  
 153 while only saving the  $i^{th}$  neuron  $\mathcal{N}_{(i)}$ . The  $k^{th}$  paralleled edge  $\mathcal{E}_i^{(k)}$  consequently only contains one  
 154 neuron (operation):  $\mathcal{E}_i^{(k)} = \{0 \times \mathcal{N}_{(1)}, 0 \times \mathcal{N}_{(2)}, \dots, \mathcal{N}_{(i)}, \dots, 0 \times \mathcal{N}_{(M)}\}$ . When saving the other  
 155 paralleled edges  $\{\mathcal{E}_{(j)}\}_{j \neq k}$ ,  $\mathcal{N}_{(i)}$  denotes any single sub-architecture (a neuron) in the supernet  $\mathcal{S}$   
 156 with tuned weights  $\omega_{\mathcal{S}}$  of the supernet. Formally, the one-shot neuron selection for  $k_{th}$  paralleled  
 157 edge is defined as:

$$\mathcal{N}^* = \operatorname{argmax}(\mathcal{F}(\{\mathcal{V}(\mathcal{N}_{(i)}, \omega_{\mathcal{S}})\})) \quad \forall \mathcal{N}_{(i)} \in \mathcal{E}^{(k)} \quad (4)$$

158 where validation accuracy  $\mathcal{V}$  is measured by an intrinsic inductive bias function  $\mathcal{F}$  such as a rein-  
 159 forcement learning policy  $\pi$  [31, 32].  $\mathcal{V}(\mathcal{N}_{(i)}) = \mathcal{V}(\{\mathcal{E}^{(1)}, \mathcal{E}^{(2)}, \dots, \mathcal{E}_i^{(k)}, \dots, \mathcal{E}^{(N)}\})$  in practise.

160 In zero-shot NAS, the weight information is agnostic, which is impacted by a latent factor  $\mathcal{Z}$  as  
 161 shown in Figure 1. [4]. The latent variable obeys a distribution  $\mathcal{P}$  in dimension  $\Lambda$ :

$$\mathcal{Z} \sim \mathcal{P}^{\Lambda} \quad (5)$$

162 When we sample larger enough numbers of impacts, the sample of factor  $\mathcal{Z}$  obeys a Gaussian  
 163 distribution by the central limit theorem (CLT). The causal zero-shot neural architecture search  
 164 (Causal-Znas) that searches in imperfect messages is defined as:

$$\mathcal{N}^* = \operatorname{argmax}(\mathcal{F}(\{\mathcal{V}(\mathcal{N}_{(i)}, \omega)\}|\mathcal{Z})) \quad \forall \mathcal{N}_{(i)} \in \mathcal{E}^{(k)} \quad (6)$$

165 for  $i = 1, 2, \dots, M$ . In this Equation 6,  $\mathcal{Z}$  means the latent information to impact agnostic-weights  
 166 (such as a random initialization [5, 10]) and consequently validation accuracies  $\mathcal{V}$ . Therefore, we get a  
 167 causal information set of singular neuron representation  $\{\mathcal{V}(\mathcal{N}_i)|\mathcal{Z}\}$  for  $i = 1, 2, \dots, M$ . For each  
 168 paralleled edge (layer)  $\mathcal{E}$  that contains  $M$  neurons  $\mathcal{N}$ :  $\mathcal{E} = \{\mathcal{N}_{(1)}, \mathcal{N}_{(2)}, \dots, \mathcal{N}_{(M)}\}$ . We calculate  
 169 the information of singular neuron  $\mathcal{N}_i$  on edge  $\mathcal{E}^{(j)}$  by freezing the other layers (ensembles/edges)  
 170  $\{\mathcal{E}^{(k)}\}_{k \neq j}$  so that the causal information is only impacted by the current neurons due to the same  
 171 condition (in the same paralleled edge). Then the causal information set of a paralleled edge  $\mathcal{E}$  is as:

$$\{\mathcal{V}(\mathcal{E})|\mathcal{Z}\} = \{\mathcal{N}_{(1)}(\mathcal{X}|\mathcal{Z}), \mathcal{N}_{(2)}(\mathcal{X}|\mathcal{Z}), \dots, \mathcal{N}_{(M)}(\mathcal{X}|\mathcal{Z})\} \quad (7)$$

172 In a Causal-Znas, a prediction function  $\mathcal{F}$  is able to measure the selected architectures from the  
 173 un-trained supernet. To avoid the improper introduction of inductive biases, we use an identity  
 174 function to measure the importance of neurons.

### 175 3.3 Gaussian intervention

176 Most existing NAS approaches use observational data and make assumptions on the architectural de-  
 177 pendencies to achieve provable representation identification. However, in our causal zero-shot neural  
 178 architecture search, there is a wealth of interventional data available. To perfect the observational  
 179 validation accuracies  $\mathcal{V}^{observed}$  in  $\mathcal{D}$ , we sample  $\mathcal{V}^{ven}$  from an interventional distribution  $\mathcal{D}(\mathcal{Z})$  to be in  
 180 substitute for the ones derived by the observation  $\mathcal{V}^{observed}$ . Formally, we have:  $\mathcal{V}^{ven} \sim \mathcal{D}(\mathcal{Z})$ . Though  
 181 pure architectural information is imperfectly observed, we can use an interventional function  $\mathcal{I}$  (**do**  
 182 **intervn** [1]) to replenish data from a one-shot perspective:

$$\mathcal{V} = \mathcal{I}_p^{\mathcal{D}(\mathcal{Z})} \mathcal{V}^{ven} \cup \mathcal{I}_{1-p}^{\mathcal{D}} \mathcal{V}^{observed} \quad (8)$$

183 Ming et al. [10] assume the inputs obey Gaussian distribution and get comparable results with  
 184 one-shot methods [12, 16]. What we use as the input for each neuron is a Gaussian image which also  
 185 obeys the assumption of Gaussian inputs of Zen-NAS [10].

186 **Lemma 1.** *Given a Gaussian input  $\mathcal{X} \sim \mathcal{N}(\mu, \sigma^2)$ , the output of a neuron  $\mathcal{N}$  in the first layer is*  
 187 *Gaussian.*

188 *Proof.* Assuming each neuron is a distinct convolution denoted as  $Conv_i$  for  $i = 1, 2, \dots, M$ , then  
 189 the output of this edge is:

$$\mathcal{O} = \sum_{i=1}^M (\{Conv_{(1)}(\mathcal{X}, \mathcal{W}_{(1)}), Conv_{(2)}(\mathcal{X}, \mathcal{W}_{(2)}), \dots, Conv_{(M)}(\mathcal{X}, \mathcal{W}_{(M)})\}) \quad (9)$$

190 where  $\mathcal{X} \sim \mathcal{N}(\mu, \sigma^2)$  and  $\mathcal{W}_{(i)} \sim \mathcal{N}(\mu_w, \sigma_w^2)$  for  $i = 1, 2, \dots, M$ . Given the i.i.d. inputs  
 191 and weights, the output score (validation accuracy) of the neural network layer is Gaussian since  
 192 the Convolution of a Gaussian (random variable) is still a Gaussian (random variable). We have  
 193 Gaussian weights  $\mathcal{W}_{(i)}$  and  $Conv_{(i)}(\mathcal{X}, \mathcal{W}_{(i)}) \sim \mathcal{N}(\mu_{(i)}, \sigma_{(i)}^2)$ . Then  $\sum_i Conv_{(i)}(\mathcal{X}, \mathcal{W}_{(i)}) \sim$   
 194  $\mathcal{N}(\sum \mu_{(i)}, \sum \sigma_{(i)}^2)$ .  $\square$

195 **Lemma 2.** *Given a Gaussian input  $\mathcal{X} \sim \mathcal{N}(\mu, \sigma^2)$ , the output of a neuron  $\mathcal{N}$  in any layer is*  
 196 *Gaussian.*

197 *Proof.* Apparently, any weighted summation of random variables that obey two distinct Gaussian is  
 198 still a Gaussian. In neural networks, the layers are stacked. Based on Lemma 1, in the latter layer,  
 199 the outputs also obey the Gaussian, whose inputs are the former layer's outputs. The convolution  
 200 (neuron)  $Conv'_{(i)}$  of the next layer with output of latter layer  $\mathcal{O}$  (in Equation 9) has  $Conv'_{(i)}(\mathcal{O}) \sim$   
 201  $\mathcal{N}(\mu'_{(i)}, \sigma'_{(i)}{}^2)$ .  $\square$

202 **Corollary 2.1.** *Given a Gaussian input  $\mathcal{X} \sim \mathcal{N}(\mu, \sigma^2)$ , the output of any neuron ensemble*  
 203  *$\{\mathcal{N}_{(i)}\}_{i \in \mathcal{M}}$  is Gaussian.*

204 Formally, we have  $\mathcal{O}^{(i)} \sim \mathcal{N}^{(i)}(\mu', \sigma'^2)$ .  $\tilde{\mathcal{O}} = \{\mathcal{O}^{(1)}, \mathcal{O}^{(2)}, \dots, \mathcal{O}^{(K)}\}$  where  $\tilde{\mathcal{O}}$  denotes all the  
 205 outputs across edges  $\overbrace{\mathcal{E}_{(1)}, \mathcal{E}_{(2)}, \dots, \mathcal{E}_{(K)}}$ . Based on Lemma 1 and Lemma 2, we get the Corollary 2.1  
 206 to select edges (topology preferences).

207 *Proof.* By Lemma 1, we have any neuron  $\mathcal{N}_{(i)}$  has a Gaussian output  $\mathcal{O}^{(i)} \sim \mathcal{N}(\mu_{(i)}, \sigma_{(i)}^2)$ . Any  
 208 ensemble of neurons has an output  $\sum_i \mathcal{O}^{(i)}$ . Then we have  $\sum_i \mathcal{O}^{(i)} \sim \mathcal{N}(\sum \mu_{(i)}, \sum \sigma_{(i)}^2)$ .  $\square$

209 As demonstrated in Equation 8, we propose an intervention function  $\mathcal{I}^{\mathcal{D}}$  to facilitate the imper-  
 210 fect causal representation of the validation information. We propose that the ideal information is  
 211 distributed in the information set by a probability  $p$ . The distribution  $\mathcal{D}$  is  $\mathcal{N}(\mu, \sigma)$  in the context.

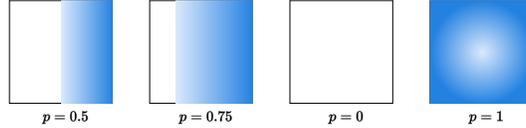


Figure 2: Illustration of intervention of observational data. The blue denotes interventional data while the white denotes observational data.

212 Herein, we question to what extent, the imperfectness can be interventionaly refined [1]. We use  
 213 the parameter  $p$  to asymmetricly flipping the random Gaussian  $\mathcal{I}_p^{\mathcal{N}(\mu, \sigma^2)}$  [15] to understand the  
 214 imperfect information in dimension  $\Lambda$  which is mapped to a vanilla Gaussian (in Equation 5). As  
 215 shown in Figure 2, it compares the information difference between the observational information set  
 216 and interventional information set impacted by the parameter  $p$ . In different environments, the data  
 217 of interventional data combined with observation obeys a distinct Gaussian, which implies strong  
 218 coherence and robustness. When  $p = 1$ , the causality is perfectly achieved due to breaking the  
 219 dependency of validation accuracy  $\mathcal{V}$  on weights  $\omega$ ; otherwise, it is imperfect. The mean and variance  
 220 coefficients of the additional notion of intervention are derived by sampling validation accuracy of  
 221 one-shot prior. We propose that setting of  $p$  is conditional on the fraction of the mean of latent factor  
 222 to the difference of the mean of observational data and the mean of interventional data.

223 **Proposition 1.** When  $p \rightarrow \frac{\mu_{\mathcal{Z}}}{\mu_{\text{observer}} - \mu_{\text{ven}}}$ , the mean of the intervened data  $\tilde{\mu} \rightarrow \mu_{\text{true}}$ .

224 As demonstrated in Proposition 1, a sufficient condition of the mean of intervened data is getting  
 225 closer to the true mean of the validation accuracy is that the  $p$  is closer to 1 and interventional data is  
 226 closer to the true data.

### 227 3.4 Causal zero-shot neural architecture search

228 We formulate the zero-shot NAS into ensemble selection and neuron selection. There are  $K$  neuron  
 229 ensembles  $\{\mathcal{N}_{(i)}\}_{i \in \mathcal{M}}^{(1)}, \{\mathcal{N}_{(i)}\}_{i \in \mathcal{M}}^{(2)}, \dots, \{\mathcal{N}_{(i)}\}_{i \in \mathcal{M}}^{(K)}$ . For each ensemble, there are  $M$  neurons  
 230 (operations). The ensemble selection is the selection of an ensemble  $\{\mathcal{N}_{(i)}\}_{i \in \mathcal{M}}^{(j)}$   
 231 the  $K$  ensembles ( $j \in \mathcal{K}$ ), while neuron selection follows the same formula and selects a neuron  $\mathcal{N}_{(i)}$   
 232 from a neuron ensemble  $\{\mathcal{N}_{(i)}\}_{i \in \mathcal{M}}^{(j)}$ .

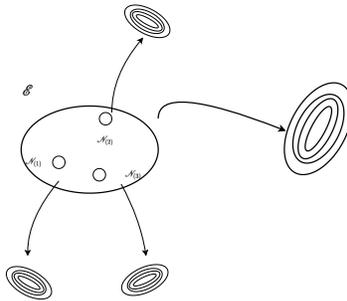


Figure 3: The distribution plate of three neurons and a big distribution plate of ensemble of them.

---

**Algorithm 1** Causal zero-shot neuron selection.

---

Initialize supernet weights  $\omega$ ;

For  $i = 1, 2, \dots, M$ :

    Calculate validate accuracy  $\mathcal{V}^{observed}(\mathcal{N}_{(i)}(\omega))$ ;

**do interv**n by  $p$ ;

    Maximize the  $\mathcal{V}$  and select the  $\mathcal{N}^*$ .

---

233 As is shown in Figure 3, the validation accuracy of both a neuron and a neuron ensemble obey  
234 Gaussian distributions respectively. From a macro perspective it is an ensemble selection while from  
235 a minor perspective, it is a neuron selection. Thus we talk about both types in the same formula.

236 As demonstrated in Equation 6, the final outcome neurons are derived by maximizing their validation  
237 accuracies according to the latent factor. Given the Gaussian intervention in Equation 8, we further  
238 modify the formula of the causal neuron selection by doing intervention (without the additional  
239 inductive bias [20]):

$$\tilde{\mathcal{N}}^* = \operatorname{argmax}(\{\tilde{\mathcal{V}}(\mathcal{N}_{(i)})\}_{i \in \mathcal{M}}) \quad (10)$$

240 , where  $\tilde{\mathcal{V}}$  is the validation accuracy with intervention.

241 The methodology of neuron selection is given in Algorithm 1. The search process of neuron ensemble  
242 follows the same formulation as mentioned in this Section. **do interv**n represents to do intervention.  
243 At first, the weight  $\omega$  of the supernet is randomly initialized [10]. Second, validation scores  $\mathcal{V}$  on the  
244 validation set are prepared for the calculation of the neurons  $\mathcal{N}$  which adopts probability  $p$  to do the  
245 intervention. At last, the maximum of values is compared to select the best neuron (operation). In  
246 practice, when the probability  $p$  is close to 1, the validation accuracy of observation has less need to  
247 compute.

248 Equation 6 reveals a universal formula for causal neural architecture search in the zero-shot settings.  
249 The measure function  $\mathcal{F}$  measures the importance [25] (“responsibility”) of a neuron and Shapley  
250 value is proposed to be ideal for the selection of a neuron [7] or ensemble [19].

$$\mathcal{N}^* = \operatorname{argmax}(\{\mathcal{G}_{(i)}(\{\tilde{\mathcal{V}})\}_{i \in \mathcal{M}}\}) \quad (11)$$

251 We use the game-theoretic inductive bias to extract the valuable information [20, 7].  $\mathcal{G}$  represent the  
252 Shapely value [21]. Given Corollary 2.1, we know that any the neuron ensemble obeys a Gaussian  
253 distribution. The information set of Shapley value is thus build on top of an ensemble of Gaussian  
254 variables. However, we could not guarantee a Gaussian distribution of the Shapley value [24]. As  
255 a consequence, we use a Gaussian distribution to do intervention on validation accuracy and then  
256 calculate the Shapely value of the intervened validation accuracy. At last, the Shapley value is  
257 maximized whose associated neuron is supposed to be more expressive [7].

### 258 3.5 Weight-agnostic weights

259 In the assumptions of various methods, weights are initialized as Gaussian. However, in our frame-  
260 work, we demonstrate that this strong assumption is not a must. Supernet can be initialized in different  
261 ways: i) with Gaussian [10], ii) Uniform [5], and iii) Constant number [5].

262 **Corollary 2.2.** *Given a Gaussian input  $\mathcal{X} \sim \mathcal{N}(\mu, \sigma^2)$ , if the initial weights are Uniform or Constant  
263 number  $C$ , the output of any neuron ensemble  $\{\mathcal{N}_{(i)}\}_{i \in \mathcal{M}}$  is not Gaussian.*

264 *Proof.* Apparently, the convolution of a Gaussian input with constant or uniform weights obeys a  
265 difference of CDF  $\Phi$  of the Gaussian in the range of constant or uniform.  $\square$

266 In the previous work [5], it is proposed that weights are supposed to be initialized by a distribution  
267 but not a constant ( $C$ ). To be more precise, we propose that the constant value could not represent the  
268 agnostic weights and thus could not reflect the latent information while a uniform distribution can  
269 guarantee the randomness. By training on a “wide range” of uniform weight samples, Gaier et al.  
270 propose that “the best performing values were outside of this training set” [5]. We propose that this  
271 phenomenon is essentially resulted from a distribution shift of the Gaussian validation accuracy which  
272 causes the change of search procedure. To solve the distribution shift, we could use the difference of  
273 CDF of Gaussian ( $\Phi$ ) to conduct intervention. Even in a broader view, if the weights distributions are  
274 totally unknown, we can use Bayesian method to approximate a distribution  $\mathcal{D}(\mathcal{Z})$  in Equation 8.

275 **4 Experiments**

276 We present the results and all experiment details of our method in this section. A robustness analysis  
 277 is included to examine the stability of our method, which also explains the time efficiency. Results  
 278 are given on the benchmark datasets, NAS-Bench-201 and CIFAR-10.

279 **4.1 Experimental details**

280 We use the search space of DARTS [12] for fair comparisons with the state-of-the-art NAS approaches.  
 281 During the searching process, we follow adopting the **same** and hyper-parameters as DARTS [12]  
 282 to initialize the supernet on the CIFAR-10 and NAS-Bench-201 datasets for a fair comparison with  
 283 DARTS-variants (one-shot methods). All the training is conducted on a single 2080Ti GPU.

284 **4.2 Results on CIFAR-10**

Table 1: Comparison with state-of-the-art NAS methods on CIFAR-10.

Algorithm	Test Error (%)	Params (M)	Search Cost (GPU seconds)	Search Strategy
DenseNet-BC [6]	3.46	25.6	-	manual
NASNet-A + cutout [32]	2.65	3.3	$1.73 \times 10^8$	RL
AmoebaNet-A [17]	$3.34 \pm 0.06$	3.2	$2.72 \times 10^8$	GA
AmoebaNet-B [17]	$2.55 \pm 0.05$	2.8	$2.72 \times 10^8$	GA
PNAS [11]	$3.41 \pm 0.09$	3.2	$1.94 \times 10^7$	SMBO
ENAS [16]	2.89	4.6	43200	RL
DARTS(1st) [12]	$3.00 \pm 0.14$	3.3	34560	gradient
DARTS(2nd) [12]	$2.76 \pm 0.09$	3.3	86400	gradient
BayesNAS [30]	$2.81 \pm 0.04$	3.4	17280	gradient
DrNAS [3]	<b>2.54 ± 0.03</b>	4.0	34560	gradient
ISTA-NAS [26]	$2.54 \pm 0.05$	3.3	4320	gradient
DARTS+PT [25]	$2.61 \pm 0.10$	3.0	69120	gradient
TE-NAS [2]	$2.63 \pm 0.06$	3.8	4320	NTK
NASI-FIX [22]	$2.79 \pm 0.01$	3.9	864	NTK
NASI-ADA [22]	$2.90 \pm 0.01$	3.7	864	NTK
Causal-Znas( $p = 0.5$ )	$2.89 \pm 0.08$	<b>2.6</b>	142	causal
Causal-Znas( $p = 1$ )	$2.75 \pm 0.10$	3.2	<b>8</b>	causal
Causal-Znas-G( $p = 1$ )	$2.61 \pm 0.04$	3.1	30	causal

285 As shown in Table 1, we compare the proposed Causal-Znas and game-version Causal-Znas-G with  
 286 the state-of-the-art methods. The comparisons are made with respect to the informatics of the model,  
 287 including test accuracy on the test set (Test Error), the number of parameters (Params), the search  
 288 costs, and the search strategies. As shown, our results set the new state-of-the-art search speed with a  
 289 competitive test error rate. Compared to DARTS [12], our method is 10000× faster with comparable  
 290 accuracy (2.75% v.s. 2.76%). Compared to DARTS+PT [25], our model is much simpler without  
 291 introducing the perturbation-based inductive bias [20] and achieves a similar test error rate (2.61%  
 292 v.s. 2.61%). DrNAS [3] and ISTA-NAS [26] are not only precise (2.54%) but also theoretically sound  
 293 approaches. ISTA-NAS [26] is extremely fast in one-shot NAS while ours are more competitive  
 294 (500× faster) in search efficiency.

295 We compare our method with other zero-shot NAS approaches in Table 1. It demonstrates that the  
 296 TE-NAS [2] which is the first algorithm that reaches 4 GPU hours search cost is experimentally  
 297 awesome. TE-NAS uses the neural tangent kernel to approximate the train so it largely reduces  
 298 the cost of training the neural networks. Compared to TE-NAS, our proposed approach is 500×  
 299 faster and our game-based result (-G) gets a comparable test error rate (2.61% v.s. 2.63%) with a  
 300 smaller number of parameters (3.1M v.s. 3.8M). We also surpass the current state-of-the-art zero-shot  
 301 (training-free) method (NASI [22]) by more than 100× in search efficiency and get fewer errors in  
 302 both settings (2.75% v.s. 2.79%; 2.89% v.s. 2.90%).

303 **4.3 Results on NAS-Bench-201**

304 NAS-Bench-201 is a pure-architecture-aware dataset where the neural architectures are trained in the  
 305 same settings, and the info such as performance, parameters, architecture topologies, and operations

are available. Compared to NAS-Bench-101 [28], NAS-Bench-201 adopts a different search space and gets results on various datasets such as CIFAR-10, CIFAR-100, and ImageNet16-120.

As shown in Table 2, it compares our proposed method with the state-of-the-art methods on NAS-Bench-201. Compared to NASWOT(N=10) [13], NASWOT(N=100) and NASWOT(N=1000) are much more accurate due to enlarged sample amounts. However, it also cause  $10\times$  and  $100\times$  waste of search costs. NASI [22] also enlarges its search cost to get much more precise results with extension of 90s. Our approach gets the same search cost with NASWOT (3s) while being much more precise on CIFAR-10 (90.03% v.s. 89.14%, 93.49% v.s. 92.44), CIFAR-100 (70.18% v.s. 68.50%, 71.18% v.s. 68.62%) and ImageNet 16-120 (43.83% v.s. 41.09%, 44.43% v.s. 41.31). A 9s extension of search cost (**Ours-G**) by neuron games gets even better results than NASWOT and NASI for their extreme results.

Table 2: Comparison with the state-of-the-art methods on NAS-Bench-201.

Algorithm	Search Cost	CIFAR-10		CIFAR-100		ImageNet 16-120	
	GPU seconds	Val (%)	Test (%)	Val (%)	Test (%)	Val (%)	Test (%)
ResNet [8]	-	90.83	93.97	70.42	70.86	44.53	43.63
<b>Optimal</b>	-	<b>91.61</b>	<b>94.37</b>	<b>73.49</b>	<b>73.51</b>	<b>46.77</b>	<b>47.31</b>
RSPS [9]	7587	84.16 ± 1.69	87.66 ± 1.69	45.78 ± 6.33	46.60 ± 6.57	31.09 ± 5.65	30.78 ± 6.12
DARTS(1st) [12]	10890	39.77 ± 0.00	54.30 ± 0.00	15.03 ± 0.00	15.61 ± 0.00	16.43 ± 0.00	16.32 ± 0.00
DARTS(2nd) [12]	29902	39.77 ± 0.00	54.30 ± 0.00	15.03 ± 0.00	15.61 ± 0.00	16.43 ± 0.00	16.32 ± 0.00
NASWOT(N=10) [13]	<b>3</b>	89.14 ± 1.14	92.44 ± 1.13	68.50 ± 2.03	68.62 ± 2.04	41.09 ± 3.97	41.31 ± 4.11
NASWOT(N=100) [13]	30	89.55 ± 0.89	92.81 ± 0.99	69.35 ± 1.70	69.48 ± 1.70	42.81 ± 3.05	43.10 ± 3.16
NASWOT(N=1000) [13]	300	89.69 ± 0.73	92.96 ± 0.81	69.86 ± 1.21	69.98 ± 1.22	43.95 ± 2.05	44.44 ± 2.10
NASI(T) [22]	30	-	93.08 ± 0.24	-	69.51 ± 0.59	-	40.87 ± 0.85
NASI(4T) [22]	120	-	93.55 ± 0.10	-	71.20 ± 0.14	-	44.84 ± 1.41
<b>Ours</b>	<b>3</b>	90.03 ± 0.61	93.49 ± 0.71	70.18 ± 1.38	71.18 ± 1.41	43.83 ± 2.10	44.43 ± 2.11
<b>Ours-G</b>	12	90.12 ± 0.52	93.59 ± 0.67	70.54 ± 1.29	71.50 ± 1.31	45.77 ± 1.20	45.73 ± 1.21

316

#### 4.4 Results on ImageNet with the DARTS search space

As shown in Table 3, we report the searched results on ImageNet. The validation size of the observation data batch is 1024. On ImageNet, the number of classes is 1000 so a large data batch is necessary. Compared to NASI [22], and TE-NAS [2], our search costs are faster when  $p = 1$ . The larger batches for evaluation enlarge the search cost for observational data resulting in a slightly larger search cost when  $p = 0.5$ . **Ours(p=1)** gets a competitive test error rate (25.0%) in the table and NASI-ADA [22] gets similar result (24.8%) but NASI-ADA has a larger search cost (864s v.s. 8s).

Table 3: Comparisons with the state-of-the-art on ImageNet.

Algorithm	Search Cost (GPU seconds)	Test Error (%)	Params (M)
DARTS [12]	$8.64 \times 10^5$	26.7	4.7
DARTS+PT [25]	$2.94 \times 10^5$	25.5	<b>4.6</b>
DrNAS [3]	$3.37 \times 10^5$	<b>24.2</b>	5.2
TE-NAS [2]	4320	26.2	5.0
TE-NAS [2]	14688	24.5	5.4
NASI-ADA [22]	864	24.8	5.2
NASI-FIX [22]	864	24.3	5.5
<b>Ours(p=0.5)</b>	1020	25.5	4.9
<b>Ours(p=1)</b>	<b>8</b>	25.0	5.2
<b>Ours-G</b>	31	24.8	5.4

323

## 5 Conclusion

In this work, we interpret the zero-shot NAS as a causal representation learning and solve it by interventional data from one-shot NAS. Besides, our work is dedicated to displaying the inheriting relationship among the latent variables. We demonstrate that the neural architectures can be evaluated and selected by a Gaussian distribution given Gaussian inputs. Experiments on benchmark datasets reveal awesome efficiency and competitive accuracy.

## References

- 330
- 331 [1] Kartik Ahuja, Divyat Mahajan, Yixin Wang, and Yoshua Bengio. Interventional causal repre-  
332 sentation learning. *arXiv preprint arXiv:2209.11924*, 2022.
- 333 [2] Wuyang Chen, Xinyu Gong, and Zhangyang Wang. Neural architecture search on imagenet in  
334 four gpu hours: A theoretically inspired perspective. *arXiv preprint arXiv:2102.11535*, 2021.
- 335 [3] Xiangning Chen, Ruochen Wang, Minhao Cheng, Xiaocheng Tang, and Cho-Jui Hsieh. Drnas:  
336 Dirichlet neural architecture search. *arXiv preprint arXiv:2006.10355*, 2020.
- 337 [4] Frederick Eberhardt and Richard Scheines. Interventions and causal inference. *Philosophy of*  
338 *Science*, 74(5):981–995, 2007.
- 339 [5] Adam Gaier and David Ha. Weight agnostic neural networks. *Advances in neural information*  
340 *processing systems*, 32, 2019.
- 341 [6] Huang Gao, Liu Zhuang, LVD Maaten, and Kilian Q Weinberger. Densely connected convolu-  
342 tional networks. In *CVPR*, volume 1, page 3, 2017.
- 343 [7] Amirata Ghorbani and James Y Zou. Neuron shapley: Discovering the responsible neurons.  
344 *Advances in Neural Information Processing Systems*, 33:5922–5932, 2020.
- 345 [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
346 recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR*  
347 *2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778, 2016.
- 348 [9] Liam Li and Ameet Talwalkar. Random search and reproducibility for neural architecture search.  
349 In *Uncertainty in Artificial Intelligence*, pages 367–377. PMLR, 2020.
- 350 [10] Ming Lin, Pichao Wang, Zhenhong Sun, Heseng Chen, Xiuyu Sun, Qi Qian, Hao Li, and Rong  
351 Jin. Zen-nas: A zero-shot nas for high-performance image recognition. In *2021 IEEE/CVF*  
352 *International Conference on Computer Vision (ICCV)*, pages 337–346, 2021.
- 353 [11] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei,  
354 Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In  
355 *Proceedings of the European conference on computer vision (ECCV)*, pages 19–34, 2018.
- 356 [12] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search.  
357 *arXiv preprint arXiv:1806.09055*, 2018.
- 358 [13] Joe Mellor, Jack Turner, Amos Storkey, and Elliot J Crowley. Neural architecture search  
359 without training. In *Proceedings of the International Conference on Machine Learning*, pages  
360 7588–7598. PMLR, 2021.
- 361 [14] Guido F Montufar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio. On the number of  
362 linear regions of deep neural networks. *Advances in neural information processing systems*, 27,  
363 2014.
- 364 [15] Yookoon Park, Sangho Lee, Gunhee Kim, and David M. Blei. Unsupervised representation  
365 learning via neural activation coding. *arXiv preprint arXiv:2112.04014*, 2021.
- 366 [16] Hieu Pham, Melody Guan, Barret Zoph, Quoc Le, and Jeff Dean. Efficient neural architecture  
367 search via parameters sharing. In *International Conference on Machine Learning*, pages  
368 4095–4104. PMLR, 2018.
- 369 [17] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image  
370 classifier architecture search. In *Proceedings of the aaai conference on artificial intelligence*,  
371 volume 33, pages 4780–4789, 2019.
- 372 [18] Hans Reichenbach. *The Direction of Time*. Dover Publications, 1956.
- 373 [19] Benedek Rozemberczki and Rik Sarkar. The shapley value of classifiers in ensemble games.  
374 *arXiv preprint arXiv:2101.02153*, 2021.

- 375 [20] Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner,  
376 Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the*  
377 *IEEE*, 109(5):612–634, 2021.
- 378 [21] LS Shapley. Quota solutions op n-person games1. *Edited by Emil Artin and Marston Morse*,  
379 page 343, 1953.
- 380 [22] Yao Shu, Shaofeng Cai, Zhongxiang Dai, Beng Chin Ooi, and Bryan Kian Hsiang Low.  
381 Nasi: Label-and data-agnostic neural architecture search at initialization. *arXiv preprint*  
382 *arXiv:2109.00817*, 2021.
- 383 [23] Jin Tian and Judea Pearl. Causal discovery from changes. *arXiv preprint arXiv:1301.2312*,  
384 2013.
- 385 [24] Isabella Verdinelli and Larry Wasserman. Feature importance: A closer look at shapley values  
386 and loco. *arXiv preprint arXiv:2303.05981*, 2023.
- 387 [25] Ruochen Wang, Minhao Cheng, Xiangning Chen, Xiaocheng Tang, and Cho-Jui Hsieh. Re-  
388 thinking architecture selection in differ-entiable nas. In *International Conference on Learning*  
389 *Representations*, 2021.
- 390 [26] Yibo Yang, Hongyang Li, Shan You, Fei Wang, Chen Qian, and Zhouchen Lin. Ista-  
391 nas: Efficient and consistent neural architecture search by sparse coding. *arXiv preprint*  
392 *arXiv:2010.06176*, 2020.
- 393 [27] Quanming Yao, Mengshuo Wang, Hugo Jair Escalante, Isabelle Guyon, Yi-Qi Hu, Yu-Feng Li,  
394 Wei-Wei Tu, Qiang Yang, and Yang Yu. Taking human out of learning applications: A survey  
395 on automated machine learning. *CoRR*, abs/1810.13306, 2018.
- 396 [28] Chris Ying, Aaron Klein, Eric Christiansen, Esteban Real, Kevin Murphy, and Frank Hutter.  
397 Nas-bench-101: Towards reproducible neural architecture search. In *International Conference*  
398 *on Machine Learning*, pages 7105–7114. PMLR, 2019.
- 399 [29] Arber Zela, Thomas Elsken, Tonmoy Saikia, Yassine Marrakchi, Thomas Brox, and Frank  
400 Hutter. Understanding and robustifying differentiable architecture search. In *8th International*  
401 *Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30,*  
402 *2020*, 2020.
- 403 [30] Hongpeng Zhou, Minghao Yang, Jun Wang, and Wei Pan. Bayesnas: A bayesian approach for  
404 neural architecture search. In *International Conference on Machine Learning*, pages 7603–7613.  
405 PMLR, 2019.
- 406 [31] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv*  
407 *preprint arXiv:1611.01578*, 2016.
- 408 [32] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable  
409 architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer*  
410 *vision and pattern recognition*, pages 8697–8710, 2018.