

ZiCo: ZERO-SHOT NAS VIA INVERSE COEFFICIENT OF VARIATION ON GRADIENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural Architecture Search (NAS) is widely used to automatically design the neural network with the best performance among a large number of candidate architectures. To reduce the search time, zero-shot NAS aims at designing *training-free* proxies that can predict the test performance of a given architecture. However, as shown recently, none of the zero-shot proxies proposed to date can actually work consistently better than a naive proxy, namely, the number of network parameters (#Params). To improve this state of affairs, as the main theoretical contribution, we first reveal how some specific gradient properties across different samples impact the convergence rate of neural networks. Based on this theoretical analysis, we propose a new zero-shot proxy, *ZiCo*, the first proxy that works consistently better than #Params. We demonstrate that *ZiCo* works better than State-Of-The-Art (SOTA) proxies on several popular NAS-Benchmarks (NASBench101, NATSBench-SSS/TSS) for multiple datasets (CIFAR10/100, ImageNet16-120). Finally, we demonstrate that the optimal architectures found via *ZiCo* are as competitive as the ones found by one-shot and multi-shot NAS methods, but with much less search time. For example, *ZiCo*-based NAS can find optimal architectures with 78.1%, 79.4%, and 80.4% test accuracy under inference budgets of 450M, 600M, and 1000M FLOPs on ImageNet within 0.4 GPU days.

1 INTRODUCTION

During the last decade, deep learning has achieved great success in many areas, such as computer vision and natural language modeling Krizhevsky et al. (2012); Liu & Deng (2015); Huang et al. (2017); He et al. (2016); Dosovitskiy et al. (2021); Brown et al. (2020); Vaswani et al. (2017). In recent years, neural architecture search (NAS) has been proposed to search for optimal architectures, while reducing the trial-and-error (manual) network design efforts Baker et al. (2017); Zoph & Le (2017); Elsken et al. (2019). Moreover, the neural architectures found via NAS show better performance than the manually-designed networks in many mainstream applications Real et al. (2017); Gong et al. (2019); Xie et al. (2019); Wu et al. (2019); Wan et al. (2020); Li & Talwalkar (2020); Kandasamy et al. (2018); Yu et al. (2020b); Liu et al. (2018b); Cai et al. (2018); Zhang et al. (2019a); Zhou et al. (2019); Howard et al. (2019).

Despite these advantages, most of the existing NAS approaches involve a time-consuming and resource-intensive search process. For example, multi-shot NAS uses a controller or an accuracy predictor to conduct the search process and it requires training of multiple networks; thus, multi-shot NAS is extremely time-consuming (typically thousands of GPU hours) Real et al. (2019); Chiang et al. (2019). Alternatively, one-shot NAS merges all possible networks from the search space into a supernet and thus only needs to train the supernet once Dong & Yang (2019); Zela et al. (2020); Chen et al. (2019); Cai et al. (2019); Stamoulis et al. (2019); Chu et al. (2021); Guo et al. (2020); this enables one-shot NAS to find a good architecture with much less search time. Though the one-shot NAS has significantly improved the time efficiency of NAS, training is still required during the search process.

In the last few years, the *zero-shot* approaches have been proposed to liberate NAS from training entirely Wu et al. (2021); Zhou et al. (2022, 2020); Ingolfsson et al. (2022); Tran & Bae (2021); Do & Luong (2021); Tran et al. (2021); Shu et al. (2022). Essentially, zero-shot NAS utilizes some proxy that can predict the test performance of a given network *without training*. Moreover, the design of the proxy in zero-shot NAS is usually based on some theoretical analysis of deep networks. Hence,

zero-shot approaches can not only significantly improve the time efficiency of NAS, but also deepen the theoretical understanding on why certain networks work well. Nonetheless, as revealed in Ning et al. (2021); White et al. (2022), the zero-shot proxies proposed to date cannot work consistently better than a naive proxy, namely, the number of parameters (#Params); in fact, #Params often achieves the best performance on most of popular NAS benchmarks. These results may undermine the effectiveness of zero-shot NAS approaches.

To address the limitations of existing zero-shot proxies, we target the following **key questions**:

1. How do some specific gradient properties, i.e., mean value and standard deviation across different samples, impact the training convergence of neural networks?
2. Can we use these two gradient properties to design a new theoretically-grounded proxy that works better than #Params consistently?

To this end, we first theoretically analyze how the mean value and standard deviation of gradients across different training batches impact the training convergence of neural networks. Based on our analysis, we propose *ZiCo*, a new proxy for zero-shot NAS. We demonstrate that, compared to all existing proxies (including #Params), ZiCo has either a higher or at least on-par correlation with the test accuracy on popular NAS-Benchmarks (NASBench101, NATS-Bench-SSS/TSS) for multiple datasets (CIFAR10/100, ImageNet16-120). Finally, we demonstrate that ZiCo enables a zero-shot NAS framework that can efficiently find the network architectures with highest test accuracy compared to other zero-shot baselines. In fact, our ZiCo-based zero-shot NAS framework achieves competitive FLOPs-accuracy tradeoffs compared to multiple one-shot and multi-shot NAS, but with much lower time costs. To summarize, we make the following **major contributions**:

- We theoretically reveal how the mean value and standard deviation of gradients across multiple samples impact the training convergence of neural networks.
- We propose a new zero-shot proxy (ZiCo), the first zero-shot proxy that consistently works better than #Params on popular NAS-Benchmarks (NASBench101, NATS-Bench-SSS/TSS) for multiple datasets (CIFAR10/100, ImageNet16-120).
- We demonstrate that our proposed zero-shot NAS achieves competitive test accuracy with representative one-shot and multi-shot NAS with much less search time.

The rest of the paper is organized as follows. We discuss related work in Section 2. In Section 3, we introduce our theoretical analysis. We introduce our proposed zero-shot proxy (ZiCo) and the NAS framework in Section 4. Section 5 validates our analysis and presents our results with the proposed zero-shot NAS. We conclude the paper in Section 6 with remarks on our main contribution.

2 RELATED WORK

2.1 ZERO-SHOT NAS

The goal of zero-shot NAS is to rank the accuracy of various candidate network architectures *without training*, such that we can replace the expensive training process in NAS with some computation-efficient proxies Xiang et al. (2021); Javaheripi et al. (2022); Bhardwaj et al. (2021); Li et al. (2021). Hence, the quality of the proxy determines the effectiveness of zero-shot NAS. Several works use the number of linear regions to approximately measure the expressivity of a deep neural network Mellor et al. (2021); Chen et al. (2021b); Bhardwaj et al. (2022). Alternatively, most of the existing proxies are derived from the gradient of deep networks. For example, Synflow, SNIP, and GraSP rely on the gradient w.r.t the parameters of neural networks; they are proved to be the different approximations of Taylor expansion of deep neural networks Abdelfattah et al. (2021); Lee et al. (2019b); Tanaka et al. (2020); Wang et al. (2020). Moreover, the Zen-score approximates the gradient w.r.t featuremaps and measures the complexity of neural networks Lin et al. (2021); Sun et al. (2021). Furthermore, Jacob_cov leverages the Jacobian matrix between the loss and multiple input samples to quantify the capacity of modeling the complex functions Lopes et al. (2021). Though zero-shot NAS can significantly accelerate the NAS process, it has been revealed that the naive proxy #Params generally works better than all the proxies proposed to date Ning et al. (2021); White et al. (2022). These limitations of existing proxies motivate us to look for a new proxy that can consistently work better than #Params and address the limitations of existing zero-shot NAS approaches.

2.2 KERNEL METHODS AND CONVERGENCE ANALYSIS

Kernel methods are widely explored to analyze the convergence property of networks trained with gradient descent Neal (1996); Williams (1996); Du et al. (2019a); Lu et al. (2020); Allen-Zhu et al.

(2019); Hanin & Nica (2020); Golikov et al. (2022). For example, the training of wide neural networks is proved to be equivalent to the optimization of a specific kernel function Arora et al. (2019a); Lee et al. (2019a); Chizat et al. (2019); Arora et al. (2019b); Cho & Saul (2009). Moreover, given the networks with specific width constraints, researchers proved that the training convergence of networks can be described by some corresponding kernels and the convergence rates of training are highly coupled with the eigenvalues of the kernel-based covariance matrix Mei et al. (2019); Zhang et al. (2019b); Garriga-Alonso et al. (2019); Du et al. (2019b). In our work, we extend such kernel-based analysis to reveal the relationships between the gradient properties and the training convergence for neural networks.

3 TRAINING CONVERGENCE VIA GRADIENT ANALYSIS

We consider the mean value and standard deviation of gradients across different samples and first explore how these two metrics impact the training convergence of linear regression tasks.

3.1 LINEAR REGRESSION

Inspired by Du et al. (2019b), we use the training set \mathbb{S} with M samples as follows:

$$\mathbb{S} = \{(\mathbf{x}_i, y_i), i = 1, \dots, M, \mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}, \|\mathbf{x}_i\| = 1, |y_i| \leq R, M > 1\} \quad (1)$$

where R is a positive constant and $\|\cdot\|$ denotes the $L2$ -norm of a given vector; $\mathbf{x}_i \in \mathbb{R}^d$ is the i^{th} input sample and normalized by its $L2$ -norm (i.e., $\|\mathbf{x}_i\| = 1$), and y_i is the corresponding label. We define the following linear model $f = \mathbf{a}^T \mathbf{x}$ optimized with an MSE-based loss function L :

$$\min_{\mathbf{a}} \sum_i L(y_i, f(\mathbf{x}_i; \mathbf{a})) = \min_{\mathbf{a}} \sum_i \frac{1}{2} (\mathbf{a}^T \mathbf{x}_i - y_i)^2 \quad (2)$$

where $\mathbf{a} \in \mathbb{R}^d$ is the initial weight vector of f . We denote the gradient of L w.r.t to \mathbf{a} as $g(\mathbf{x}_i)$ when taking (\mathbf{x}_i, y_i) as the training sample:

$$g(\mathbf{x}_i) = \frac{\partial L(y_i, f(\mathbf{x}_i; \mathbf{a}))}{\partial \mathbf{a}} \quad (3)$$

We denote the j^{th} element of $g(\mathbf{x}_i)$ as $g_j(\mathbf{x}_i)$. We compute the mean value (μ_j) and standard deviation (σ_j) of $g_j(\mathbf{x}_i)$ across all training samples as follows:

$$\mu_j = \frac{1}{M} \sum_i g_j(\mathbf{x}_i) \quad \sigma_j = \sqrt{\frac{1}{M} \sum_i (g_j(\mathbf{x}_i) - \mu_j)^2} \quad (4)$$

Theorem 3.1. We denote the updated weight vector as $\hat{\mathbf{a}}$ and denote $\sum_{ij} [g_j(\mathbf{x}_i)]^2 = G$. Assume we use the accumulated gradient of all training samples and learning rate η to update the initial weight vector \mathbf{a} , i.e., $\hat{\mathbf{a}} = \mathbf{a} - \eta \sum_i g(\mathbf{x}_i)$. If the learning rate $0 < \eta < 2$, then the total training loss is bounded as follows:

$$\sum_i L(y_i, f(\mathbf{x}_i; \hat{\mathbf{a}})) \leq \frac{G}{2} - \frac{\eta}{2} M^2 (2 - \eta) \sum_j \mu_j^2 \quad (5)$$

In particular, if the learning rate $\eta = \frac{1}{M}$, then $L(\hat{\mathbf{a}})$ is bounded by:

$$\sum_i L(y_i, f(\mathbf{x}_i; \hat{\mathbf{a}})) \leq \frac{M}{2} \sum_j \sigma_j^2 \quad (6)$$

We provide the proof in Appendix A and the experimental results to validate this theorem in Sec 5.2.

Remark 3.1 Intuitively, Theorem. 3.1 tells us that the higher the gradient absolute mean across different training samples, the lower the training loss the model converges to; i.e., the network converges at a faster rate. Similarly, if $\eta M < 1$, the smaller the gradient standard deviation across different training samples/batches, the lower the training loss the model can achieve.

3.2 MLPs WITH RELU

In this section, we generalize the linear model to a network with ReLU activation functions. We primarily consider the standard deviation of gradients in the Gaussian kernel space. We still focus on the regression task on the training set \mathbb{S} defined in Eq. 1. We consider a neural network in the same form as Du et al. (2019b):

$$h(\mathbf{x}; \mathbf{s}, \mathbf{W}) = \frac{1}{\sqrt{m}} \sum_i^m s_r \text{ReLU}(\mathbf{w}_r^T \mathbf{x}) \quad (7)$$

where m is the number of output neurons of the first layer; s_r is the r^{th} element in the output weight vector \mathbf{s} ; $\mathbf{W} \in \mathbb{R}^{m \times d}$ is the input weight matrix, and $\mathbf{w}_r \in \mathbb{R}^d$ is the r^{th} row weight vector in \mathbf{W} .

For training on the dataset \mathbb{S} with M samples defined in Eq. 1, we minimize the following loss function:

$$L(\mathbf{s}, \mathbf{W}) = \sum_{i=1}^M \frac{1}{2} (h(\mathbf{x}_i; \mathbf{s}, \mathbf{W}) - y_i)^2 \quad (8)$$

Following the common practice Du et al. (2019b), we fix the second layer (\mathbf{s}) and use gradient descent to optimize the first layer (\mathbf{W}) with a learning rate η :

$$\mathbf{w}_r(t) = \mathbf{w}_r(t-1) - \eta \sum_{i=0}^t \frac{\partial L(\mathbf{s}, \mathbf{W}(t-1))}{\partial \mathbf{w}_r(t-1)} \quad (9)$$

where $\mathbf{W}(t-1)$ denote the input weight matrix after $t-1$ training steps; $\mathbf{w}_r(t)$ denote the r^{th} row weight vector after t training steps.

Definition 1. (Gram Matrix) A Gram Matrix $H(t) \in \mathbb{R}^{M \times M}$ on the training set $\{(\mathbf{x}_i, y_i), i = 1, \dots, M\}$ after t training steps is defined as follows:

$$H_{ij}(t) = \frac{1}{m} \mathbf{x}_i^T \mathbf{x}_j \sum_{r=1}^m \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(t) \geq 0, \mathbf{x}_j^T \mathbf{w}_r(t) \geq 0\} \quad (10)$$

where \mathbb{I} is the indicator function and $\mathbb{I}\{\mathcal{A}\} = 1$ if and only if event \mathcal{A} happens. We denote the $\lambda_{\min}(H)$ as the minimal eigenvalue of a given matrix H . We denote the $\lambda_0 = \lambda_{\min}(H(\infty))$.

Theorem 3.2. Given a neural network with ReLU activation function optimized by minimizing Eq. 8, we assume that each initial weight vector $\{\mathbf{w}_r(0), r = 1, \dots, n\}$ is i.i.d. generated from $\mathcal{N}(0, I)$ and the gradient for each weight follows i.i.d. $\mathcal{N}(0, \sigma)$, where the σ is measured across different training steps. For some positive constants δ and ϵ , if the learning rate η satisfies $\eta < \frac{\lambda_0 \sqrt{\pi} \delta}{2M^2 \sqrt{2\Phi(1-\epsilon)} t \sigma}$, then with probability at least $(1-\delta)(1-\epsilon)$, the following holds true: for any $r \in [m]$, $\|\mathbf{w}_r(0) - \mathbf{w}_r(t)\| \leq C = \eta t \sigma \sqrt{\Phi(1-\epsilon)}$, and at training step t the Gram matrix $H(t)$ satisfies:

$$\lambda_{\min}(H(t)) \geq \lambda_{\min}(H(0)) - \frac{2\sqrt{2}M^2 \eta t \sigma \sqrt{\Phi(1-\epsilon)}}{\sqrt{\pi} \delta} > 0 \quad (11)$$

where $\Phi(\cdot)$ is the inverse cumulative distribution function for a d -degree chi-distribution $\chi^2(d)$.

We provide the proof in Appendix B. We now introduce the following conclusion from Du et al. (2019b) to further help our analysis.

Lemma 1. Du et al. (2019b) Assume we set the number of output neurons of the first layer $m = \Omega(\frac{M^6}{\lambda_0^4 \delta^3})$ and we i.i.d. initialize $\mathbf{w}_r \sim \mathcal{N}(0, I)$ and $s_r \sim \text{uniform}[-1, 1]$, for $r \in [m]$. When minimizing the loss function in Eq. 8 on the training set \mathbb{S} in Eq. 1, with probability at least $1 - \delta$ over the initialization, the training loss after t training steps is bounded by:

$$L(\mathbf{s}, \mathbf{W}(t)) \leq e^{-\lambda_{\min}(H(t))} L(\mathbf{s}, \mathbf{W}(t-1)) \quad (12)$$

Theorem 3.3. Under the assumptions of Theorem 3.2 and Lemma 1, with probability at least $(1-\delta)(1-\epsilon)$, the following inequality holds true:

$$L(\mathbf{s}, \mathbf{W}(t)) \leq e^{-\lambda_{\min}(H(0))} e^{\frac{2\sqrt{2}M^2 \eta t \sigma \sqrt{\Phi(1-\epsilon)}}{\sqrt{\pi} \delta}} L(\mathbf{s}, \mathbf{W}(t-1)) \quad (13)$$

The proof consists of replacing $\lambda_{\min}(H(t))$ in Eq. 12 with its lower bound given by Theorem 3.2.

Remark 3.2 Theorem. 3.3 shows that after some training steps t , the network with a smaller standard deviation (σ) of gradients will have a smaller training loss; i.e., the network has a faster convergence rate at each training step. We further validate this theorem in Sec. 5.2.

3.3 SUMMARY OF OUR THEORETICAL ANALYSIS

Theorem 3.1 and Theorem 3.3 tell us that the network with high training convergence speed should have *high absolute mean values* and *low standard deviation values* for the gradient, w.r.t the parameters across different training samples/batches. Inspired by these theoretical insights, we next propose a proxy that jointly considers both absolute mean and standard deviation values.

4 NEW ZERO-SHOT PROXY AND NAS FRAMEWORK

In this section, we first define our proxy (ZiCo) and then introduce our zero-shot NAS framework. Following the standard practice, we consider convolutional neural networks (CNNs) as candidate networks.

4.1 PROPOSED ZERO-SHOT PROXY: ZiCo

Definition 2. Given a neural network with D layers and loss function L , the Zero-shot inverse Coefficient of Variation (ZiCo) is defined as follows:

$$\text{ZiCo} = \sum_{l=1}^D \log\left(\sum_{\omega \in \theta_l} \frac{|\mathbb{E}[\nabla_{\omega} L(\mathbf{X}_i, \mathbf{y}_i; \Theta)]|}{\sqrt{\text{Var}(\nabla_{\omega} L(\mathbf{X}_i, \mathbf{y}_i; \Theta))}}\right), \quad i \in \{1, \dots, N\} \quad (14)$$

where Θ denote the initial parameters of the given network; θ_l denote the parameters of the l^{th} layer of the network, and ω represents each element in θ_l ; \mathbf{X}_i and \mathbf{y}_i are the i^{th} input batch and corresponding labels from the training set; N is number of training batches used to compute ZiCo. We incorporate \log to stabilize the computation by regularizing the extremely large or small values.

Of note, our metric is applicable to general CNNs; i.e., there’s no restriction w.r.t. the neural architecture when calculating ZiCo. As discussed in Section 3.3, the networks with higher ZiCo tend to have better convergence rates. Hence, the architectures with higher ZiCo are better architectures.

We remark that the loss values in Eq. 14 are all computed with the **initial** parameters Θ ; that is, we **never** update the value of the parameters when computing ZiCo for a given network (hence it follows the basic principle of zero-shot NAS, i.e., never train, and only use the initial parameters). In practice, two batches are enough to make ZiCo achieve the SOTA performance among all previously proposed accuracy proxies (see Sec. 5.5). Hence, we use only two input batches ($N = 2$) to compute ZiCo; this makes ZiCo highly time efficient for a given network.

4.2 ZiCo-BASED ZERO-SHOT NAS

We use an Evolutionary Algorithm (EA) to conduct the zero-shot NAS because it is concise and easy to implement¹. As shown in Algorithm 1, we search for the neural architectures with the highest ZiCo within the search space, given a specific budget B (e.g., FLOPs). We repeat the search T times; at each search step, we randomly select a structure from the candidate set \mathbb{F} and mutate its architectures (e.g., kernel size, block type, number of blocks, and layer width) to generate a new network $F_i \in \mathcal{S}$. If the generated network F_i meets the inference budget B , we calculate its ZiCo on \mathbb{Z} and add F_i to the candidate set \mathbb{F} . We remove the network with the smallest ZiCo from \mathbb{F} , if the number of architectures in \mathbb{F} exceeds the threshold E . After T steps, we select the network with the largest ZiCo as the final (optimal) architecture F_P .

Algorithm 1 ZiCo-based zero-shot NAS framework

INPUT: Number of search steps T
 Inference budget B , Search space \mathcal{S}
 Set of input batch $\mathbb{Z} = \{(\mathbf{X}_i, \mathbf{y}_i), i = 1, 2\}$
 Population size E , Initial network $F_0 \in \mathcal{S}$

OUTPUT: Optimal network F_P

SEARCH:
 Initialize $\mathbb{F} = \{F_0\}$
for $i = 1$ **to** T **do**
 Randomly sample network F_i from \mathbb{F}
 F_i = randomly mutated architecture based on F_i from \mathcal{S}
 if F_i meets the inference budget B **then**
 Compute ZiCo for F_i on \mathbb{Z} by Eq. 14
 Add F_i to \mathbb{F}
 if $|\mathbb{F}| > E$ **then**
 Remove network with the smallest ZiCo from \mathbb{F}
 end if
 end if
end for
 F_P = the network of the highest ZiCo in \mathbb{F} .

¹One can also use other methods like genetic algorithms or Reinforcement Learning to perform the search.

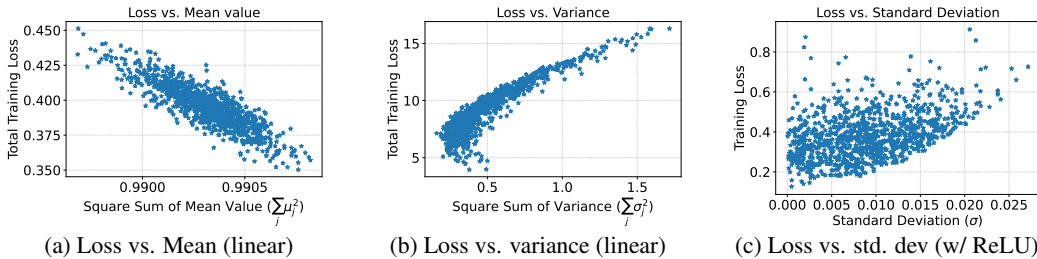


Figure 1: (a). Training loss vs. square sum of mean gradients ($\sum_j \mu_j^2$ in Eq. 5) for randomly sampled 1000 linear networks on MNIST after one epoch. Clearly, smaller mean gradient values lead to lower loss values. (b). Training loss vs. sum of gradients variances ($\sum_j \sigma_j^2$ in Eq. 6) for randomly sampled 1000 linear networks on MNIST after one training epoch. Networks with higher $\sum_j \sigma_j^2$ have lower loss values. (c). Training loss vs. standard deviation of gradients (σ in Eq. 13) for randomly sampled 1000 two-layer MLPs with ReLU on MNIST after one training epoch. Networks with higher σ tend to have lower loss values.

5 EXPERIMENTAL RESULTS

5.1 EXPERIMENTAL SETUP

We conduct the following types of experiments: (i) Empirical validation of Theorem 3.1 and Theorem 3.3, (ii) Evaluation of the proposed ZiCo on multiple NAS benchmarks, (iii) Illustration of ZiCo-based zero-shot NAS on ImageNet.

For the experiments (i), to validate Theorem 3.1, we optimize a linear model as in Eq. 2 on the MNIST dataset, the mean gradient values and the standard deviation vs. the total training loss. Moreover, we also optimize the model defined by Eq. 7 on MNIST and report the training loss vs. the standard deviation in order to validate Theorem 3.2.

For experiments (ii), we compare our proposed ZiCo against existing proxies on two mainstream NAS benchmarks: *NATSBench* is a popular cell-based search space with two different search spaces: (1) NATSBench-TSS consisting of 15625 total architectures with different cell structures trained on CIFAR10, CIFAR100, and ImageNet16-120 (Img16-120) datasets, which is just renamed from NASBench-201 Dong & Yang (2020); (2) NATSBench-SSS contains includes 32768 architectures (which differ only in the width values of each layer) and is also trained on the same three above datasets Dong et al. (2021). *NASBench101* provides users with 423k neural architectures with their test accuracy on CIFAR10 dataset; the architectures are built by stacking the same cell multiple times Ying et al. (2019).

For experiments (iii), we use ZiCo to conduct the zero-shot NAS (see Algorithm 1) on ImageNet. We first use Algorithm 1 to find the networks with the highest ZiCo under various FLOPs budgets. We conduct the search for 100k steps; this takes 10 hours on a single NVIDIA 3090 GPU (i.e., 0.4 GPU days). Then, we train the obtained network with the exact same training setup as Lin et al. (2021). Specifically, we train the neural network for 480 epochs with the batch size 512 and input resolution 224. We also use the distillation-based training loss functions by taking Efficient-B3 as the teacher. Finally, we set the initial learning rate as 0.1 with a cosine annealing scheduling scheme.

5.2 VALIDATION OF THEOREM 3.1 AND THEOREM 3.3

To empirically validate Theorem 3.1, we first create the training set \mathbb{S} by normalizing randomly sampled 1000 training samples in MNIST and normalizing them with their L_2 -norm. We compute the gradient w.r.t. the network parameters for each individual training sample. Next, as discussed in Theorem 3.1, we use the accumulated gradient over these samples to update the network parameters with learning rate $\eta = 1$. Then, we calculate the square sum of mean gradients and the total training loss. We repeat the above process 1000 times on the same \mathbb{S} . As shown in Fig. 1(a), we plot the total training loss vs. square sum of mean gradients as defined in Eq. 5. Clearly, the networks with the higher square sum of mean gradients values tend to have lower training loss. In comparison, Fig. 1(b) shows that networks with a lower square sum of variance value tend to have lower training loss values, which coincides with the conclusion drawn from Eq. 6. These results empirically validate our Theorem 3.1.

Table 1: The correlation coefficients between various zero-cost proxies and two naive proxies (#Params and FLOPs) vs. test accuracy on NATSBench-SSS and NATSBench-TSS (KT and SPR represent Kendall’s τ and Spearman’s ρ , respectively). The results in italics represent the values of #Params’ correlation coefficients. The results better than #Params are shown with bold fonts. Clearly, our proposed ZiCo is the only proxy that works consistently better than #Params and is generally the best among all these proxies.

NATSBench-TSS (NASBench201)							
Dataset		CIFAR10		CIFAR100		Img16-120	
Proxy	Correlation	KT	SPR	KT	SPR	KT	SPR
Grad_norm Abdelfattah et al. (2021)		0.46	0.63	0.47	0.63	0.43	0.58
SNIP Lee et al. (2019b)		0.46	0.63	0.46	0.63	0.43	0.58
GraSP Wang et al. (2020)		0.37	0.54	0.36	0.51	0.40	0.56
Fisher Liu et al. (2021)		0.40	0.55	0.41	0.55	0.37	0.50
Synflow Tanaka et al. (2020)		0.07	0.11	0.06	0.08	-0.17	-0.24
Zen-score Lin et al. (2021)		0.29	0.38	0.28	0.36	0.29	0.40
FLOPs		0.54	0.73	0.51	0.71	0.49	0.67
#Params		<i>0.57</i>	<i>0.75</i>	<i>0.55</i>	<i>0.73</i>	<i>0.52</i>	<i>0.69</i>
ZiCO		0.61	0.80	0.61	0.81	0.60	0.79

NATSBench-SSS							
Dataset		CIFAR10		CIFAR100		Img16-120	
Proxy	Correlation	KT	SPR	KT	SPR	KT	SPR
Grad_norm Abdelfattah et al. (2021)		0.35	0.51	0.34	0.49	0.49	0.67
SNIP Lee et al. (2019b)		0.42	0.59	0.46	0.62	0.57	0.76
GraSP Wang et al. (2020)		-0.09	-0.13	0.01	0.01	0.29	0.42
Fisher Liu et al. (2021)		0.30	0.44	0.41	0.55	0.33	0.47
Synflow Tanaka et al. (2020)		0.61	0.81	0.60	0.80	0.39	0.57
Zen-score Lin et al. (2021)		0.50	0.69	0.52	0.71	0.69	0.87
FLOPs		0.19	0.28	0.21	0.30	0.38	0.53
#Params		<i>0.53</i>	<i>0.72</i>	<i>0.54</i>	<i>0.73</i>	<i>0.65</i>	<i>0.84</i>
ZiCO		0.54	0.73	0.55	0.75	0.70	0.88

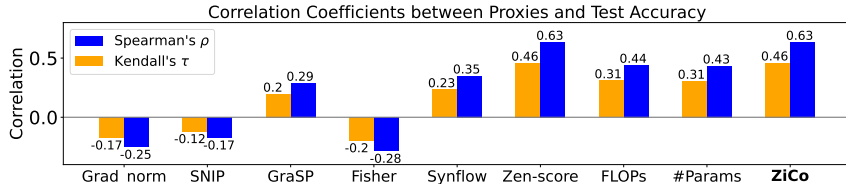


Figure 2: The correlation coefficients between various zero-cost proxies vs. test accuracy on NAS-Bench101 search space for CIFAR10 dataset. As shown, our proposed ZiCo correlates best with the real test accuracy and is significantly better than all other proxies except for Zen-score.

Moreover, to optimize a two-layer MLP with ReLU activation functions as defined in Eq. 7, we use the entire training set of MNIST and apply the gradient descent (Eq. 9) to update the weights. We set the batch size as 256 and measure the standard deviation of gradients (σ) w.r.t. parameters across different training batches. We set a very small learning rate $\eta = 10^{-8}$ to satisfy the assumption in Theorem 3.2. We plot the training loss after one training epoch vs. standard deviation of gradients (σ) in Fig. 1(c). Clearly, the results show that if a network has a lower gradient standard deviation, then it tends to have lower training loss values, and thus, a faster convergence rate. These results empirically prove our claims in Theorem 3.3.

5.3 ZiCO vs. OTHER PROXIES ON NAS BENCHMARKS

We first calculate the correlation coefficients between various proxies and the test accuracy on CIFAR10, CIFAR100, and ImageNet16-120 datasets for NATSBench. As shown in Table. 1, ZiCo achieves the highest correlation with the real test accuracy, except for CIFAR10/100 on NATSBench-SSS. Moreover, we observe that most of the previously proposed proxies work well for some specific scenarios, but do *not* generalize well to other scenarios. For example, though Synflow is slightly better than ZiCo for CIFAR10/100 on NATSBench-SSS, it has poor correlation scores in all other cases (e.g., 0.07 Kendall’s τ score on CIFAR10 on NATSBench-TSS). Similarly, Zen-score performs well for Img16-120 on NATSBench-SSS, but it doesn’t work well on NATSBench-TSS. In contrast, ZiCo has either the highest or second highest correlation coefficients under *all* these scenarios in Table 1. We provide more results in Appendix D.

Table 2: The test accuracy of optimal architectures obtained by various zero-shot proxies (average on 5 runs) on NATSBench-TSS search space. The best results are shown with bold fonts.

Dataset	Groud Truth	Grad_norm	SNIP	GraSP	Fisher	Jacob_cov	Synflow	Zen-score	#Params	FLOPs	ZiCo
CIFAR100	73.5	60.0	60.0	60.0	60.0	68.9	62.2	68.1	71.1	71.1	71.1±0.3
Img16-120	47.3	29.3	29.3	5.5	29.3	25.1	26.1	40.8	41.4	41.4	41.8±0.3
CIFAR10	94.5	89.5	89.5	89.5	89.5	88.4	88.5	90.6	93.7	93.7	94.0±0.4

Table 3: Comparison of Top-1 accuracy of our ZiCo-based NAS against SOTA NAS methods on ImageNet under various FLOP budgets. For the ‘Method’ column, ‘MS’ represents multi-shot NAS; ‘OS’ is short for one-shot NAS; Scaling represents network scaling methods; ‘ZS’ is short for zero-shot NAS. OFA[‡] is trained from scratch and reported in Moons et al. (2021).

Budget (maximal #FLOPs)	Approach	FLOPs	Top-1 [%]	Method	Costs [GPU Days]
450M	EfficientNet-B0 Tan & Le (2019)	390M	77.1	Scaling	3800
	MnasNet-A3 Tan et al. (2019)	403M	76.7	MS	-
	OFA [‡] Cai et al. (2020)	406M	77.7	OS	50
	BN-NAS Chen et al. (2021a)	470M	75.7	MS	0.8
	RLNAS Zhang et al. (2021)	473M	75.6	OS	-
	NASNet-B Zoph et al. (2018)	488M	72.8	MS	1800
	CARS-D Yang et al. (2020)	496M	73.3	MS	0.4
	DONNA Moons et al. (2021)	501M	78.0	OS	405
	#Params	451M	63.5	ZS	0.02
ZiCo (Ours)	448M	78.1	ZS	0.4	
600M	DARTS Liu et al. (2019)	574M	73.3	OS	4
	NAO Luo et al. (2018)	584M	75.5	MS	58.3
	PC-DARTS Xu et al. (2019)	586M	75.8	OS	3.8
	BigNAS-L Yu et al. (2020a)	586M	79.5	OS	2304 (TPU days)
	PNAS Liu et al. (2018a)	588M	74.2	MS	224
	CARS-I Yang et al. (2020)	591M	75.2	MS	0.4
	EnTranNAS Yang et al. (2021)	594M	76.2	OS	2.1
	ProxlessNAS Cai et al. (2019)	595M	76.0	OS	8.3
	RLNAS Zhang et al. (2021)	597M	75.9	OS	-
	MAGIC-AT Xu et al. (2022)	598M	76.8	OS	2
	SemiNAS Luo et al. (2020)	599M	76.5	MS	4
	DONNA Moons et al. (2021)	599M	78.4	OS	405
	Zen-score Lin et al. (2021)	611M	79.1	ZS	0.5
	OFA [‡] Cai et al. (2020)	662M	78.7	OS	50
	EfficientNet-B1 Tan & Le (2019)	700M	79.1	Scaling	3800
ZiCo (Ours)	603M	79.4	ZS	0.4	
1000M	sharpDARTS Hundt et al. (2019)	950M	76.0	OS	-
	EfficientNet-B2 Tan & Le (2019)	1000M	80.1	Scaling	3800
	ZiCo (Ours)	1005M	80.5	ZS	0.4

For NASBench101, as shown in Fig. 2, ZiCo has a significantly higher correlation score with the real test accuracy than all the other proxies, except Zen-score. For example, ZiCo has a 0.46 Kendall’s τ score, while #Params is only 0.31. In general, ZiCo has the highest correlation coefficients among all existing proxies for various search spaces and datasets of NATSBench and NASBench101.

Beside the correlation coefficients, we also report the optimal architectures found with various proxies. As shown in Table 2, the architectures found via ZiCo have the highest test accuracy on all these three datasets. To our best knowledge, ZiCo is the first proxy that shows a consistently higher correlation coefficient compared to #Params.

The above results validate the effectiveness of our proposed ZiCo; thus, ZiCo can be directly used to search for optimal networks for various budgets. Next, we describe the search results in detail.

5.4 ZICO ON IMAGENET

Search Space We use the commonly used MobileNetv2-based search space where the candidate networks are built by stacking multiple Inverted Bottleneck Blocks (IBNs) with SE modules Sandler et al. (2018); Pham et al. (2018); Lin et al. (2021). As for each IBN, the kernel size of the depth-wise convolutional layer is sampled from $\{3,5,7\}$ and the expansion ratio is randomly selected from $\{1,2,4,6\}$. We primarily consider ReLU as the activation function. We use standard Kaiming_Init to initialize all linear and convolution layers for every candidate networks He et al. (2015). More details of the search space are given in Appendix C.

We use Algorithm 1 to search networks under various FLOPs budgets (450M, 600M, and 1000M) within the above search space. As shown in Table 3, ZiCo outperforms most previous NAS ap-

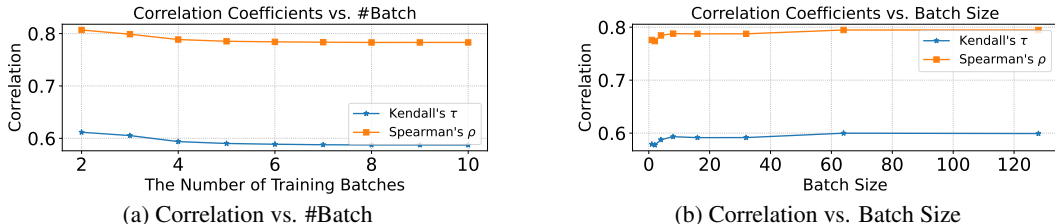


Figure 3: Ablation study. (a) The correlation coefficients between ZiCo computed under varying number of batches and real test accuracy. (b) The correlation coefficients between ZiCo computed with varying batch size and real test accuracy.

proaches by a large margin. For example, when the FLOPs budget is around 450M, ZiCo achieves 78.1% Top-1 accuracy, which is competitive with one of the SOTA NAS methods (DONNA), but with fewer FLOPs and 648 \times faster search speed Moons et al. (2021). Moreover, if the FLOPs is 600M, ZiCo achieves 2.6% higher Top-1 Accuracy than the latest one-shot NAS method (MAGIC-AT) with a 3 \times reduction in terms of search time Xu et al. (2022).

To make further comparison with #Params, we also use #Params as the proxy and Algorithm 1 to conduct the search under a 450M FLOPs budget. As shown in Table 3, the obtained network by #Params has a 14.6% lower accuracy than ours (63.5% vs. 78.1%). Hence, even though the correlations for ZiCo and #Params in Table 1 and the optimal networks in Table 2 are similar for small-scale datasets, ZiCo significantly outperforms naive baselines like #Params for large datasets like ImageMet. To conclude, ZiCo achieves SOTA results for Zero-Shot NAS and outperforms naive methods, existing zero-shot proxies, as well as several one-shot and multi-shot methods.

We remark that these results demonstrate two benefits of our proposed ZiCo: (i) **Lightweight computation costs.** As discussed in Sec 3, during the search process, to evaluate a given architecture, we only need to conduct the backward propagation twice (only takes 0.3s on an NVIDIA 3090 GPU). The computation efficiency and exemption of training enable ZiCo to significantly reduce the search time of NAS. (ii) **High correlation with the real test accuracy.** As demonstrated in Sec 5.3, ZiCo has a very high correlation score with real accuracy for architectures from various search spaces and datasets. Hence, ZiCo can accurately predict the test accuracy of diverse neural architectures, thus helping find the optimal architectures with the best test performance.

5.5 ABLATION STUDY

Number of batches We randomly select 2000 networks from NATSBench-TSS on CIFAR100 dataset and compute ZiCo under varying number of training batches (N in Eq. 14) from $\{2, \dots, 10\}$. We then calculate the correlation coefficients between ZiCo computed under different number of training batches and the real test accuracy. As shown in Fig. 3(a), using two batches to compute ZiCo generates the highest score. Therefore, in our work, we always use two batches ($N = 2$) to compute ZiCo since it is both accurate and time-efficient.

Batch size Similarly, we randomly select 2000 networks from NATSBench-TSS on CIFAR100 and compute ZiCo with two batches under varying batch size $\{1, 2, 4, 8, 16, 32, 64, 128\}$. We then calculate the correlation coefficients between ZiCo computed under various batch sizes and the real test accuracy. As shown in Fig. 3(b), batch size 64 is enough to stabilize the coefficient. Therefore, we set the batch size as 128 and use two batches to compute ZiCo.

6 CONCLUSION

In this work, we have proposed ZiCo, a new SOTA proxy for zero-shot NAS. As the main theoretical contribution, we first reveal how the mean value and standard deviation of gradients impact the training convergence of a given architecture. Based on this theoretical analysis, we have shown that ZiCo works better than all zero-shot NAS proxies proposed so far on multiple popular NAS-Benchmarks (NASBench101, NATSBench-SSS/TSS) for multiple datasets (CIFAR10/100, ImageNet16-120). In particular, we have demonstrated that ZiCo is consistently better than (#Params) and existing zero-shot proxies. Moreover, ZiCo enables us to find architectures with competitive test performance to representative one-shot and multi-shot NAS methods, but with much lower search costs. For example, ZiCo-based NAS can find the architectures with 78.1%, 79.4%, and 80.4% test accuracies under 450M, 600M, and 1000M FLOPs budgets on ImageNet within 0.4 GPU days.

REFERENCES

- Mohamed S. Abdelfattah, Abhinav Mehrotra, Lukasz Dudziak, and Nicholas Donald Lane. Zero-cost proxies for lightweight NAS. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via overparameterization. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 242–252. PMLR, 2019.
- Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 8139–8148, 2019a.
- Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 322–332. PMLR, 2019b.
- Bowen Baker, Otakrist Gupta, Nikhil Naik, and Ramesh Raskar. Designing neural network architectures using reinforcement learning. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- Kartikeya Bhardwaj, Guihong Li, and Radu Marculescu. How does topology influence gradient propagation and model performance of deep networks with densenet-type skip connections? In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 13498–13507. Computer Vision Foundation / IEEE, 2021.
- Kartikeya Bhardwaj, James Ward, Caleb Tung, Dibakar Gope, Lingchuan Meng, Igor Fedorov, Alex Chalfin, Paul N. Whatmough, and Danny Loh. Restructurable activation networks. *CoRR*, abs/2208.08562, 2022.
- Tom B. Brown et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- Han Cai, Tianyao Chen, Weinan Zhang, Yong Yu, and Jun Wang. Efficient architecture search by network transformation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- Han Cai, Ligeng Zhu, and Song Han. ProxylessNAS: Direct neural architecture search on target task and hardware. In *International Conference on Learning Representations*, 2019.
- Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network and specialize it for efficient deployment. In *International Conference on Learning Representations*, 2020.
- Boyu Chen, Peixia Li, Baopu Li, Chen Lin, Chuming Li, Ming Sun, Junjie Yan, and Wanli Ouyang. BN-NAS: neural architecture search with batch normalization. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pp. 307–316. IEEE, 2021a.
- Wuyang Chen, Xinyu Gong, and Zhangyang Wang. Neural architecture search on imagenet in four GPU hours: A theoretically inspired perspective. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021b.
- Xin Chen, Lingxi Xie, Jun Wu, and Qi Tian. Progressive differentiable architecture search: Bridging the depth gap between search and evaluation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1294–1303, 2019.

- Wei-Lin Chiang et al. Cluster-gcn: An Efficient Algorithm for Training Deep and Large Graph Convolutional Networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 257–266, 2019.
- Lénaïc Chizat, Edouard Oyallon, and Francis R. Bach. On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 2933–2943, 2019.
- Youngmin Cho and Lawrence K. Saul. Kernel methods for deep learning. In *Advances in Neural Information Processing Systems 22: 23rd Annual Conference on Neural Information Processing Systems 2009. Proceedings of a meeting held 7-10 December 2009, Vancouver, British Columbia, Canada*, pp. 342–350. Curran Associates, Inc., 2009.
- Xiangxiang Chu, Bo Zhang, and Ruijun Xu. Fairnas: Rethinking evaluation fairness of weight sharing neural architecture search. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 12239–12248, 2021.
- Tu Do and Ngoc Hoang Luong. Training-free multi-objective evolutionary neural architecture search via neural tangent kernel and number of linear regions. In *International Conference on Neural Information Processing*, pp. 335–347. Springer, 2021.
- Xuanyi Dong and Yi Yang. Searching for a robust neural architecture in four gpu hours. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1761–1770, 2019.
- Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- Xuanyi Dong, Lu Liu, Katarzyna Musial, and Bogdan Gabrys. Nats-bench: Benchmarking nas algorithms for architecture topology and size. *IEEE transactions on pattern analysis and machine intelligence*, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- Simon S. Du, Jason D. Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1675–1685. PMLR, 2019a.
- Simon S. Du, Xiyu Zhai, Barnabás Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019b.
- Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *The Journal of Machine Learning Research*, 20(1):1997–2017, 2019.
- Adrià Garriga-Alonso, Carl Edward Rasmussen, and Laurence Aitchison. Deep convolutional networks as shallow gaussian processes. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019.
- Eugene Golikov, Eduard Pokonechnyy, and Vladimir Korviakov. Neural tangent kernel: A survey. *CoRR*, abs/2208.13614, 2022.
- Xinyu Gong, Shiyu Chang, Yifan Jiang, and Zhangyang Wang. Autogan: Neural architecture search for generative adversarial networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3224–3234, 2019.

- Zichao Guo, Xiangyu Zhang, Haoyuan Mu, Wen Heng, Zechun Liu, Yichen Wei, and Jian Sun. Single path one-shot neural architecture search with uniform sampling. In *European conference on computer vision*, pp. 544–560. Springer, 2020.
- Boris Hanin and Mihai Nica. Finite depth and width corrections to the neural tangent kernel. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pp. 1026–1034. IEEE Computer Society, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1314–1324, 2019.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely Connected Convolutional Networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- Andrew Hundt, Varun Jain, and Gregory D. Hager. sharpdarts: Faster and more accurate differentiable architecture search. *CoRR*, abs/1903.09900, 2019.
- Thorir Mar Ingolfsson, Mark Vero, Xiaying Wang, Lorenzo Lamberti, Luca Benini, and Matteo Spallanzani. Reducing neural architecture search spaces with training-free statistics and computational graph clustering. In *CF '22: 19th ACM International Conference on Computing Frontiers, Turin, Italy, May 17 - 22, 2022*, pp. 213–214. ACM, 2022.
- Mojan Javaheripi, Shital Shah, Subhabrata Mukherjee, Tomasz L. Religa, Caio C. T. Mendes, Gustavo H. de Rosa, Sébastien Bubeck, Farinaz Koushanfar, and Debadeepta Dey. Littransformersearch: Training-free on-device search for efficient autoregressive language models. *CoRR*, abs/2203.02094, 2022.
- Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabás Póczos, and Eric P. Xing. Neural architecture search with bayesian optimisation and optimal transport. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 2020–2029, 2018.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, 2012.
- Jaehoon Lee, Lechao Xiao, Samuel S. Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 8570–8581, 2019a.
- Namhoon Lee, Thalaiyasingam Ajanthan, and Philip Torr. SNIP: SINGLE-SHOT NETWORK PRUNING BASED ON CONNECTION SENSITIVITY. In *International Conference on Learning Representations*, 2019b.
- Guihong Li, Sumit K. Mandal, Ümit Y. Ogras, and Radu Marculescu. FLASH: fast neural architecture search with hardware optimization. *ACM Trans. Embed. Comput. Syst.*, 20(5s):63:1–63:26, 2021.
- Liam Li and Ameet Talwalkar. Random search and reproducibility for neural architecture search. In *Uncertainty in artificial intelligence*, pp. 367–377. PMLR, 2020.

- Ming Lin, Pichao Wang, Zhenhong Sun, Heseng Chen, Xiuyu Sun, Qi Qian, Hao Li, and Rong Jin. Zen-nas: A zero-shot nas for high-performance image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 347–356, 2021.
- Chenxi Liu et al. Progressive Neural Architecture Search. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 19–34, 2018a.
- Hanxiao Liu, Karen Simonyan, Oriol Vinyals, Chrisantha Fernando, and Koray Kavukcuoglu. Hierarchical representations for efficient architecture search. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018b.
- Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019.
- Liyang Liu, Shilong Zhang, Zhanghui Kuang, Aojun Zhou, Jing-Hao Xue, Xinjiang Wang, Yimin Chen, Wenming Yang, Qingmin Liao, and Wayne Zhang. Group fisher pruning for practical network compression. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 7021–7032. PMLR, 2021.
- Shuying Liu and Weihong Deng. Very deep convolutional neural network based image classification using small training sample size. In *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*, pp. 730–734, 2015.
- Vasco Lopes, Saeid Alirezazadeh, and Luís A Alexandre. Epe-nas: Efficient performance estimation without training for neural architecture search. In *International Conference on Artificial Neural Networks*, pp. 552–563. Springer, 2021.
- Yiping Lu, Chao Ma, Yulong Lu, Jianfeng Lu, and Lexing Ying. A mean-field analysis of deep resnet and beyond: Towards provable optimization via overparameterization from depth. *CoRR*, abs/2003.05508, 2020.
- Renqian Luo, Fei Tian, Tao Qin, Enhong Chen, and Tie-Yan Liu. Neural architecture optimization. *Advances in neural information processing systems*, 31, 2018.
- Renqian Luo, Xu Tan, Rui Wang, Tao Qin, Enhong Chen, and Tie-Yan Liu. Semi-supervised neural architecture search. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory, COLT 2019, 25-28 June 2019, Phoenix, AZ, USA*, volume 99 of *Proceedings of Machine Learning Research*, pp. 2388–2464. PMLR, 2019.
- Joe Mellor, Jack Turner, Amos Storkey, and Elliot J Crowley. Neural architecture search without training. In *International Conference on Machine Learning*, pp. 7588–7598. PMLR, 2021.
- Bert Moons, Parham Noorzad, Andrii Skliar, Giovanni Mariani, Dushyant Mehta, Chris Lott, and Tijmen Blankevoort. Distilling optimal neural networks: Rapid search in diverse spaces. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pp. 12209–12218. IEEE, 2021.
- Radford M Neal. Priors for infinite networks. In *Bayesian Learning for Neural Networks*, pp. 29–53. Springer, 1996.
- Xuefei Ning, Changcheng Tang, Wenshuo Li, Zixuan Zhou, Shuang Liang, Huazhong Yang, and Yu Wang. Evaluating efficient performance estimators of neural architectures. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 12265–12277, 2021.

- Hieu Pham, Melody Guan, Barret Zoph, Quoc Le, and Jeff Dean. Efficient neural architecture search via parameters sharing. In *International conference on machine learning*, pp. 4095–4104. PMLR, 2018.
- Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *Proceedings of the aaai conference on artificial intelligence*, pp. 4780–4789, 2019.
- Esteban Real et al. Large-scale Evolution of Image Classifiers. In *International Conference on Machine Learning*, pp. 2902–2911. PMLR, 2017.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted Residuals and Linear Bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4510–4520, 2018.
- Yao Shu, Zhongxiang Dai, Zhaoxuan Wu, and Bryan Kian Hsiang Low. Unifying and boosting gradient-based training-free neural architecture search. *CoRR*, abs/2201.09785, 2022.
- Dimitrios Stamoulis, Ruizhou Ding, Di Wang, Dimitrios Lymberopoulos, Bodhi Priyantha, Jie Liu, and Diana Marculescu. Single-path NAS: designing hardware-efficient convnets in less than 4 hours. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2019, Würzburg, Germany, September 16-20, 2019, Proceedings, Part II*, volume 11907 of *Lecture Notes in Computer Science*, pp. 481–497. Springer, 2019.
- Zhenhong Sun, Ming Lin, Xiuyu Sun, Zhiyu Tan, and Rong Jin. Revisiting efficient object detection backbones from zero-shot neural architecture search. *arXiv preprint arXiv:2111.13336*, 2021.
- Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6105–6114. PMLR, 2019.
- Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2820–2828, 2019.
- Hidenori Tanaka, Daniel Kunin, Daniel L Yamins, and Surya Ganguli. Pruning neural networks without any data by iteratively conserving synaptic flow. In *Advances in Neural Information Processing Systems*, volume 33, pp. 6377–6389. Curran Associates, Inc., 2020.
- Linh Tam Tran and Sung-Ho Bae. Training-free hardware-aware neural architecture search with reinforcement learning. *Journal of Broadcast Engineering*, 26(7):855–861, 2021.
- Linh-Tam Tran, Muhammad Salman Ali, and Sung-Ho Bae. A feature fusion based indicator for training-free neural architecture search. *IEEE Access*, 9:133914–133923, 2021.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Alvin Wan, Xiaoliang Dai, Peizhao Zhang, Zijian He, Yuandong Tian, Saining Xie, Bichen Wu, Matthew Yu, Tao Xu, Kan Chen, et al. Fbnetv2: Differentiable neural architecture search for spatial and channel dimensions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12965–12974, 2020.
- Chaoqi Wang, Guodong Zhang, and Roger B. Grosse. Picking winning tickets before training by preserving gradient flow. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- Colin White, Mikhail Khodak, Renbo Tu, Shital Shah, Sébastien Bubeck, and Debadeepta Dey. A deeper look at zero-cost proxies for lightweight nas. In *ICLR Blog Track, 2022*. URL <https://iclr-blog-track.github.io/2022/03/25/zero-cost-proxies/>.

- Christopher K. I. Williams. Computing with infinite networks. In *Advances in Neural Information Processing Systems 9, NIPS, Denver, CO, USA, December 2-5, 1996*, pp. 295–301. MIT Press, 1996.
- Bichen Wu, Xiaoliang Dai, Peizhao Zhang, Yanghan Wang, Fei Sun, Yiming Wu, Yuandong Tian, Peter Vajda, Yangqing Jia, and Kurt Keutzer. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10734–10742, 2019.
- Meng-Ting Wu, Hung-I Lin, and Chun-Wei Tsai. A training-free genetic neural architecture search. In *ACM ICEA '21: 2021 ACM International Conference on Intelligent Computing and its Emerging Applications, Jinan, China, December 28 - 29, 2022*, pp. 65–70. ACM, 2021.
- Lichuan Xiang, Lukasz Dudziak, Mohamed S. Abdelfattah, Thomas Chau, Nicholas D. Lane, and Hongkai Wen. Zero-cost proxies meet differentiable architecture search. *CoRR*, abs/2106.06799, 2021.
- Sirui Xie, Hehui Zheng, Chunxiao Liu, and Liang Lin. SNAS: stochastic neural architecture search. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019.
- Jin Xu, Xu Tan, Kaitao Song, Renqian Luo, Yichong Leng, Tao Qin, Tie-Yan Liu, and Jian Li. Analyzing and mitigating interference in neural architecture search. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 24646–24662. PMLR, 2022.
- Yuhui Xu, Lingxi Xie, Xiaopeng Zhang, Xin Chen, Guo-Jun Qi, Qi Tian, and Hongkai Xiong. Pc-darts: Partial channel connections for memory-efficient architecture search. *arXiv preprint arXiv:1907.05737*, 2019.
- Yibo Yang, Shan You, Hongyang Li, Fei Wang, Chen Qian, and Zhouchen Lin. Towards improving the consistency, efficiency, and flexibility of differentiable neural architecture search. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 6667–6676. Computer Vision Foundation / IEEE, 2021.
- Zhaohui Yang, Yunhe Wang, Xinghao Chen, Boxin Shi, Chao Xu, Chunjing Xu, Qi Tian, and Chang Xu. CARS: continuous evolution for efficient neural architecture search. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pp. 1826–1835. Computer Vision Foundation / IEEE, 2020.
- Chris Ying, Aaron Klein, Eric Christiansen, Esteban Real, Kevin Murphy, and Frank Hutter. Nas-bench-101: Towards reproducible neural architecture search. In *International Conference on Machine Learning*, pp. 7105–7114. PMLR, 2019.
- Jiahui Yu, Pengchong Jin, Hanxiao Liu, Gabriel Bender, Pieter-Jan Kindermans, Mingxing Tan, Thomas Huang, Xiaodan Song, Ruoming Pang, and Quoc Le. Bignas: Scaling up neural architecture search with big single-stage models. In *European Conference on Computer Vision*, pp. 702–717. Springer, 2020a.
- Kaicheng Yu, Christian Sciuto, Martin Jaggi, Claudiu Musat, and Mathieu Salzmann. Evaluating the search phase of neural architecture search. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020b.
- Arber Zela, Thomas Elsken, Tonmoy Saikia, Yassine Marrakchi, Thomas Brox, and Frank Hutter. Understanding and robustifying differentiable architecture search. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*, 2020.
- Chris Zhang, Mengye Ren, and Raquel Urtasun. Graph hypernetworks for neural architecture search. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*, 2019a.

- Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu. Learning one-hidden-layer relu networks via gradient descent. In *The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan*, volume 89 of *Proceedings of Machine Learning Research*, pp. 1524–1534. PMLR, 2019b.
- Xuanyang Zhang, Pengfei Hou, Xiangyu Zhang, and Jian Sun. Neural architecture search with random labels. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 10907–10916. Computer Vision Foundation / IEEE, 2021.
- Dongzhan Zhou, Xinchu Zhou, Wenwei Zhang, Chen Change Loy, Shuai Yi, Xuesen Zhang, and Wanli Ouyang. Econas: Finding proxies for economical neural architecture search. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pp. 11396–11404, 2020.
- Hongpeng Zhou, Minghao Yang, Jun Wang, and Wei Pan. Bayesnas: A bayesian approach for neural architecture search. In *International conference on machine learning*, pp. 7603–7613. PMLR, 2019.
- Qinqin Zhou, Kekai Sheng, Xiawu Zheng, Ke Li, Xing Sun, Yonghong Tian, Jie Chen, and Rongrong Ji. Training-free transformer architecture search. *CoRR*, abs/2203.12217, 2022.
- Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning Transferable Architectures for Scalable Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8697–8710, 2018.

A PROOF OF THEOREM 3.1

Theorem 3.1 We denote the updated weight vector as $\hat{\mathbf{a}}$ and $\sum_{i,j}[g_j(\mathbf{x}_i)]^2 = G$. Assume we use the accumulated gradient of all training samples and learning rate η to update the initial weight vector \mathbf{a} , i.e., $\hat{\mathbf{a}} = \mathbf{a} - \eta \sum_i g(\mathbf{x}_i)$. If the learning rate $0 < \eta < 2$, then the total training loss is bounded as follows:

$$\sum_i L(y_i, f(\mathbf{x}_i; \hat{\mathbf{a}})) \leq \frac{G}{2} - \frac{\eta}{2} M^2 (2 - \eta) \sum_j \mu_j^2 \quad (15)$$

In particular, if the learning rate $\eta = \frac{1}{M}$, then $L(\hat{\mathbf{a}})$ is bounded by:

$$\sum_i L(y_i, f(\mathbf{x}_i; \hat{\mathbf{a}})) \leq \frac{M}{2} \sum_j \sigma_j^2 \quad (16)$$

Proof. Given each training sample (\mathbf{x}_i, y_i) the gradient of L w.r.t to \mathbf{a} when taking (\mathbf{x}_i, y_i) as the input is as follows:

$$g(\mathbf{x}_i) = \frac{\partial L(y_i, f(\mathbf{x}_i; \mathbf{a}))}{\partial \mathbf{a}} = \mathbf{x}_i \mathbf{x}_i^T \mathbf{a} - y_i \mathbf{x}_i \quad (17)$$

We note that:

$$\begin{aligned} (\mathbf{a} - g(\mathbf{x}_i))^T \mathbf{x}_i - y_i &= \mathbf{a}^T \mathbf{x}_i - \mathbf{a}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{x}_i + y_i \mathbf{x}_i^T \mathbf{x}_i - y_i \\ &= \mathbf{a}^T \mathbf{x}_i - (\mathbf{a}^T \mathbf{x}_i)(\mathbf{x}_i^T \mathbf{x}_i) \\ &= \mathbf{a}^T \mathbf{x}_i - \mathbf{a}^T \mathbf{x}_i \\ &= 0 \implies y_i = (\mathbf{a} - g(\mathbf{x}_i))^T \mathbf{x}_i \end{aligned} \quad (18)$$

Then the total training loss among all training samples is given by:

$$\sum_{i=1}^M \frac{1}{2} (\hat{\mathbf{a}}^T \mathbf{x}_i - y_i)^2 \quad (19)$$

By using Eq. 18, we can rewrite Eq. 19 as follows:

$$\begin{aligned} \sum_{i=1}^M \frac{1}{2} (\hat{\mathbf{a}}^T \mathbf{x}_i - y_i)^2 &= \sum_{i=1}^M \frac{1}{2} (\hat{\mathbf{a}}^T \mathbf{x}_i - (\mathbf{a} - g(\mathbf{x}_i))^T \mathbf{x}_i)^2 \\ &= \sum_{i=1}^M \frac{1}{2} ((\hat{\mathbf{a}} - \mathbf{a} + g(\mathbf{x}_i))^T \mathbf{x}_i)^2 \end{aligned} \quad (20)$$

Recall the assumption that $\hat{\mathbf{a}} = \mathbf{a} - \eta \sum_i g(\mathbf{x}_i)$; we rewrite Eq. 20 as follows:

$$\sum_{i=1}^M \frac{1}{2} (\hat{\mathbf{a}}^T \mathbf{x}_i - y_i)^2 = \sum_{i=1}^M \frac{1}{2} (g(\mathbf{x}_i) - \eta \sum_i g(\mathbf{x}_i))^T \mathbf{x}_i)^2 \quad (21)$$

According to the Cauchy–Schwarz inequality and $\|\mathbf{x}_i\| = 1$, the total training loss is bounded by:

$$\begin{aligned}
\sum_{i=1}^M \frac{1}{2} (\hat{\mathbf{a}}^T \mathbf{x}_i - y_i)^2 &\leq \frac{1}{2} \sum_{i=1}^M \|(g(\mathbf{x}_i) - \eta \sum_i g(\mathbf{x}_i))\|^2 * \|\mathbf{x}_i\|^2 \\
&= \frac{1}{2} \sum_{i=1}^M \|(g(\mathbf{x}_i) - \eta \sum_i g(\mathbf{x}_i))\|^2 \\
&= \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^d ((g_j(\mathbf{x}_i) - \eta M \mu_j)^2) \\
&= \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^d ([g_j(\mathbf{x}_i)]^2 - 2\eta M \mu_j g_j(\mathbf{x}_i) + \eta^2 M^2 \mu_j^2) \\
&= \frac{1}{2} \sum_{ij} [g_j(\mathbf{x}_i)]^2 + \sum_j \eta^2 M^2 \mu_j^2 - 2 \sum_j (\eta M \mu_j \sum_i g_j(\mathbf{x}_i)) \\
&= \frac{1}{2} \sum_{ij} [g_j(\mathbf{x}_i)]^2 + \sum_j \eta^2 M^2 \mu_j^2 - 2 \sum_j (\eta M \mu_j M \mu_j) \\
&= \frac{1}{2} G + \sum_j (\eta^2 M^2 \mu_j^2 - 2\eta M^2 \mu_j^2) \\
&= \frac{1}{2} G - \eta M^2 (2 - \eta) \sum_j \mu_j^2
\end{aligned} \tag{22}$$

Since $\sum_{i=1}^M \frac{1}{2} (\hat{\mathbf{a}}^T \mathbf{x}_i - y_i)^2$ is always non-negative, the above upper bound of training loss satisfies:

$$\frac{1}{2} G - \eta M^2 (2 - \eta) \sum_j \mu_j^2 \geq \sum_{i=1}^M \frac{1}{2} (\hat{\mathbf{a}}^T \mathbf{x}_i - y_i)^2 \geq 0 \tag{23}$$

Note that, if $0 < \eta < 2$, then $\eta(2 - \eta) > 0$. Therefore, the larger $\sum_j \mu_j^2$ term would make the upper bound of training loss in Eq. 22 closer to 0. In other words, the higher the gradient absolute mean values across different training samples/batches, the lower the training loss values the model converges to; i.e., the network converges at a faster rate.

In particular, if $\eta = \frac{1}{M}$, the Eq. 22 can be rewritten as:

$$\begin{aligned}
\sum_{i=1}^M \frac{1}{2} (\hat{\mathbf{a}}^T \mathbf{x}_i - y_i)^2 &\leq \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^d ((g_j(\mathbf{x}_i) - \mu_j)^2) \\
&= \frac{1}{2} \sum_j M \sigma_j^2 \\
&= \frac{M}{2} \sum_j \sigma_j^2
\end{aligned} \tag{24}$$

This completes our proof. \square

B PROOF OF THEOREM 3.2

Theorem 3.2 *Given a neural network with ReLU activation function optimized by minimizing Eq. 8, we assume that each initial weight vector $\{\mathbf{w}_r(0), r = 1, \dots, n\}$ is i.i.d. generated from $\mathcal{N}(0, I)$ and the gradient for each weight follows a i.i.d. $\mathcal{N}(0, \sigma)$. For some positive constants δ and ϵ , if the learning rate η satisfies $\eta < \frac{\lambda_0 \sqrt{\pi} \delta}{2M^2 \sqrt{2\Phi(1-\epsilon)} t \sigma}$, then with with probability at least $(1 - \delta)(1 - \epsilon)$, the following holds true: for any $r \in [m]$, $\|\mathbf{w}_r(0) - \mathbf{w}_r(t)\| \leq C = \eta t \sigma \sqrt{\Phi(1 - \epsilon)}$, and at training step t the Gram matrix $H(t)$ satisfies:*

$$\lambda_{\min}(H(t)) \geq \lambda_{\min}(H(0)) - \frac{2\sqrt{2}M^2\eta t\sigma\sqrt{\Phi(1-\epsilon)}}{\sqrt{\pi}\delta} > 0 \quad (25)$$

where $\Phi(\cdot)$ is the inverse cumulative distribution function for a d -degree chi-distribution $\chi^2(d)$.

Proof. We first compute the probability of $\|\mathbf{w}_r(0) - \mathbf{w}_r(t)\| \leq C$. Based on the assumption $w_i(0), i = 1, \dots, n$ follows i.i.d. $\mathcal{N}(0, I)$ and the gradient for each weight follows i.i.d. $\mathcal{N}(0, \sigma)$, considering the weight updating rule defined in Eq. 9, each element in $\mathbf{w}_r(0) - \mathbf{w}_r(t)$ follows a i.i.d. $\mathcal{N}(0, \eta t\sigma)$. Therefore, $\frac{\|\mathbf{w}_r(0) - \mathbf{w}_r(t)\|^2}{\eta^2 t^2 \sigma^2}$ follows the chi-distribution with d degrees of freedom $\chi^2(d)$.

$$\begin{aligned} P(\|\mathbf{w}_r(0) - \mathbf{w}_r(t)\| \leq C) &= P(\|\mathbf{w}_r(0) - \mathbf{w}_r(t)\|^2 \leq C^2) \\ &= P\left(\frac{\|\mathbf{w}_r(0) - \mathbf{w}_r(t)\|^2}{\eta^2 t^2 \sigma^2} \leq \frac{C^2}{\eta^2 t^2 \sigma^2}\right) \\ &= P\left(\frac{\|\mathbf{w}_r(0) - \mathbf{w}_r(t)\|^2}{\eta^2 t^2 \sigma^2} \leq \Phi(1 - \epsilon)\right) \\ &= 1 - \epsilon \end{aligned} \quad (26)$$

Given an input sample \mathbf{x}_i and a weight vector $\mathbf{w}_r(t)$ from $\mathbf{W}(t)$, we define the following event:

$$\mathcal{A}_{ir} = \{\|\mathbf{w}_r(t) - \mathbf{w}_r(0)\| \leq C\} \cap \{\mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(0) \geq 0\} \neq \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(t) \geq 0\}\} \quad (27)$$

If $\|\mathbf{w}_r(t) - \mathbf{w}_r(0)\| \leq C$ holds true,

$$\begin{aligned} \mathbf{x}_i^T \mathbf{w}_r(t) &= \mathbf{x}_i^T (\mathbf{w}_r(t) - \mathbf{w}_r(0)) + \mathbf{x}_i^T \mathbf{w}_r(0) \\ &= \text{sign}(\mathbf{x}_i^T (\mathbf{w}_r(t) - \mathbf{w}_r(0))) \|\mathbf{w}_r(t) - \mathbf{w}_r(0)\| + \text{sign}(\mathbf{x}_i^T \mathbf{w}_r(0)) \|\mathbf{w}_r(0)\| \end{aligned} \quad (28)$$

Eq. 28 tells us that if $\|\mathbf{w}_r(0)\|$ is larger than $\|\mathbf{w}_r(t) - \mathbf{w}_r(0)\|$, then $\mathbf{x}_i^T \mathbf{w}_r(0)$ determines the sign value of $\mathbf{x}_i^T \mathbf{w}_r(t)$; in other words, $\mathbf{x}_i^T \mathbf{w}_r(t)$ always has the same sign values with $\mathbf{x}_i^T \mathbf{w}_r(0)$; i.e., $\mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(0) \geq 0\} = \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(t) \geq 0\}$. That is, if $\|\mathbf{w}_r(t) - \mathbf{w}_r(0)\| \leq C$ and $\mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(0) \geq 0\} \neq \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(t) \geq 0\}$ hold true, then $\|\mathbf{w}_r(0)\| \leq C$. Therefore, the probability of event \mathcal{A}_{ir} :

$$P(\mathcal{A}_{ir}) \leq P(\{\|\mathbf{w}_r(0)\| \leq C\}) \quad (29)$$

By anti-concentration inequality of Gaussian distribution Du et al. (2019b), we have:

$$P(\mathcal{A}_{ir}) \leq P(\{\|\mathbf{w}_r(0)\| \leq C\}) \leq \frac{\sqrt{2}C}{\sqrt{\pi}} \quad (30)$$

Therefore, if any weight vector w_1, \dots, w_m satisfies $\|\mathbf{w}_r(0) - \mathbf{w}_r(t)\| \leq C$, we can bound the entry-wise deviation on the Gram matrix $H(t)$ at training step t : for any $(i, j) \in [n] \times [n]$:

$$\begin{aligned} &\mathbb{E}[|H_{ij}(0) - H_{ij}(t)|] \\ &= \mathbb{E}\left[\frac{1}{m} |\mathbf{x}_i^T \mathbf{x}_j \sum_{r=1}^m (\mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(0) \geq 0, \mathbf{x}_j^T \mathbf{w}_r(t) \geq 0\} - \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(t) \geq 0, \mathbf{x}_j^T \mathbf{w}_r(t) \geq 0\})\right] \\ &= \mathbb{E}\left[\frac{1}{m} |\mathbf{x}_i^T \mathbf{x}_j \sum_{r=1}^m (\mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(0) \geq 0\} \mathbb{I}\{\mathbf{x}_j^T \mathbf{w}_r(t) \geq 0\} - \mathbb{I}\{\mathbf{x}_i^T \mathbf{w}_r(t) \geq 0\} \mathbb{I}\{\mathbf{x}_j^T \mathbf{w}_r(t) \geq 0\})\right] \\ &\leq \mathbb{E}\left[\frac{1}{m} \sum_{r=1}^m (\mathbb{I}\{\mathcal{A}_{ir} \cup \mathcal{A}_{jr}\})\right] \leq P(\mathcal{A}_{ir}) + P(\mathcal{A}_{jr}) \\ &\leq \frac{2\sqrt{2}C}{\sqrt{\pi}} \end{aligned} \quad (31)$$

where the expectation is summing over the initial weight $w(0)$ and the weight $w(t)$ at training step t . Hence, considering all the elements in H , we have:

$$\mathbb{E}\left[\sum_{i=1, j=1}^{M, M} |H_{ij}(0) - H_{ij}(t)|\right] \leq \frac{2M^2\sqrt{2}C}{\sqrt{\pi}} \quad (32)$$

Therefore, by Markov’s inequality, given the probability $1 - \delta$, we get:

$$\sum_{i=1, j=1}^{M, M} |H_{ij}(0) - H_{ij}(t)| \leq \frac{2M^2\sqrt{2}C}{\sqrt{\pi}\delta} \quad (33)$$

In Du et al. (2019b), the authors prove that, given a small perturbation K :

$$\text{if } \left[\sum_{ij} |H_{ij}(0) - H_{ij}| \right] \leq K, \text{ then } \lambda_{\min}(H) \geq \lambda_{\min}(H(0)) - K \quad (34)$$

In our case, K in Eq. 34 is given by $\frac{2M^2\sqrt{2}C}{\sqrt{\pi}\delta}$. Therefore,

$$\lambda_{\min}(H(t)) \geq \lambda_{\min}(H(0)) - \frac{2M^2\sqrt{2}C}{\sqrt{\pi}\delta} = \lambda_{\min}(H(0)) - \frac{2\sqrt{2}M^2\eta t\sigma\sqrt{\Phi(1-\epsilon)}}{\sqrt{\pi}\delta} \quad (35)$$

We replace the term η in Eq.35 with η ’s upper bound given in the assumption of Theorem 3.2, i.e., $\eta < \frac{\lambda_0\sqrt{\pi}\delta}{2M^2\sqrt{2}\Phi(1-\epsilon)t\sigma}$, we can get that $\lambda_{\min}(H(t))$ is always larger than 0; that is:

$$\lambda_{\min}(H(t)) \geq \lambda_{\min}(H(0)) - \frac{2\sqrt{2}M^2\eta t\sigma\sqrt{\Phi(1-\epsilon)}}{\sqrt{\pi}\delta} > 0 \quad (36)$$

This completes our proof. \square

C DETAILS OF ZICO ON IMAGENET

Search space We use the commonly used MobileNetv2-based search space where the candidate networks are built by stacking multiple Inverted Bottleneck Blocks (IBNs) with SE modules Sandler et al. (2018); Pham et al. (2018); Lin et al. (2021); all the SE modules share the same `se_ratio` as 0.25. For each of IBN, we vary the kernel size of the depth-wise convolutional layer from $\{3,5,7\}$ and sample the expansion ratio from $\{1,2,4,6\}$. We primarily consider ReLU as the activation function. For each point-wise convolutional layer, the range of the number of channels is from 8 to 1024 with a step size of 8. We use standard `Kaiming_Init` to initialize all linear and convolution layers for every candidate networks He et al. (2015).

Search details: We use Algorithm 1 (main paper) to conduct the zero-shot search. We repeat the search 10^5 times with the population size $E = 512$. For each of the candidate architectures, we compute ZiCo with two batches randomly sampled from the training set of ImageNet with batch size 128. In total, it takes 10 hours on a single NVIDIA 3090 GPU for 10^5 search steps.

Training details: We use the same data augmentations configurations as in Pham et al. (2018): mix-up, label-smoothing, random erasing, random crop/resize/flip/lighting and AutoAugment. We use the SGD optimizer with momentum 0.9 and weight decay $4e-5$. We take EfficientNet-B3 as a teacher network and use the knowledge distillation method to train the network. We set the initial learning rate as 0.1 and used the cosine annealing scheme to adjust the learning rate during training. We train the obtained network 480 epochs, which takes 83 hours on a 40-core Intel Xeon CPU and 8 NVIDIA 3090 GPU-powered server.

D SUPPLEMENTARY RESULTS ON NAS BENCHMARKS

We provide some illustration figures of real test accuracy vs. various proxies on NATSBench-SSS search space for CIFAR10 (Fig. 4) and ImageNet16-120 datasets(Fig. 5). We also show the same illustrative results (real test accuracy vs. various proxies) on NASBench101 search space in Fig. 6.

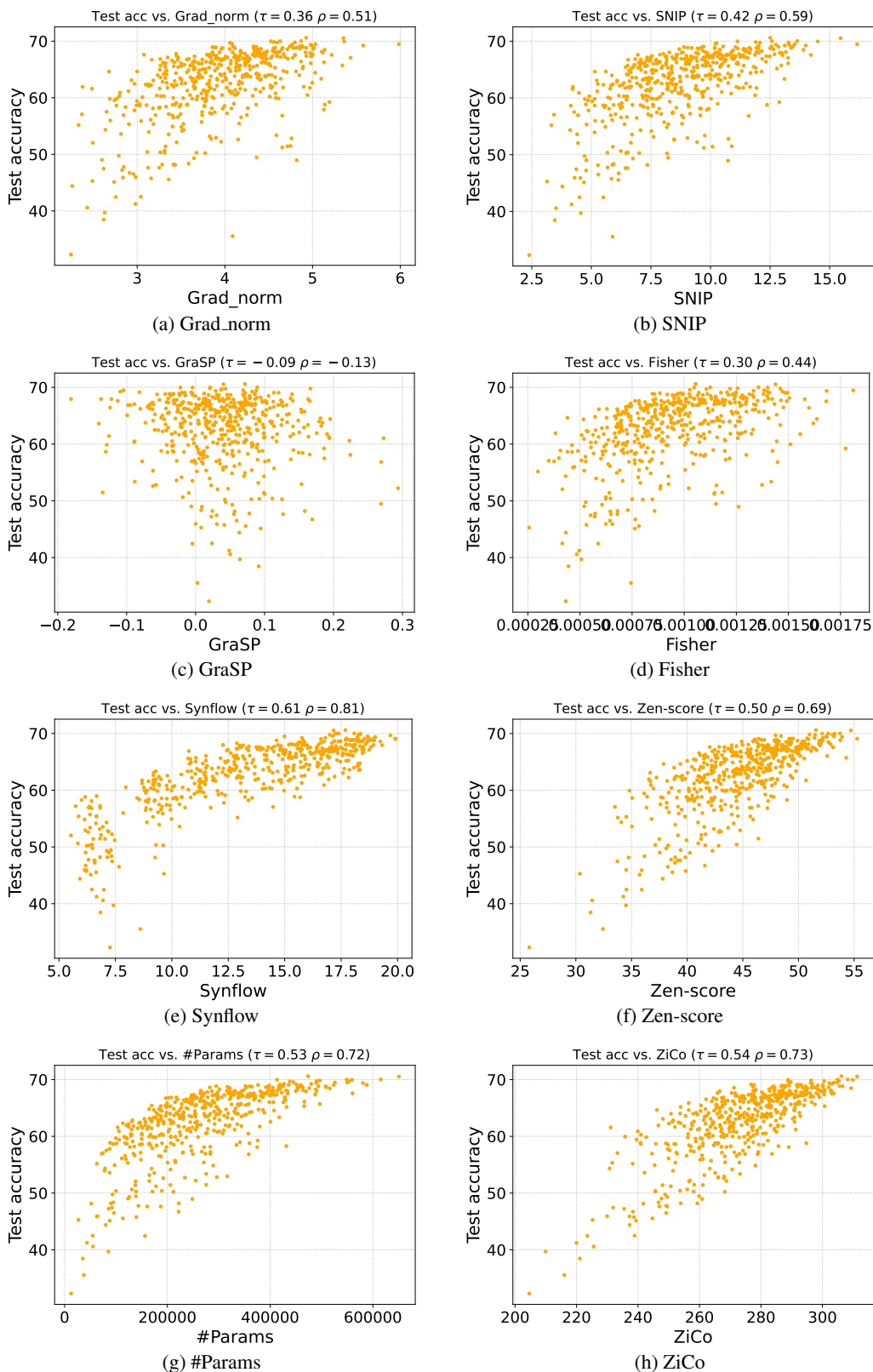


Figure 4: Real test accuracy vs. various proxies on NATSBench-SSS search space for CIFAR10 dataset. τ and ρ are short for Kendall’s τ and Spearman’s ρ , respectively.

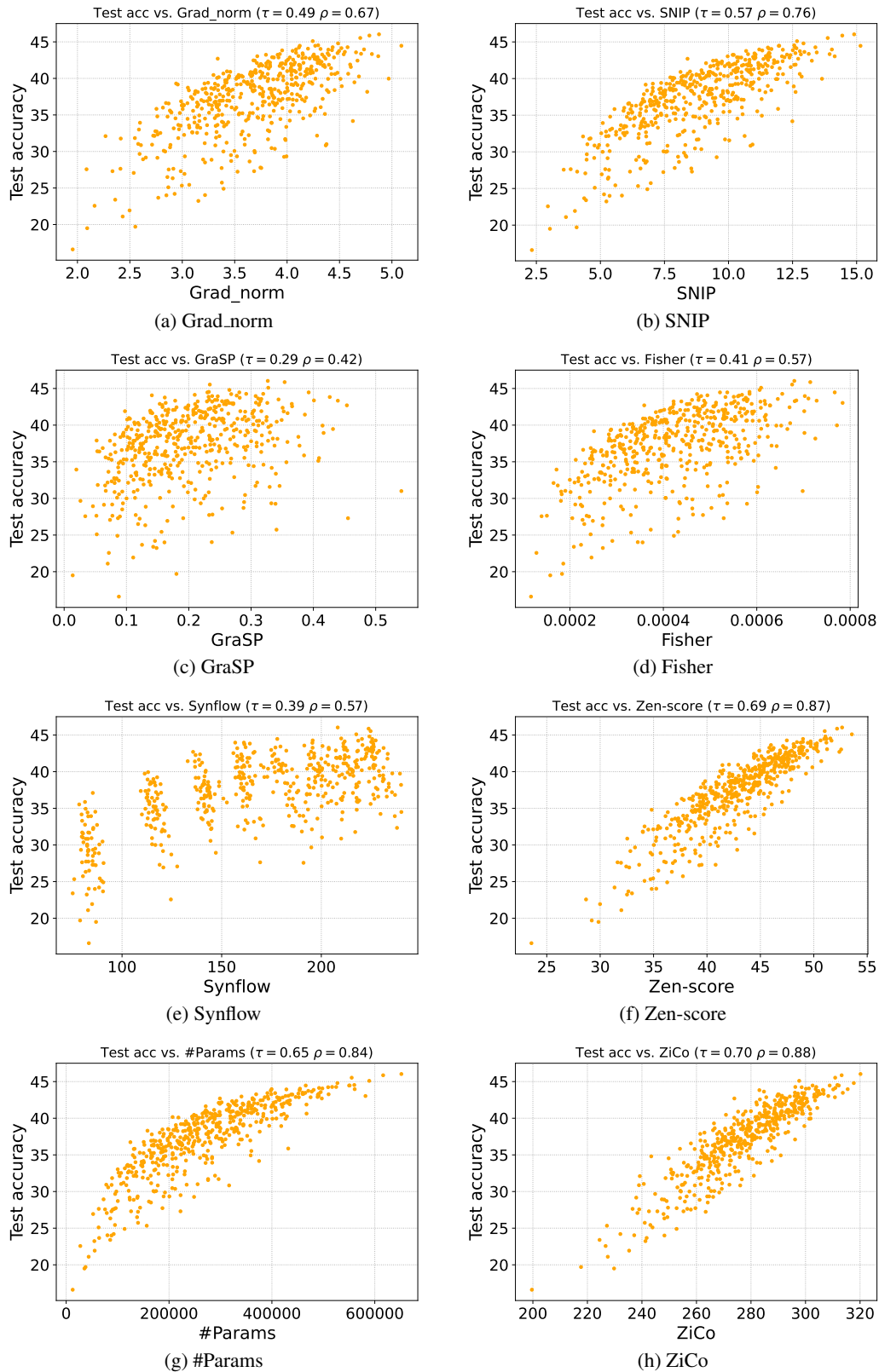


Figure 5: Real test accuracy vs. various proxies on NATSBench-SSS search space for ImageNet16-120 dataset. τ and ρ are short for Kendall’s τ and Spearman’s ρ , respectively.

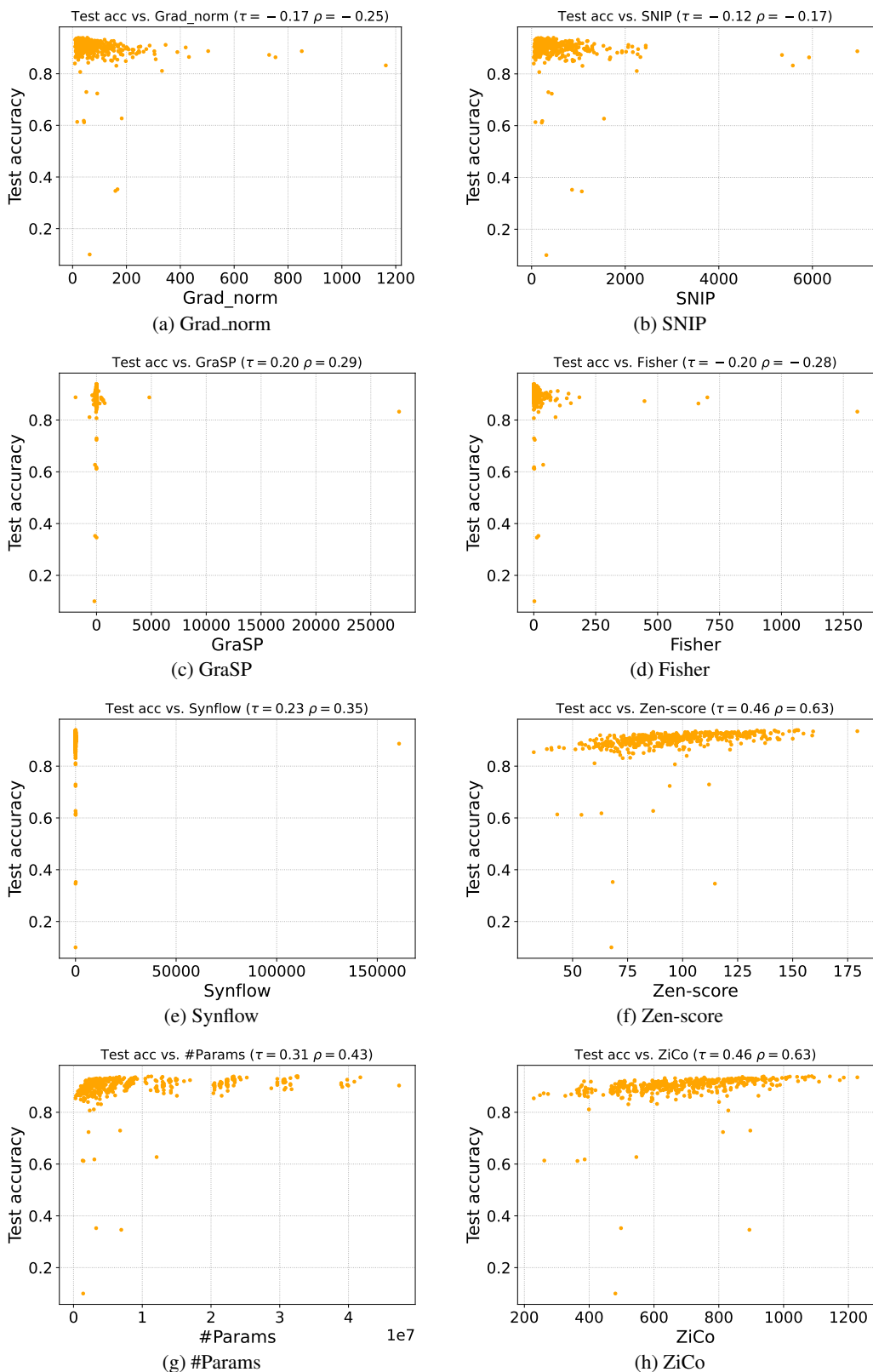


Figure 6: Real test accuracy vs. various proxies on NASBench101 search space for CIFAR10 dataset. τ and ρ are short for Kendall’s τ and Spearman’s ρ , respectively.