One-Shot Neural Architecture Search: Maximising Diversity to Overcome Catastrophic Forgetting

Miao Zhang[®], Huiqi Li[®], *Senior Member, IEEE*, Shirui Pan[®], Xiaojun Chang, Chuan Zhou[®], Zongyuan Ge[®], and Steven Su[®], *Senior Member, IEEE*

Abstract—One-shot neural architecture search (NAS) has recently become mainstream in the NAS community because it significantly improves computational efficiency through weight sharing. However, the supernet training paradigm in one-shot NAS introduces catastrophic forgetting, where each step of the training can deteriorate the performance of other architectures that contain partially-shared weights with current architecture. To overcome this problem of catastrophic forgetting, we formulate supernet training for one-shot NAS as a constrained continual learning optimization problem such that learning the current architecture does not degrade the validation accuracy of previous architectures. The key to solving this constrained optimization problem is a novelty search based architecture selection (NSAS) loss function that regularizes the supernet training by using a greedy novelty search method to find the most representative subset. We applied the NSAS loss function to two one-shot NAS baselines and extensively tested them on both a common search space and a NAS benchmark dataset. We further derive three variants based on the NSAS loss function, the NSAS with depth constrain (NSAS-C) to improve the transferability, and NSAS-G and NSAS-LG to handle the situation with a limited number of constraints. The experiments on the common NAS search space demonstrate that NSAS and it variants improve the predictive ability of supernet training in one-shot NAS with remarkable and efficient performance on the CIFAR-100, and ImageNet datasets. The results with the NAS benchmark dataset also confirm the significant improvements these one-shot NAS baselines can make.

Index Terms—AutoML, neural architecture search, continual learning, catastrophic forgetting, novelty search

1 INTRODUCTION

NEURAL architecture search (NAS) has re7cently attracted massive interest from the deep learning community because experts do not have inordinate amounts of time and labor designing neural networks [13], [18], [29], [35],

- Chuan Zhou is with the Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100081, China. E-mail: zhouchuan@amss.ac.cn.
- Zongyuan Ge is with the Monash e-Research Centre, Monash University, Clayton, VIC 3800, Australia. E-mail: Zongyuan.Ge@monash.edu.
- Steven Su is with the Faculty of Engineering and Information Technology, University of Technology Sydney, Ultimo, NSW 2007, Australia. E-mail: steven.su@uts.edu.au.

Manuscript received 15 Mar. 2020; revised 14 Sept. 2020; accepted 26 Oct. 2020. Date of publication 3 Nov. 2020; date of current version 4 Aug. 2021. (Corresponding authors: Huiqi Li and Shirui Pan.) Recommended for acceptance by S. Dickinson. Digital Object Identifier no. 10.1109/TPAMI.2020.3035351 [42], [47], [50], [71]. Early NAS methods were based on a nested approach that trained numerous separate architectures from scratch and then used reinforcement learning (RL) or an evolutionary algorithm (EA) to find the most promising architectures, based on validation accuracy [19], [46], [72]. However, these methods are so computationally-expensive as to be impractical for most machine learning practitioners. For example, it would take more than 1800 GPU days through RL to find promising architectures for the problem outlined in [72], and Real et al. [46] spent 7 days with 450 GPUs searching for promising architectures with an EA. Recent studies have shown that NAS can significantly improve computational efficiency [3], [6], [63]. Weight sharing, in particular, also called one-shot NAS [4], [45], [66], has attracted enormous attention for automating neural architecture design. This is because it not only finds state-of-the-art architectures but also significantly reduces the search hours needed. One-shot NAS encodes the search space as a supernet, where all possible architectures directly inherit weights from the supernet for evaluation without needing to be trained from scratch. Since one-shot NAS only trains the supernet for architecture searches, this learning paradigm might reduce the time a search takes from many days down to several hours.

Pioneer studies on one-shot NAS follow two sequential steps [4], [14], [30], [45]. They first adopt an architecture sampling controller to sample architectures for training the supernet. Then, a heuristic search method finds promising architectures over a discrete search space based on the trained supernet [20], [30], [45], [58]. Later studies [7], [16],

0162-8828 © 2020 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.

Miao Zhang is with the School of Information and Electronics, Beijing Institute of Technology, Beijing 100081, China, and with the Faculty of Information Technology, Monash University, Clayton, VIC 3800, Australia, and also with the Faculty of Engineering and Information Technology, University of Technology Sydney, Ultimo, NSW 2007, Australia. E-mail: Miao.Zhang-2@student.uts.edu.au.

Huiqi Li is with the School of Information and Electronics, Beijing Institute of Technology, Beijing 100081, China. E-mail: huiqili@bit.edu.cn.

Shirui Pan is with the Faculty of Information Technology, Monash University, Clayton, VIC 3800, Australia. E-mail: shirui.pan@monash.edu.

Xiaojun Chang is with the Faculty of Information Technology, Monash University, Clayton, VIC 3800, Australia, and also with the Faculty of Computing and Information Technology, King Abdulaziz University, Jeddah 21589, Saudi Arabia. E-mail: xiaojun.chang@monash.edu.



Fig. 1. Left: The general process of one-shot NAS. First, the search space is defined as a supernet containing all candidate architectures. Then a single path of the supernet (an architecture) is trained in each step of the supernet training process. Promising architectures are selected based on the validation accuracy of weights inherited from the trained supernet without the need for training from scratch. *Right*: The validation accuracy for four different architectures during the supernet training. The solid lines ("Arch") are the accuracies returned using weights inherited from the supernet; the dashed lines ("Arch-R") are the accuracies after retraining.

[38], [39], [55], [59], [64] have further employed continuous relaxation to differentiate between architectures so that the gradient descent can be used to optimize the architecture with respect to validation accuracy. The architecture parameters and supernet weights are alternatively optimized through a bilevel optimization method, and the most promising architecture is obtained once the supernet is trained.

Since one-shot NAS evaluates candidate architectures based on the validation accuracy of the weights it inherits from the supernet as opposed to training them from scratch, the success of one-shot NAS relies on a critical assumption that the validation accuracy should approximate the test accuracy after training from scratch or be highly predictive. The authors of the first study on one-shot NAS [4] observed a strong positive correlation between the validation accuracy and the test accuracy when the supernet was trained through random path dropout. Subsequent studies all rightly considered this assumption to be true for all one-shot NAS methods. However, several recent studies have revealed that this assumption may not hold in most popular one-shot NAS approaches. For instance, Sciuto et al. [61] show that there is no observable correlation between the validation and test accuracy of the weight-sharing paradigm with ENAS [45], and Adam et al. [1] show that the RNN controller in ENAS does not depend on past sampled architectures, which means its performance is the same as a random search. Similarly, Singh et al. [51] find that there is no visible progress in terms of the retrained performance for found architectures based on supernet during the architecture search phase, implying the supernet training is useless for improving the predictive ability of one-shot NAS. Further, Yang et al. [57] conducted extensive experiments that demonstrated that the current one-shot NAS techniques struggle to outperform naive baselines. Rather, the success of one-shot NAS is mostly due to the design of the search space.

Most one-shot NAS approaches [7], [14], [20], [30] adopt a single-path training method for their supernet training, where only a single path (one architecture) in the supernet is trained in each step. This is the scenario we consider. However, Benyahia *et al.* [5] observed that when training multiple models (architectures) with partially-shared weights for a single task, the training of each model may lower the performance of other models. Benyahia *et al.* [5] defined this phenomenon

as multi-model forgetting, also known as catastrophic forgetting. They also observed this catastrophic forgetting in one-shot NAS. For example, consider a large supernet containing multiple models with shared weights across them. Sequentially training each model on a single task could mean that the accuracy of each model tends to drop when training another model containing partially-shared weights [5], [61]. This multi-model forgetting in one-shot NAS is illustrated in Fig. 1 in terms of the validation accuracy of four different architectures during supernet training. What is clear from the figure is that inheriting weights makes performance deteriorate even further during supernet training.

So, although weight sharing can greatly reduce computation hours, it can also introduce catastrophic forgetting into the supernet training, which results in unreliable architecture rankings. Addressing multi-model forgetting during supernet training is an urgent issue if we are to better leverage one-shot NAS and improve the predictive ability of supernets. Hence, we have formulated supernet training as a constrained optimization problem for continual learning to avoid degrading the performance of previous architectures when training a new one.

That said, it is intractable to consider all previously visited architectures. Therefore, only the most representative subset of previous architectures is used to regularize learning of the current architecture. We have devised an efficient greedy novelty search method based on maximizing diversity to select the constraints. We have also implemented the approach in two one-shot baselines. The experimental results demonstrate that our strategy is able to relieve multi-model forgetting in one-shot NAS methods. A summary of our main contributions follows.

- We first formulate supernet training with one-shot NAS as a constrained optimization problem of continual learning, where learning the current architecture should not degrade the performance of previous architectures with partially-shared weights.
- We have also designed an efficient greedy novelty search method based on maximizing diversity to select a subset of the constraints that best approximate the feasible region formed by all previous architectures.

- With these two techniques, we then incorporate this *NSAS* loss function (novelty search-based architecture selection) into the RandomNAS [30] and GDAS [16] one-shot NAS baselines to form RandomNAS-NSAS and GDAS-NSAS, with the goal of reducing the level of multi-model forgetting during supernet training. Our best-found models from the common search space [38] returned a competitive test error of 2.59 percent on CIFAR-10 and only took 0.7 GPU days of search time.
- To improve transferability, we further devised a variant of NSAS, called *NSAS-C*, which searches for "deeper" architectures in the convolutional cell search. Experiments on the common search space demonstrate increased transferability of the found models, and competitive test errors of 16.69 percent on CIFAR-100 and 25.5 percent on ImageNet.
- A series of experiments conducted with the NAS benchmark dataset NAS-Bench-201 [17] also verify that *NSAS* significantly reduce forgetting and improve the performance of one-shot NAS baselines.

This paper outlines some significant extensions to our recent conference paper [65]. These include: (1) improvements to the transferability of the found models through a variant of novel search strategy called *NSAS-C. NSAS-C* has a depth constraint for convolutional architecture searches. Experiments with ImageNet demonstrate its efficacy; (2) a new NAS search space, NSA-Bench-201, along with comparative experiments against relevant baselines; (3) two new variants of NSAS, *NSAS-G* and *NSAS-LG*, to handle situations with a limited number of constraints. Comparison experiments with the NAS-Bench-201 dataset are provided to showcase these extensions; and (4) an impact analysis of the hyperparameters settings and constraint selection strategy in Sections 5.2.2 and 5.2.3.

2 BACKGROUND

2.1 Neural Architecture Search

The goal of NAS is to automatically design deep neural networks without human intervention. In general, the architecture of a deep neural network α is usually represented as a directed acyclic graph (DAG), which is also a subgraph of the whole search space $\alpha \in A$. A deep neural network could be defined as $\mathcal{U}(\alpha, w_{\alpha})$, where w_{α} are the weights associated with architecture α . With NAS, one tries to find the architecture with the best validation performance according to:

$$\alpha^* = \underset{\alpha \in \mathcal{A}}{\operatorname{argmin}} \ \mathcal{L}_{\operatorname{val}}(\mathcal{U}(\alpha, w_{\alpha}^*)), \tag{1}$$

where w_{α}^{*} is derived by training architecture α on the training set while minimizing the training loss function \mathcal{L}_{train} :

$$w_{\alpha}^{*} = \underset{w}{\operatorname{argmin}} \mathcal{L}_{\operatorname{train}}(\mathcal{U}(\alpha, w_{\alpha})).$$
 (2)

Early studies on NAS usually used a nested approach to finding promising architectures by training numerous architectures from scratch and leveraging EA [46] or RL [71] to reveal the promising ones. However, from a practical standpoint, it is computationally inefficient and often unaffordable to evaluate numerous architectures in this way. Therefore, more recently, researchers have shifted their attention to reducing computation costs with strategies such as performance prediction [3], [54], weights generation [6], [63], weight sharing [38], [45], and so on [72].

2.2 One-Shot Neural Architecture Search

One-shot NAS encodes a search space \mathcal{A} as a supernet $\mathcal{W}_{\mathcal{A}}$ that consumes all possible candidates. Only the supernet is trained, while all candidate architectures α directly inherit weights from the supernet without needing to be trained from scratch. Search times are therefore greatly reduced because only one neural network needs to be trained during the architecture search phase. The most promising architecture α^* is based on validation performance with weights inherited from the supernet:

$$\min_{\alpha \in \mathcal{A}} \quad \mathcal{L}_{\text{val}}(\mathcal{W}_{\mathcal{A}}^*(\alpha))$$
s.t. $\mathcal{W}_{\mathcal{A}}^*(\alpha) = \operatorname{argmin} \mathcal{L}_{\text{train}}(\mathcal{W}_{\mathcal{A}}(\alpha)).$

$$(3)$$

Eq. (3) is more than a challenging bilevel optimization problem: the discrete characteristic of the architecture space makes it impossible to use a gradient-based method to solve the formula directly. For this reason, ENAS [45] uses an LSTM controller to sample the architectures. Whereas, [20] and [30] train the supernet based on a uniform sampling strategy and the best-performing architecture from the trained supernet is found through a random search or evolutionary method.

Several state-of-the-art one-shot methods use continuous relaxation to transform discrete architectures into a continuous space A_{θ} with a **softmax** function to further improve efficiency [16], [38], [43], [55]. The supernet weights and architecture parameters can be jointly optimized through:

$$(\alpha_{\theta}^{*}, \mathcal{W}_{\mathcal{A}_{\theta}}(\alpha_{\theta}^{*})) = \underset{\alpha_{\theta}, \mathcal{W}}{\operatorname{argmin}} \mathcal{L}_{\operatorname{train}}(\mathcal{W}_{\mathcal{A}_{\theta}}(\alpha_{\theta}^{*})), \tag{4}$$

making it possible to continually optimize the architecture search. The best architecture α^* is determined through argmax based on the continuous architecture representation α^*_{θ} .

Since Eq. (4) is supposed to train the entire supernet in each step, it has a much higher memory requirement than ENAS. Hence, ProxylessNAS [7] transforms the real-valued architecture parameters into binary representations through binary gates, and only a single path is activated during the supernet training. In this way, the memory requirement for ProxylessNAS is the same as training a single architecture. GDAS [16] introduces a gradient-based sampler to sample the single path for each training step. Additionally, the distribution of architectures and the supernet weights can be jointly optimized, which means the memory requirement also equates to only training a single architecture. Yao *et al.* [59] developed a constrained optimization method to force each step of the architecture optimization process in the continuous space to arrive at a binary result, thus reducing the memory requirement of supernet training. Unlike continuous relaxation, NAO [39] uses an LSTM-based autoencoder to transform discrete neural architectures into continuous representations. A differentiable method is then used to search for architectures in the continuous space.

2.3 Multi-Model Forgetting in One-Shot NAS

Catastrophic Forgetting is a common problem in artificial general intelligence and multi-task learning. It describes the phenomenon of where a model forgets what it has learned about a previous task(s) after being trained on a new task [9], [22], [26], [28], [34], [44]. Formally, a model with optimal parameters θ_A^* for dataset \mathcal{D}_A will perform substantially less well on \mathcal{D}_A after it has trained on another dataset \mathcal{D}_B . Methods to resolve such issues are defined as *continual learning*. Some examples include learning without forgetting (LwF) [33], which adds a response from the old task as a regularization term to prevent catastrophic forgetting, and elastic weight consolidation (EWC) [26], which maximizes the likelihood of a conditional probability $p(\theta | \mathcal{D})$, where \mathcal{D} containing two independent data sets \mathcal{D}_A and \mathcal{D}_B , and \mathcal{D}_A is not available when trained on \mathcal{D}_B .

Multi-model Forgetting occurs when training multiple models with a single dataset. Unlike training a model on several tasks sequentially, one-shot NAS applies different models to a single dataset \mathcal{D} , e.g., $\theta_a = (\theta_a^p, \theta^s)$ and $\theta_b = (\theta_b^p, \theta^s)$, to a single dataset \mathcal{D} , where θ^s is the shared weight and θ_a^p and θ_b^p are private weights. Several recent studies [31], [53], [61] have shown that the interactions between networks can degrade the performance of a whole network, and that catastrophic forgetting with one-shot NAS can lower the performance of previous architectures after training a new architecture in the supernet. To alleviate this problem, Benyahia *et al.* [5] proposed a weight plasticity loss (WPL), which maximizes the posterior probability $p(\theta_a^p, \theta_b^p, \theta^s | \mathcal{D})$ as:

$$p(\theta \mid \mathcal{D}) = \frac{p(\theta_a^p, \theta_b^p, \theta^s, \mathcal{D})}{p(\mathcal{D})} = \frac{p(\theta_a^p \mid \theta_b^p, \theta^s, \mathcal{D})p(\theta_b^p, \theta^s, \mathcal{D})}{p(\mathcal{D})}$$
$$= \frac{p(\theta_a^p, \theta^s \mid \mathcal{D})p(\mathcal{D} \mid \theta_b^p, \theta^s)p(\theta_b^p, \theta^s)}{p(\theta^s, \mathcal{D})}$$
$$= \frac{p(\theta_a \mid \mathcal{D})p(\mathcal{D} \mid \theta_b)p(\theta_b)}{p(\theta^s, \mathcal{D})}.$$
(5)

However, it is intractable to calculate (θ^s, D) in Eq. (5), so Benyahia *et al.* [5] made several presuppositions to make the calculation feasible: a) that θ_a^p and θ_s are independent, and b) that the weights θ_a for previous model are in optimal points. This way, $p(\theta^s, D)$ can be estimated by the distance of θ^s to the optimal θ_s^* with the diagonal of the Fisher information defining the importance of each parameter. In WPL, the loss function to maximize the likelihood of $p(\theta_a^p, \theta_b^p, \theta^s | D)$ is calculated as:

$$\mathcal{L}_{WPL}(\theta_b) = \mathcal{L}_c(\theta_b) + \frac{\eta}{2} (\left\|\theta_b^p\right\|^2 + \left\|\theta^s\right\|^2) + \sum_{\theta_{s_i} \in \mathbf{\theta}_s} \frac{\varepsilon}{2} F_{\theta_{s_i}}(\theta_{s_i} - \theta_{s_i}^*),$$
(6)

where \mathcal{L}_c is the cross-entropy loss function, and $F_{\theta_{s_i}}$ is the diagonal element of the Fisher information matrix corresponding to parameter θ_{s_i} . $F_{\theta_{s_i}}$ is estimated by presupposing parameters (θ_a^p, θ_b^p) are independent, and that θ_s^* are the shared parameters θ^s after the previous model has been trained, which are assumed to be in the optimal points. A detailed derivation of Eq. (6) can be found in [5].

Limitations weight plasticity loss (WPL) only considers one previous architecture in each step of supernet training. This method is also based on the assumption that the shared weights are optimal. However, these two assumptions are hard to hold when training a supernet in a one-shot NAS scheme given numerous architectures shared weights with the current architecture. Plus, the shared weights are usually far away from the optimal points. To address these concerns, we formulated supernet training with one-shot NAS as a constrained optimization problem, where learning the current architecture does not degrade the performance of previously-visited architectures. We consider a subset of previous architectures as constraints to regularize the learning of the current architecture. We also demonstrate that the loss function of the posterior probability $p(\theta^p_a, \theta^p_b, \theta^s | D)$ can be calculated without assuming that the shared weights are optimal when maximizing the diversity of the selected architectures.

3 METHODOLOGY

3.1 Problem Formulation

Unlike jointly optimizing the posterior probability under the assumption that θ_a is near-optimal as per WPL [5] or keeping the shared weights fixed as per Learn to Grow [32], we formulate supernet training as a constrained optimization problem. Specifically, we enforce the architectures with inherited weights in the current step so as to perform better than the last step. Without loss of generality, we consider a typical scenario where only one architecture in the supernet is trained in each step, and the constrained optimization problem is defined as:

$$\begin{split} & \mathcal{W}_{\mathcal{A}}^{t} = \underset{\boldsymbol{\theta} \in \mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})}{\operatorname{argmin}} \quad \mathcal{L}_{\operatorname{train}}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})), \\ & \text{s.t.} \quad \mathcal{L}_{\operatorname{train}}(\mathcal{W}_{\mathcal{A}}^{t}(\boldsymbol{\alpha}^{i})) \leq \mathcal{L}_{\operatorname{train}}(\mathcal{W}_{\mathcal{A}}^{t-1}(\boldsymbol{\alpha}^{i})); \ \forall i \in \{0...t-1\}. \end{split}$$

(7)

Here, $\mathcal{L}_{train}(\mathcal{W}_{\mathcal{A}}(\alpha)) = \mathcal{L}_{c}(\mathcal{W}_{\mathcal{A}}(\alpha)) + \lambda \mathcal{R}(\mathcal{W}_{\mathcal{A}}(\alpha))$, and $\mathcal{W}_{\mathcal{A}}$ represents the total of all weights in the supernet. α^{t} is the current architecture in step *t*, and $\mathcal{W}_{\mathcal{A}}(\alpha^{t})$ is the weights of α^{t} inherited from the supernet, and only $\mathcal{W}_{\mathcal{A}}(\alpha^{t})$ is optimized in each step *t*.

Algorithm 1. Greedy Novelty Search

Input: constraints archive \mathcal{M} , recent architectures archive \mathcal{C} , selected architecture α^m , *n*.

- N(α^m, M) ← calculate the novelty score of α^m in M based on Eq. (9);
- 2: for i = 1, 2, ..., n do
- 3: randomly sample an architecture α^r from *C*;
- 4: if $N(\alpha^r, \mathcal{M}) > N(\alpha^m, \mathcal{M})$ then
- 5: replace α^m with α^r ;
- 6: **end if**
- 7: end for

3.2 Constraints Selection Based on Novelty Search

The constraints in Eq. (7) prevent the learning the current architecture from degrading the performance of previous architectures as a strategy to overcome multi-model forgetting in one-shot NAS. However, the number of constraints in Eq. (7) increases linearly with the step, which makes it intractable to consider all constraints in the optimization. Therefore, it is more practical to try and select a subset of M constraints from the previous architectures that forms as close a feasible region to the original feasible region as possible. Intuitively, maximizing the diversity of the subset is an efficient way to find the most representative samples from the previous architectures. Based on this observation and motivated by [2], we propose a surrogate for constraint selection:

maximize_{$$\mathcal{M}$$} $\sum_{\alpha^{i},\alpha^{j}\in\mathcal{M}} dis(\alpha^{i},\alpha^{j}),$
s.t. $\mathcal{M} \subset \{\alpha^{1}...\alpha^{t-1}\}; |\mathcal{M}| = M,$ (8)

where $dis(\alpha^i, \alpha^j)$ is a function to calculate the distance between architectures. Further, to solve this equation, we proposed a greedy novelty search method to maximize the diversity of the subset. Before the archive is full, all the new coming architectures are added into the subset. Once full, the most similar one to the current architecture is chosen to replace the one that maximizes the novelty score of the archive. Algorithm 1 sets out a simple implementation of our greedy novelty search method. A simple and standard novelty measurement, defined as $N(\alpha, \mathcal{M})$, measures the mean distance of its *k*-nearest neighbors in \mathcal{M} :

$$N(\alpha, \mathcal{M}) = \frac{1}{|S|} \sum_{\alpha^{j} \in S} dis(\alpha, \alpha^{j})$$

$$S = kNN(\alpha, \mathcal{M}) = \{\alpha^{1}, \alpha^{2}, \dots, \alpha^{k}\}.$$
(9)

In this paper, we measure the difference of the input edges for each node in an architecture. The input edge of the same node for two architectures is considered to be the same only when the two edges have the same input node and the same operations. M constraint architectures are then selected from |C| recent architectures rather than all previous architectures.

3.3 The NSAS Loss Function

After finding the *M* most representative architectures $\{\alpha_1, \ldots, \alpha_M\}$ by maximizing diversity, we need to forcibly optimize the learning of the current architecture in the feasible region formed by these constraints. A common approach is to convert the constraints into a soft regularization loss or apply a replay buffer [2]. The weights of these architectures in the subset are described as $\{\theta_1, \ldots, \theta_M\}$. When the selected constraints are converted to a soft regularization loss, the loss function for the constrained optimization problem in Eq. (7) could be described as:

$$\mathcal{L}_{N}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})) = (1 - \beta)(\mathcal{L}_{c}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})) + \lambda \mathcal{R}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t}))) + \frac{\beta}{M} \sum_{i=1:M} (\mathcal{L}_{c}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{i})) + \lambda \mathcal{R}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{i}))),$$
(10)

where \mathcal{L}_c is the cross-entropy loss function, \mathcal{R} is the ℓ_2 regularization term, and β are the trade-offs. The term $\mathcal{L}_N(\mathcal{W}_A(\alpha^t))$ is the **NSAS** loss function. While learning the current architecture α_t , **NSAS** protects the performance of those constraint architectures by ensuring the shared parameters stay in a region of low error for these constraints, as shown schematically in Fig. 2.



Fig. 2. *NSAS* loss function ensures that the learning of current architecture will not deteriorate the performance of previous architectures in the constraint subset.

3.4 From Weight Plasticity Loss (WPL) to NSAS

WPL [5] regularizes the learning of current architecture by maximizing the posterior probability $p(\theta_a^p, \theta_b^p, \theta^s | D)$, where $\theta_a = \{\theta_a^p, \theta^s\}$ is the weights of the last architecture, $\theta_b = \{\theta_b^p, \theta^s\}$ is the weights of the current architecture, and θ^s is their shared weights. Unlike WPL, which only considers one previous architecture, we consider a subset of previously visited architectures $-\theta_a = \{\theta_1, \dots, \theta_M\} = \{(\theta_1^p, \theta_1^s), \dots, (\theta_M^p, \theta_M^s)\}$ - where θ_i^p is the private weights, and θ_i^s is the weights shared with the current architecture. When selected constraints make the following two assumptions hold true, then Lemma 1 describes the relationship between WPL and NSAS.

- **Assumption 1.** The architectures in the constraint subset cover all weights of the current architecture α_t that $\theta_b^{(e)} \subseteq \{\theta_1^{(e)} \cup ... \cup \theta_M^{(e)}\}$ for every edge e in α_t , where $\theta_m^{(e)}$ is the weight of the operations assigned to each edge of α_m .
- **Assumption 2.** There are no shared weights among these constraint architectures, *i.e.*, $\theta_1 \cap \theta_2 = ... = \theta_{M-1} \cap \theta_M = \emptyset$.
- **Lemma 1.** When the selected constraint architectures satisfy the above two assumptions, the NSAS loss function with the WPL can be formulated as:

$$\mathcal{L}_{N}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})) = \mathcal{L}_{c}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})) + \gamma \mathcal{L}_{WPL}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})).$$
(11)

Proof. Since the weights of current architecture θ_b are shared by the constraints described in **Assumption 1**, $\theta_b^{(e)} \subseteq \{\theta_1^{(e)} \cup ... \cup \theta_M^{(e)}\}$, and for every edge e in α_t , we have $\theta_b^p = \emptyset$. Further, θ_i and θ_j (i, j = 1...M) are independent as the architectures are trained separately without shared weights, as described in **Assumption 2**. Now the posterior probability in the WPL can be written as:

$$p(\theta \mid \mathcal{D}) = p(\theta_1 \dots \theta_M, \theta_b \mid \mathcal{D}) = \frac{p(\theta_1^p \dots \theta_M^p, \theta_1^s \dots \theta_M^s, \mathcal{D})}{p(\mathcal{D})}$$
$$= \frac{p(\theta_1 \dots \theta_M, \mathcal{D})}{p(\mathcal{D})} = \frac{p(\theta_1 \mid \theta_2 \dots \theta_M, \mathcal{D})p(\theta_2 \dots \theta_M, \mathcal{D})}{p(\mathcal{D})}$$
$$= \prod_{i=1:M} p(\theta_i \mid \mathcal{D}) \propto \prod_{i=1:M} p(\mathcal{D} \mid \theta_i)p(\theta_i)$$
$$= p(\theta) \prod_{i=1:M} p(\mathcal{D} \mid \theta_i) = p(\theta^t) \prod_{i=1:M} p(\mathcal{D} \mid \theta_i),$$

$$\mathcal{L}_{WPL}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})) = \epsilon \mathcal{R}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})) + \sum_{i=1:M} \mathcal{L}_{c}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{i})), \quad (13)$$

where ϵ is the trade-off factor. And the proposed **NSAS** loss function with the WPL can be expressed as:

$$\mathcal{L}_{N}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})) = (1 - \beta)(\mathcal{L}_{c}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})) + \lambda \mathcal{R}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t}))) + \frac{\beta}{M} \sum_{i=1:M} (\mathcal{L}_{c}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{i})) + \lambda \mathcal{R}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{i})))$$
(14)
$$= \mathcal{L}_{c}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})) + \gamma \mathcal{L}_{WPL}(\mathcal{W}_{\mathcal{A}}(\boldsymbol{\alpha}^{t})).$$

Lemma 1 demonstrates that the NSAS loss function not only contains the WPL but also optimizes learning of the current architecture when the appropriate constraints have been selected. Additionally, when a specific number of constraints for a densely-connected search space are set, the strategy of selecting constraints based on maximizing diversity has the potential to see the two assumptions hold true. Take the search space of NAS-Bench-201 [17] (as defined in Section 4.2) as an example. When the number of candidate operations in each edge M = 5 and the diversity of the constraint subset is maximized, those five constraint architectures should cover all possible operations for all edges (Assumption 1), and each edge in each constraint architecture should contain different operations (Assumption 2).

Algorithm 2. One-Shot NAS-NSAS

Input: \mathcal{D}_{train} , \mathcal{D}_{val} , \mathcal{W} , constraints archive $\mathcal{M} = \emptyset$, M, neural architecture search iteration T, batch size b

for $t = 1, 2, \ldots, (T * \text{size}(\mathcal{D}_{train})/b)$ do

2: if $size(\mathcal{M}) < M$ then

sample α^t based on gradient search or random search, and update the weights $W_A(\alpha)$ by normal loss function, and add architecture α into \mathcal{M} ;

4: **else**

sample α^t based on gradient search or random search, select the architecture α^m that is most similar to α^t from \mathcal{M} , and replace α^m with α^t to maximize the diversity of \mathcal{M} based on Algorithm 1. Update the weights $\mathcal{W}_A(\alpha)$ by our proposed loss function in Eq. (10) or a replay buffer;

6: **end if**

end for

3.5 One-Shot NAS With Novelty Search Based Architecture Selection

We implemented our loss function into two popular oneshot NAS: RandomNAS [30] and GDAS [16]. Like the most weight sharing NAS methods, only a single path is trained in each step of the architecture search phase. Therefore, incorporating NSAS into RandomNAS is relatively easy. However, most gradient-based NAS methods, like DARTS [38], train the whole supernet in each step of the supernet training, which would violate both the assumptions set out above. For this reason, we chose GDAS [16] as the gradient method, which uses the Gumbel-Max trick [24], [41], [55] to relax the discrete architecture distribution so as to be continuous and differentiable. The *argmax* function reparameterizes the architecture distribution and samples only one architecture in each supernet training step during the forward pass. The *softmax* function is applied during the backward pass for architecture learning. Algorithm 2 outlines the one-shot NAS with the *NSAS* loss function, called one-shot NAS-NSAS.

4 EXPERIMENTAL SETTINGS

To evaluate the performance of NSAS loss function, we compared baseline versions of RandomNAS [30] and GDAS [16] with our NSAS implementations denoted as RandomNAS-NSAS and GDAS-NSAS. We considered two different search spaces: a common search space adopted by most state-ofthe-art one-shot NAS methods [30], [38], and a second NAS-Bench-201 space [17]. The NAS-Bench-201 dataset was specifically designed for one-shot NAS research, so it comes with a guarantee of fair comparison between one-shot NAS methods. The NAS-Bench-201 search space is much smaller than the common search space, and, accordingly, the best test performances for all candidate architectures on all datasets were reported with this search space, relieving the computational concern in the further analysis of one-shot NAS approaches. We first apply RandomNAS-NSAS and GDAS-NSAS to search for promising neural architectures in the common search space and compared the results with the two baselines and many other current one-shot NAS algorithms. We then further analyzed NSAS loss function with the NAS-Bench-201 benchmark dataset. Fig. 3 illustrates the differences between the two different search spaces. In the next two subsections, we describe the experimental settings for each of these search spaces.

4.1 One-Shot NAS Common Search Space

The design on the search space is important with NAS, and a common search space [38] is typically regarded as best for a fair comparison. The cell structure in this space contains eight different types of operations: 3×3 max pooling and average pooling operation, 3×3 and 5×5 separable convolution operation, 3×3 and 5×5 dilated separable convolutions operation, identity, and zero. There are seven nodes in each cell: two input nodes, four operation nodes, and one output node. The inputs to a cell are the outputs of two of its former cells, and the output of the cell is the sum of the outputs of all operation nodes. In a CNN structure, there are two types of cells with the same search space: a normal cell α_{normal} and a reduction cell α_{reduce} . The reduction cells are located in the 1/3 and 2/3 depths of the network as residual blocks. The cell structure is repeatedly stacked to form the final CNN structure.

The number of cells for the CNN supernet training was set to only 8. We used the momentum SGD optimizer for supernet to learn the weights, and an Adam optimizer to optimize the architecture parameters. The initial learning

^{8:} Obtain α^* based on Eq. (3) (RandomNAS-NASA) or Eq. (4) (GDAS-NASA).



Fig. 3. The search spaces for the two different datasets. (a) The cell structure of a CNN with a common search space takes the two previous cells' output as the cell inputs. There are four operation nodes in each cell, and each operation node selects two outputs from the former nodes associated with those operations as inputs. The output of this cell is the sum of the outputs of all operation nodes. (b) The cell in NAS-Bench-201 is a densely-connected structure, where the operation nodes and output nodes select all former nodes7 with the applied operations as their inputs.

rate for SGD was set to 0.025 with a cosine schedule to annealed down to 0. The momentum of SGD and the weight decay were set to 0.9 and 3×10^{-4} , respectively. The initial learning rate for Adam was set to 3×10^{-4} , and the momentum and weight decay were set to (0.5,0.999) and 10^{-3} , respectively. The dropout probability was 0.4 with 16 initial channels. The CIFAR-10 dataset was used as the training set, divided into two halves at the architecture search stage. One half was used for the supernet weight learning, and the other for the architecture parameter optimization. Training took place with a batch size of 64 to derive the most promising cell structure. After supernet training and selecting the promising cells, we stacked 20 cells for full training with a batch size of 96.

The best-found cell structure with CIFAR-10 was then transferred to CIFAR-100 with the same hyperparameter settings as with the CIFAR-10 experiments. We also transferred the the best-found cell structure to ImageNet following the mobile settings in [30], [38], [55], and restricted the number of FLOPs to less than 600M. The weight decay is 3×10^{-5} , and the initial SGD learning rate is 0.1 with a decayed factor of 0.97. The network is stacked by 14 cells with batch size 128 with 250 epochs training.

As described in the previous Section 3.2 and outlined in Algorithm 1, one-shot NAS-NSAS selects M architectures from |C| previously visited architectures based on Algorithm 1. We set M = 8 and |C| = 50 for the common CNN search space. The trade-off in Eq. (10) is set as $\beta = 0.5$.

4.2 NAS-Bench-201

The NAS-Bench-201 search space is similar to the recent cell-based NAS methods [30], [38], which repeatedly stack computational cells to form the final structure. The architectural skeleton of this search space contains three stages connected by a basic residual block [21]] with a stride of 2 between them. In each stage, the cell structure was stacked N = 5 times. In this search space, the cell structure is represented as a *densely* connected directed acyclic graph (DAG) with four nodes. There are six different edges between these nodes, and each edge is associated with five candidate operations, resulting in $5^6 = 15625$ candidate cell structures. The candidate operations included: a 1×1 convolution, 3×3 convolution, 3×3 average pooling, skip connection, and zero. The zero helps to drop the associated edge. The initial channel *c* for the supernet was set to 16 and trained with an SDG optimizer. The learning rate was decayed from 0.025 to 0.001 with a cosine schedule, and the weight decay and the momentum were set to 0.0005 and 0.9, respectively. We followed the experimental setup in [17], and trained the supernet with a batch size of 64. After the architecture search phase and the most promising architectures in hand, we directly indexed the test performance of each architecture based on NAS-Bench-201 dataset [17] without training from scratch.

As described in Section 3.2, another important hyperparameter in our proposed method is the number of constraints in Eq. (10). Since the search space of NAS-Bench-201 is very small, we default set M = 2 in this search space, and the trade-off for the constraints regularization term as $\beta = 0.2$.

5 EXPERIMENTS AND RESULTS

Our first set of experiments was to conduct a neural architecture search on the common search space and compare with all methods, including our implementations, the baselines and a range of current and state-of-the-arts methods. The in-depth experiments with *NSAS* loss function and baselines on the NAS-Bench-201 dataset [17] followed.

5.1 Experimental Results on Common Search Space

5.1.1 Architecture Search on CIFAR-10

With this experiment, we searched for micro-cell structures in the search space and formed the final structure by stacking the cells in series. To compare the performance of one-shot NAS-NSAS with state-of-the-art NAS methods, we follow DARTS's experimental setting in [38]. We conducted the architecture search several times with different random seeds to obtain the architectures, and then retrained them to pick the best architectures based on retraining validation performance. The comparison results are provided in Table 1 and can be summarized as follows:

- Compared to RandomNAS and GDAS, Random-NAS-NSAS and GDAS-NSAS greatly improve the search results. The **NSAS** loss function decreased the test errors from 2.85 percent for RandomNAS to 2.59 percent, from 2.93 percent for GDAS to 2.75 percent, demonstrating the effectiveness of NSAS at improving the predictive ability of the supernet.
- RandomNAS-NSAS' results were competitive compared to the other NAS methods, with a 2.59 percent test error and only 489M FLOPs. This is an inspiring result to validate our strategy for overcoming multimodel forgetting.

TABLE 1 Results With the Existing NAS Approaches on CIFAR-10 and CIFAR-100

Method	Test E	rror (%)	Param.	FLOPs	Search Cost	Memory	Search
	CIFAR-10	CIFAR-100	(M)	(M)	(GPU Days)	Consumption	Method
NASNet-A [71]	2.65	17.81	3.3	-	1800	Single path	RL
AmoebaNet-A [46]	$3.34{\pm}0.06$	-	3.2	-	3150	Single path	EA
Hierarchical Evo [37]	3.75 ± 0.12	-	15.7	-	300	Single path	EA
PNAS [36]	$3.41 {\pm} 0.09$	17.63	3.2	-	225	P Single path	SMBO
IRLAS [19]	2.60	-	3.91	-	-	Single path	RL
IRLAS-differential [19]	2.71	-	3.43	-	-	Single path	RL
NAO [39]	3.18	-	10.6	-	1000	Single path	Gradient
NAO-WS [39]	3.53	-	2.5	-	-	Single path	Gradient
SETN (T=1K) [15]	2.69	17.25	4.6	606	1.8	Single path	Gradient
ENAS [45]	2.89	18.91	4.6	-	0.5	Single path	RL
SNAS [55]	$2.85 {\pm} 0.02$	20.09	2.8	422	1.5	Whole Supernet	Gradient
PARSEC [8]	$2.86 {\pm} 0.06$	-	3.6	485	0.6	Single path	Gradient
BayesNAS [70]	$2.81{\pm}0.04$	-	3.4	-	0.2	Whole Supernet	Gradient
RENAS [12]	$2.88 {\pm} 0.02$	-	3.5	-	6	-	RL&EA
MdeNAS [69]	2.55	17.61	3.8	599	0.16	Single path	MDL
DSO-NAS [67]	$2.84{\pm}0.07$	-	3.0	-	1	Whole Superne	Gradient
WPL [5]	3.81	-	-	-	-	Single path	RL
XNAS [43]	2.57±0.09*	16.34	3.7	596	0.3	-	Gradient
PDARTS [11]	2.50	16.63	3.4	532	0.3	-	Gradient
PC-DARTS [56]	$2.57{\pm}0.07$	17.11	3.6	557	0.3	-	Gradient
Random baseline [38]	$3.29 {\pm} 0.15$	-	3.2	-	4	-	Random
DARTS (1st) [38]	2.94	-	2.9	501	1.5	Whole Supernet	Gradient
DARTS (2nd) [38]	$2.76 {\pm} 0.09$	17.54	3.4	528	4	Whole Supernet	Gradient
GDAS [16]	2.93	18.38	3.4	519	0.21	Single path	Gradient
GDAS-NSAS	$2.75 {\pm} 0.08$	18.02 ± 0.05	3.5	528	0.4	Single path	Gradient
GDAS-NSAS-C	2.70 ± 0.07	$16.70{\pm}0.08$	3.3	520	0.4	Single path	Gradient
RandomNAS [30]	2.85 ± 0.08	17.63	4.3	612	2.7	Single path	Random
RandomNAS-NSAS	$2.59{\pm}0.06$	$17.56 {\pm} 0.05$	3.1	489	0.7	Single path	Random
RandomNAS-NSAS-C	$2.65{\pm}0.05$	$16.69{\pm}0.06$	3.5	552	0.7	Single path	Random

The first block contains the NAS methods without weight sharing. The approaches in the second block are the one-shot NAS methods. "*" indicates the results were reproduced with the best-reported cell structures in the original paper but with the same experimental settings as all the other comparators. Methods with "-" in the CIFAR-100 experiment were not reproduced because either they had different search spaces or did not report their best structures. All models were trained for 600 epochs, and we trained our best-searched architecture with 3 different random seeds to get the statistical results. "P Single path" means that the search space progressively increases during the architecture search, while only a single path is trained at each step of supernet training.

 Our NSAS evaluate more architectures during supernet training, so the search cost is slightly higher than the baselines. However, it still efficient in the sense that the supernet training in RandomNAS-NSAS only took 0.7 GPU days for and only 0.4 GPU days for GDAS-NSAS.

5.1.2 Convolutional Cell Search With Depth Constraint to Improve Transferability

In the next experiments, we transferred the best-found architectures from CIFAR-10 to CIFAR-100 and ImageNet datasets to evaluat their transferability. The results on CIFAR-100 are reported in Table 1 and ImageNet are reported in Table 2.

From Tables 1 and 2, we can see that the *NSAS* loss function improves the performance of RandomNAS and GDAS with both datasets. However, although the *NSAS* methods yielded remarkable performance with CIFAR-10, the performance was not as impressive with CIFAR-100 and ImageNet. For example, **NSAS** decrease RandomNAS' test error from 2.85 to 2.59 percent on CIFAR-10, but only from 17.63 to 17.56 percent with CIFAR-100. Similarly, with ImageNet, the improvement was only 27.1 to 26.1 percent. More importantly though, the architectures returned by RandomNAS-NSAS on

CIFAR-10 were competitive, while XNAS [43] was the superior method with CIFAR100 and ImageNet, thus demonstrating better transferability. XNAS [43] suggests that the architectures with "deeper" cell structures should provide superior performance with the ImageNet dataset. The authors also observe that most NAS methods usually return shallower cells with a larger width after searching CIFAR-10, noting that a visualization of the CNN models found from all one-shot baselines is provided in Appendix 2, which can be found on the Computer Society Digital Library at http://doi. ieeecomputersociety.org/10.1109/TPAMI.2020.3035351. For example, the architectures found by PC-DARTS and NSAS on CIFAR-10 are extremely shallow, which gave excellent results with CIFAR-10 but poor results with ImageNet. Conversely, the architectures found by XNAS, PDARTS, and PC-DARTS on ImageNet were much deeper and the results were impressive. A recent study on neural network optimization [49] gives a hint as to why most NAS methods prefer wider networks. The authors observe that width is a key factor affecting the convergence speed of neural networks, and therefore wider networks are easier to train. Based on this observation, the wider (shallower) architecture in the NAS search space reduces the loss with limited supernet training epochs, and has a higher probability of being chosen.

TABLE 2 Results With Existing Manual-Designed Architectures and NAS Approaches on the ImageNet Dataset

Method	Test Error (%)	Param.	FLOPs
	ImageNet	(M)	(M)
Inception-v1 [52]	30.2	6.6	1448
MobileNet [23]	29.4	4.2	569
MobileNet V2 [10], [48]	25.3	6.9	585
ShuffleNet $2 \times (V1)$ [68]	26.4	5	524
ShuffleNet $2 \times (V2)$ [40]	25.1	5	591
NASNet-A [72]	26.0	5.3	564
AmoebaNet-A [46]	25.5	5.1	555
PNAS [36]	25.8	5.1	588
SNAS [55]	27.3	4.3	522
SETN [15]	25.7	5.4	599
PARSEC [8]	26.3	5.5	-
BayesNAS [70]	26.5	3.9	-
MdeNAS [69]	26.8	6.1	595
DSO-NAS [67]	26.2	4.7	571
PDARTS [11]	25.9*	4.9	557
XNAS [43]	25.3*(24.7 [†])	5.3	590
PC-DARTS [56]	25.7* (25.1†)	5.3	586
DARTS (2nd) [38]	26.7	4.7	574
GDAS [16]	27.5	4.4	497
GDAS-NSAS	26.7	5.1	564
GDAS-NSAS-C	25.9	5.2	565
RandomNAS [30]	27.1	5.4	595
RandomNAS-NSAS	26.1	5.2	581
RandomNAS-NSAS-C	25.5 (24.65†)	5.4	593

The first block contains manually-designed architectures and the NAS methods without weight sharing. The second block contains one-shot NAS methods. We trained RandomNAS-NSAS with 52 initial channels C and RandomNAS-NSAS-C with C = 50. GDAS-NSAS and GDAS-NSAS-C were set to C = 50, and the FLOPs were restricted to less than 600M. * indicates that the architecture evaluation is reproduced following the common DARTS [38] setting, same as remaining methods. † indicates that the architecture evaluation settings are following PC-DARTS [56], with a warm-up linear learning rate scheduler.

These findings suggest that encouraging NAS methods to search for "deeper" architectures could improve transferability; hence, our variant of NSAS-C, NSAS with depth constraint. The depth constraint in NSAS-C force NAS to search for "deeper" architectures. Simply put, the structure of the architectures are "fixed" to a depth so that the inputs of each node are the outputs of its previous node and the output of the previous cell. Figs. 4c and 4d show an example. This way, we only need to determine the operation of each edge in an architecture. All remaining experimental settings stay the same. Tables 1 and 2 report the transferability of the models found by NSAS-C with CIFAR-100 and Image-Net. Compared to NSAS, the results are excellent. The best found cells by NSAS-C are shown in Figs. 4c and 4d, with competitive performance of 2.69, 16.58, and 25.5 percent test errors on CIFAR-10, CIFAR-100 and ImageNet, respectively. The codes and trained models are available online.¹ From Table 1, we can see that restricting the architecture depth is a very effective way of improving the transferability of NAS methods, e.g., 16.75 percent for GDAS-NSAS compared to 18.02 percent for GDAS with CIFAR-100, and RandomNAS-NSAS from 17.56 to 16.69 percent. Similarly, as shown in the ImageNet results in Table 2, NSAS-C again improves



Fig. 4. The best found cells with NSAS and NSAS-C on CIFAR-10.

transferability, with improving GDAS-NSAS from 26.7 to 25.9 percent, and RandomNAS-NSAS from 26.2 to 25.5 percent.

5.1.3 Supernet Predictive Ability Comparison

Multi-Model Forgetting in One-Shot NAS. To demonstrate catastrophic forgetting in a neural architecture search, we conducted experiments with a convolutional cell search task. The results show the differences between weight sharing and a retraining-based architecture ranking strategy. We tracked the validation accuracy of inheriting weights for several fixed sampled architectures with GDAS and also plotted the validation accuracy over 100 epochs when retraining these separate architectures from scratch in Fig. 1. From the results, we find that the validation accuracy of the architectures that directly inherit weights from the supernet fluctuate tremendously, making it hard to verify the quality of the architecture. What is worse is that the architecture ranking results completely violate the primary hypothesis of weight sharing NAS, i.e., that architectures with higher validation performance based on weight sharing should yield better retraining performance. It is worth noting that the performance of the architectures that inherited weights gets even worse during the supernet training, as shown in Fig. 1.

We also tracked the validation accuracy of weight sharing and retraining during the supernet training with RandomNAS-NSAS and GDAS-NSAS. The results are given in Fig. 5. We find that the NSAS loss function substantially alleviates multi-model forgetting with one-shot NAS. The plots of the validation accuracy with the inherited weight methods are much smoother, especially for architectures 2,



Fig. 5. The validation accuracy during supernet training for four different architectures with RandomNAS-NSAS and GDAS-NSAS. The solid lines ("Arch") indicate the validation accuracy with weights inherited from the supernet, and the dashed lines ("Arch-R") represent the validation accuracy after retraining.

3, and 4. Moreover, performance does not decrease during supernet training. This is clearly a more reliable method.

Supernet Predictive Ability Comparison. RandomNAS-NSAS and GDAS-NSAS should also alleviate ranking errors. The experiments we conducted to verify the architecture ranking predictions are shown in Figs. 6 and 7. For these experiments, we sampled four of the best architectures over four rounds with RandomNAS and RandomNAS-NSAS, and four randomly sampled from the previous experiment. Then we individually trained these 12 architectures from scratch and calculated the correlation between the architecture ranking and the validation accuracy for each of the weight sharing and retraining approaches. Fig. 6 presents the Kendall Tau (τ) metric [25], [69] of the architecture rankings based on weight sharing and retraining. The results show the difference in rankings between the normal crossentropy loss function and the NSAS loss function. Fig. 7a gives the final Kendall Tau (τ) metric values for RandomNAS and GDAS with different loss functions after supernet training. Here, the normal loss function has poor supernet predictive ability, with only $\tau = 0.0909$ and $\tau = -0.1818$ for RandomNAS and GDAS, respectively. Although the supernet trained with the NSAS loss function was not able to provide identical architecture rankings, the positive correlations matched the Kendall Tau metrics ($\tau = 0.4242$ and $\tau = 0.3030$ for RandomNAS-NSAS and GDAS-NSAS, respectively). From this we surmise that a supernet with better predictive ability tends to provide architectures with better retraining performance. Fig. 7b plots the mean retraining validation accuracy of the sampled architectures with various methods. We found that RandomNAS-NSAS achieved better results than RandomNAS, further verifying its effectiveness.

5.2 Experimental Results on NAS-Bench-201

Evaluating architectures in one-shot NAS is much more computationally intensive than an architecture search, so most



Fig. 6. The Kendall Tau metric (τ) of architecture ranking based on weight sharing and retraining.



(a) Architectures ranking difference af- (b) Retraining validation accuracy ter supernet training



state-of-the-art one-shot NAS methods only report the results of their best-found architectures. Comprehensive statistical analyses of the results are usually also overlooked due to computational limitations. Several concurrent studies [17], [27], [60], [62] have tried to address this problem by building benchmark datasets for NAS. With these datasets, researchers can analyze their one-shot NAS methods without evaluating numerous architectures. To analyze our approach in this way, we chose NASBench-201 [17] as a benchmark evaluation set. NAS-Bench-201 is easy to use and can be directly applied to most one-shot NAS algorithms. It also reports the performance of all candidate architectures on CIFAR-10, CIFAR-100, and ImageNet, making it sufficient to evaluate one-shot NAS algorithms. We did not restrict the width of architectures in the NAS-Bench-201 search space because the architectures are densely connected and have the same depth, making that constraint somewhat moot. To verify and further analyze the effectiveness of the NSAS loss function, we conducted three sets of experiments with this search space: 1) a comparison of the baselines; 2) a study of the hyperparameter settings; and 3) a study of the constraint selection strategies.

5.2.1 Empirical Comparison With Baselines

The results of the comparison study are presented in Table 3, and all experimental settings follow [17]. The statistical results were calculated from independent searches with four different *random seeds*. We found the NSAS loss function significantly improved the performance of the two baselines. RandomNAS-NSAS, in particular, achieved a test accuracy of 92.61% \pm 0.10 on CIFAR-10 compared to RandomNAS at only 88.14% \pm 0.21. Similarly, GDAS-NSAS yielded a test accuracy of 93.40% \pm 0.49 of GDAS. Furthermore, the architectures searched by RandomNAS-NSAS and GDAS-NSAS also performed better when transferred to the larger CIFAR-100 and ImageNet datasets.

5.2.2 Hyperparameter Study

As described in Eq. (10), the trade-off β and the number of constraints M are important hyperparameters for the NSAS loss function \mathcal{L}_N . we studied the impact of β concurrent with M.

First, we considered to fix the number of constraints, and investigated the impact of four different settings of β on

TABLE 3 Results of One-Shot NAS Baselines on NAS-Bench-201

Method	CIFA	AR-10	CIFA	R-100	ImageNet-16-120	
	Valid(%)	Test(%)	Valid(%)	Test(%)	Valid(%)	Test(%)
ENAS [45] DARTS (1st) [38] DARTS (2nd) [38] SETN [15] RandomNAS [30]	37.51±3.19 39.77±0.00 39.77±0.00 84.04±0.28 80.42±3.58 90.00±0.21	53.89 ± 0.58 54.30 ± 0.00 54.30 ± 0.00 87.64 ± 0.00 84.07 ± 3.61 92.51 ± 0.12	$13.37\pm2.35 \\ 15.03\pm0.00 \\ 15.03\pm0.00 \\ 58.86\pm0.06 \\ 52.12\pm5.55 \\ 71.14\pm0.27 \\ 1.14\pm0.27 $	$13.96\pm2.33 \\ 15.61\pm0.00 \\ 15.61\pm0.00 \\ 59.05\pm0.24 \\ 52.31\pm5.77 \\ 70.61\pm0.26 \\ 10.00 $	15.06 ± 1.95 16.43 ± 0.00 16.43 ± 0.00 33.06 ± 0.02 27.22 ± 3.24 41.70 ± 1.26	$14.84\pm2.10\\16.32\pm0.00\\16.32\pm0.00\\32.52\pm0.21\\26.28\pm3.09\\41.84\pm0.00$
RandomNAS* [30] RandomNAS-NSAS GDAS* [16] GDAS-NSAS	$\begin{array}{r} 85.30 \pm 0.21 \\ \hline 85.30 \pm 0.59 \\ \hline 89.20 \pm 0.31 \\ \hline 89.88 \pm 0.33 \\ \hline 89.99 \pm 0.29 \end{array}$	$\begin{array}{r} 88.14 \pm 0.21 \\ \textbf{92.61} \pm \textbf{0.10} \\ 93.40 \pm 0.49 \\ \textbf{93.55} \pm \textbf{0.16} \end{array}$	$\begin{array}{c} 62.60 \pm 3.56 \\ \textbf{68.62} \pm \textbf{1.94} \\ 70.95 \pm 0.78 \\ \textbf{71.17} \pm \textbf{0.44} \end{array}$	$\begin{array}{r} 63.40 \pm 4.52 \\ \textbf{68.47} \pm \textbf{1.73} \\ 70.33 \pm 0.87 \\ \textbf{70.69} \pm \textbf{0.33} \end{array}$	$\begin{array}{r} 33.60 \pm 4.36 \\ \textbf{41.17} \pm \textbf{2.16} \\ 41.28 \pm 0.46 \\ \textbf{41.85} \pm \textbf{1.71} \end{array}$	$\begin{array}{r} 41.84 \pm 0.09 \\ \hline 33.83 \pm 3.17 \\ \textbf{41.68} \pm \textbf{1.91} \\ 41.47 \pm 0.21 \\ \textbf{42.14} \pm \textbf{1.40} \end{array}$

"*" indicates that we reproduce the results with same random seeds as our approaches. All results in the first block are from [17]. The hyperparameters M and β were set to 5 and 0.5 for RandomNAS-NSA, 2 and 0.2 for GDAS-NSAS. We run each scenario for 4 independent times with random seed { 0, 1, 100, 101 } following the experimental settings in [17].

TABLE 4 Analysis of One-Shot NAS With Various Settings for β and M on the NAS-Bench-201 Dataset

Method	β	CIFAR-10		CIFAR-100		ImageNet-16-120	
	r	Valid Acc(%)	Test Acc(%)	Valid Acc(%)	Test Acc(%)	Valid Acc(%)	Test Acc(%)
	0	85.30±0.59	88.14±0.21	62.60±3.56	$63.40 {\pm} 4.52$	33.60 ± 4.36	33.83±3.17
RandomNAS-NSAS	0.2	86.38 ± 4.35	89.58 ± 2.82	63.64 ± 6.36	64.72 ± 4.47	34.68 ± 7.80	34.19 ± 5.72
(M=2)	0.5	85.13 ± 1.43	88.09 ± 1.23	58.73 ± 6.82	60.77 ± 3.85	31.67 ± 3.58	$30.35 {\pm} 4.18$
	0.8	86.82 ± 2.44	$90.14{\pm}2.35$	64.41 ± 4.34	64.27 ± 3.26	$35.06 {\pm} 4.81$	$34.82{\pm}6.06$
	0	$85.30 {\pm} 0.59$	$88.14{\pm}0.21$	62.60 ± 3.56	$63.40{\pm}4.52$	33.60 ± 4.36	33.83±3.17
RandomNAS-NSAS	0.2	$88.38 {\pm} 4.35$	90.74 ± 1.31	63.64 ± 6.36	$64.83 {\pm} 4.05$	36.53 ± 6.50	$36.08 {\pm} 4.68$
(M=3)	0.5	87.93±1.63	91.23±1.03	66.03 ± 3.46	$66.17 {\pm} 4.12$	$38.14{\pm}2.76$	38.72 ± 3.11
	0.8	$85.58 {\pm} 2.59$	88.78 ± 2.23	$64.12 {\pm} 4.55$	65.06 ± 3.35	$34.80{\pm}4.24$	$34.37 {\pm} 5.57$
	0	85.30±0.59	88.14±0.21	62.60±3.56	$63.40{\pm}4.52$	33.60±4.36	33.83±3.17
RandomNAS-NSAS	0.2	86.73±1.30	$90.64 {\pm} 0.99$	$62.38 {\pm} 4.57$	66.42 ± 1.47	36.68 ± 3.80	37.59 ± 3.72
(M = 4)	0.5	87.13 ± 1.43	$91.04{\pm}0.43$	$64.43 {\pm} 4.82$	64.77 ± 3.61	36.86 ± 3.70	$36.35 {\pm} 4.15$
	0.8	$88.52 {\pm} 0.74$	$92.04{\pm}0.50$	$67.40 {\pm} 2.22$	$67.62 {\pm} 1.94$	$39.91 {\pm} 4.50$	$40.61 {\pm} 3.51$
	0	85.30±0.59	88.14±0.21	62.60±3.56	$63.40{\pm}4.52$	33.60±4.36	33.83±3.17
RandomNAS-NSAS	0.2	$88.45 {\pm} 0.47$	$91.36 {\pm} 0.74$	65.79 ± 0.59	$65.58 {\pm} 0.42$	37.69±1.23	37.31 ± 2.42
(M = 5)	0.5	89.20±0.31	92.61±0.10	68.62±1.94	68.47±1.73	$41.17{\pm}2.16$	$41.68{\pm}1.91$
	0.8	88.42 ± 0.30	$91.56{\pm}0.36$	$66.77 {\pm} 4.83$	$66.58 {\pm} 4.99$	$39.53 {\pm} 4.85$	$38.64{\pm}5.13$

multi-model forgetting with one-shot NAS. In this experiment, we took the RandomNAS-NSAS with M = 5 as example. The results, shown in the last four rows of Table 4, indicate that using constraints to regularize the supernet training can greatly improve test performance and, additionally, that RandomNAS-NSAS is somewhat sensitive to β . More important, with different number of constraints (M = 2, 3, 4), the results all demonstrated our regularization method can enhance the performance, where our Random-NAS-NSAS with different M and β all outperformed the baseline.

We then fixed $\beta = 0.2$ and varied M, also to analyzed the impact on forgetting. Given there are only five candidate operations in the NAS-Bench-201 search space, there are no shared weights among constraints only when $M \leq 5$. The four settings for M in this experiment were 2, 3, 4, and 5. From the results, we found that, again, RandomNAS-NSAS seemed sensitive to the number of constraints, and the larger M = 5 gave much better results than the other scenarios for RandomNAS-NSAS with CIFAR10, CIFAR-100, and

ImageNet. More interestingly, there was a large performance gain between M = 4 and M = 5 with RandomNAS-NSAS. One underlying reason may be that M = 5 has the potential to ensure the two assumptions hold true, as discussed in Section 3.3. Similarly, the tendency also exists in the remain 2 different β , that increasing the number of constrained architectures can enhance the performance. More details on these results can be found in Table 4.

Overall, Table 4 considered three settings for β and four settings for M, and presented the results of RadnomNAS-NSAS on CIFAR-10, CIFAR-100, and ImageNet-16-120. In general, a larger β and a larger M provided the better results. In the next subsection, we discuss the benefits of holding to the two assumptions with NSAS, with respect to the constrained architecture selection strategy.

5.2.3 Analysis of Constraints Selection

Although we demonstrate the theoretical benefits of the NSAS loss function in relieving catastrophic forgetting

TABLE 5 Analysis of the One-Shot NAS Methods With Various Constraint Selection Strategies on CIFAR-10

Method	β	NSAS	NSAS-G	NSAS-LG	RG	LoW	LoW-R
		Test Acc(%)	Test Acc (%)	Test Acc(%)	Test Acc(%)	Test Acc(%)	Test Acc (%)
RandomNAS Test Acc(%) (88.14±0.21)	0.2 0.5 0.8	$\begin{array}{c} 89.58{\pm}2.82\\ 88.09{\pm}1.23\\ 90.14{\pm}2.35\end{array}$	91.74±1.16 90.37±2.14 92.52±0.48	91.11±2.16 91.19±0.79 91.18±1.73	$\begin{array}{c} 89.24{\pm}2.24\\ 77.30{\pm}19.76\\ 89.53{\pm}0.24\end{array}$	90.13±4.25 89.86±2.91 89.39±3.88	$\begin{array}{c} 89.72{\pm}3.64\\ 86.85{\pm}10.29\\ 85.54{\pm}7.18\end{array}$
GDAS Test Acc(%) (93.40±0.49)	0.2 0.5 0.8	93.55±0.16 93.49±0.24 93.32±0.18	93.37±0.27 93.51±0.14 93.29±0.29	93.52±0.30 93.46±0.27 93.55±0.20	93.29 ± 0.19 93.31 ± 0.30 93.40 ± 0.28	93.51±0.14 93.47±0.29 93.55±0.20	93.40±0.26 93.36±0.21 93.55±0.16

TABLE 6

Analysis of the One-Shot NAS Methods With Various Constraint Selection Strategies on CIFAR-100

Method	β	NSAS	NSAS-G	NSAS-LG	RG	LoW	LoW-R
		Test Acc(%)	Test Acc (%)	Test Acc(%)	Test Acc(%)	Test Acc(%)	Test Acc (%)
RANDOMNAS Test Acc(%) (63.40±4.52)	0.2 0.5 0.8	64.72±4.47 60.77±3.85 64.27±3.26	67.22±2.20 65.30±3.49 67.83±1.66	66.58±3.43 66.74±2.66 66.01±2.94	63.66 ± 3.97 47.28 ± 27.17 64.13 ± 0.56	$\begin{array}{c} 64.06{\pm}7.89\\ 64.27{\pm}6.43\\ 61.37{\pm}8.95\end{array}$	60.68 ± 9.04 59.47 ± 13.13 58.32 ± 8.82
GDAS Test Acc(%) (70.33±0.87)	0.2 0.5 0.8	70.69 ± 0.33 70.53 ± 0.30 70.33 ± 0.41	70.49 ± 0.61 70.57 ± 0.23 70.28 ± 0.58	$70.86 {\pm} 0.90$ $70.69 {\pm} 0.67$ $70.80 {\pm} 0.55$	$70.43 {\pm} 0.56 \\ 70.21 {\pm} 0.44 \\ 70.40 {\pm} 0.60$	70.53 ± 0.34 70.10 ± 0.70 70.78 ± 0.19	$70.40 {\pm} 0.51$ $70.25 {\pm} 0.38$ $70.38 {\pm} 0.45$

 TABLE 7

 Analysis of the One-Shot NAS Methods With Various Constraint Selection Strategies on ImageNet-16-120

Method	β	NSAS	NSAS-G	NSAS-LG	RG	LoW	LoW-R
		Test Acc(%)	Test Acc (%)	Test Acc(%)	Test Acc(%)	Test Acc(%)	Test Acc (%)
RandomNAS Test Acc(%) (33.83±3.17)	0.2 0.5 0.8	34.19 ± 5.72 30.35 ± 4.18 34.82 ± 6.06	39.58±2.60 35.14±3.33 40.14±2.60	37.79±5.11 39.81±2.81 38.34±4.65	34.38 ± 5.02 31.96 ± 26.53 34.94 ± 2.29	36.32±8.07 36.81±5.47 33.75±7.91	30.64 ± 13.19 29.11 ± 17.77 28.38 ± 9.84
GDAS Test Acc(%) (41.47±0.21)	0.2 0.5 0.8	$\begin{array}{c} 42.14{\pm}1.40\\ 42.20{\pm}1.31\\ 41.78{\pm}0.89\end{array}$	42.26±0.20 42.16±1.30 42.21±0.16	41.71±0.57 42.29±1.00 41.64±1.01	41.35 ± 0.13 42.45 ± 1.07 41.68 ± 0.96	42.16 ± 1.30 41.32 ± 1.56 41.84 ± 1.01	$\begin{array}{c} 41.68{\pm}1.18\\ 42.20{\pm}1.25\\ 41.64{\pm}1.21\end{array}$

Section 3.3, and the experiments in Sections 5.2.1 and 5.2.2 support these theories, at least for one-shot NAS, is it still open to debate as to whether these improvements are due to the constraint selection strategy or simply because of the regularization. In this section, we further conduct an ablation study to investigate the impact of different architecture selection strategies.

Directly maximizing the diversity of the constraint subset, as per Section 3.4, easily holds Assumption 2, but it does not guarantee that Assumption 1 will hold. Therefore, we devised two variants of the NSAS loss function, both of which strictly observe the assumptions when selecting constraints. These are NSAS-G and NSAS-LG. The difference between the two concerns treatment of the last architecture. More specifically, with NSAS-G, the constraints are generated randomly, maximizing diversity, but the last constraint θ_M is generated by complementing the operations contained in the current architecture α^{t} that have not been covered in the previous constraints. This means all selected architectures $\{\theta_1, \ldots, \theta_M\}$ covering all parameters of α_t such that $\theta_t \subseteq \{\theta_1 \cup ... \cup \theta_M\}$. With *NSAS*-*LG*, however, the last architecture α^{t-1} is first added into the subset, and remaining constraints are generated as following NSAS-G. This is to test the common thinking on catastrophic forgetting that the last architecture deteriorates performance the most.

We evaluated all three loss functions - *NSAS*, *NSAS-G*, and *NSAS-LG* - along with three naive architecture selection methods added to the RandomNAS and GDAS baseline to regularize the supernet training. Thus, the six loss functions were:

- *NSAS* which selects constraints through maximizing diversity.
- *NSAS-G* a variant of *NSAS* described above.
- *NSAS-LG* a variant of *NSAS* described above.
- *RG* randomly generates architectures to form the constraint subspace.
- *LoW* only adds the last architecture α^{t-1} to the constraint subset.
- LoW-R adds the last architecture α^{t-1} to the constraint subset plus randomly generated constraints.

Tables 5, 6, and 7 show the test accuracies for the CIFAR-10, CIFAR-100, and ImageNet-16-120 datasets. We set the number of constraints M = 2 in this experiment, to more precisely investigate the effect of architecture selection and quantitatively analyze the constraint selection strategies. The results for one-shot NAS without relieving forgetting are shown in the first column of Tables 5, 6, and 7. All NSAS methods improved performance but, interestingly, some of the naive constraint selection methods did as well, which indicates that overcoming catastrophic forgetting is a promising research direction for one-shot NAS. For example, LoW improved performance simply by including the last visited architecture in the regularization. However, these results also show the importance of constraint selection strategy, as randomly selecting constraints can reduce performance. We observed that RG was the worst method in all cases, with a reduction in test accuracy from 88.14 ± 0.21 to $77.30\% \pm$ 19.76 for RandomNAS, and from 93.40 ± 0.49 to 93.29 ± 0.19 for GDAS in CIFAR-10. Moreover, the LoW-R strategy of adding more randomly generated constraints into the replay buffer yielded even worse results than LoW in most cases. These results suggest that randomly selecting constraints does not alleviate multi-model forgetting with one-shot NAS.

As for the *NSAS* and its variants, *NSAS-G* and *NSAS-LG*, all improved performance significantly. It is interesting that *NSAS* and *NSAS-G* achieved similar results with GDAS. This indicates that Assumption 1, which requires the constraints to cover all parameters of α_t , may not be so important for relieving catastrophic forgetting with gradient-based one-shot NAS while holding to this assumption with RandomNAS did help. Overall, *NSAS-G* achieved much better results than *NSAS* and, in most cases, *NSAS-G* and *NSAS-LG* produced the best results. Thus, simultaneously considering the last visited architecture and maximizing the diversity of constraints combined are the two key factors that need to be addressed to relieve catastrophic forgetting with both random sampling-based and gradient-based one-shot NAS.

5.3 Discussion

We can draw several conclusions from this series of experiments.

- RandomNAS tends to achieve better performance than GDAS with a common search space, whereas GDAS outperforms RandomNAS with the NAS-Bench-201 space no matter the loss function. This may be because gradient-based methods typically arrive at the local optimal solution once the supernet is trained. RandomNAS, however, must perform a subsequent model selection process using either a random search or an EA to find a global optimal solution from the trained supernet. Since common search spaces are much more complicated than the one in NAS-Bench-201, a global optimization method will usually outperform a gradient method, while gradient-based NAS is more efficient and effective with simple search spaces.
- It is clear that the NSAS loss function can increase the predictive ability of the supernet, which, in turn, greatly improves the performance of the architectures found by RandomNAS and GDAS. However, supernet training in one-shot NAS is still a problem with much room for further advancements. Devising a more appropriate loss function than the status quo appears to be a promising direction for improving the performance of one-shot NAS methods.

 Lastly, the ablation study indicates that adding recently visited architectures into the constraint subset and maximizing its diversity are two efficient ways to mitigate catastrophic forgetting with oneshot NAS.

6 CONCLUSION AND FUTURE WORKS

In this paper, we formulated supernet training as a constrained optimization problem to reduce some of the negative impacts of catastrophic forgetting with one-shot NAS, and multi-model forgetting in particular. Our strategy is to select a representative subset of constraints with a greedy novelty search method. Then the supernet training is regularized in a feasible region with a new novelty search-based architecture selection loss function, i.e., *NSAS* to overcome multi-model forgetting.

We implemented *NSAS* into two one-shot NAS baselines -RandomNAS and GDAS - and compared the quality of the architecture selections with and without the new loss function. The results of experiments on the common search space of a neural architecture show *NSAS* and two of its variants improve the predictive ability of the supernet with both convolutional and recurrent cell search. Experiments with the NAS-Bench-201 dataset also suggest that *NSAS* can substantially offset performance degradation due to forgetting with one-shot NAS. In future research, we plan to focus on searching on a latent space by transforming discrete architectures into continuous representations. Further, we will look to leveraging expert knowledge with DNN searches to design architectures with greater transferability.

ACKNOWLEDGMENTS

This work was supported in part by the NSFC under Grant No. 61702415 and No. 61972315, the Australian Research Council (ARC) under a Discovery Early Career Researcher Award (DECRA) No. DE190100626, the Air Force Research Laboratory, DARPA under Agreement No. FA8750- 19-2-0501, and the Youth Innovation Promotion Association CAS (No. 2017210).

REFERENCES

- [1] G. Adam and J. Lorraine, "Understanding neural architecture search techniques," 2019, *arXiv: 1904.00438*.
- [2] R. Aljundi, M. Lin, B. Goujaud, and Y. Bengio, "Gradient based sample selection for online continual learning," in *Proc. Conf. Neural Inf. Process. Syst.*, 2019, pp. 11816–11825.
- [3] B. Baker, O. Gupta, R. Raskar, and N. Naik, "Accelerating neural architecture search using performance prediction," 2018.
- [4] G. Bender, P.-J. Kindermans, B. Zoph, V. Vasudevan, and Q. Le, "Understanding and simplifying one-shot architecture search," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 550–559.
- [5] Y. Benyahia *et al.*, "Overcoming multi-model forgetting," in *Proc.* 36th Int. Conf. Mach. Learn., 2019, pp. 594–603.
- [6] A. Brock, T. Lim, J. M. Ritchie, and N. J. Weston, "Smash: Oneshot model architecture search through hypernetworks," in *Proc. Int. Conf. Learni. Representations*, 2018.
- [7] H. Cai, L. Zhu, and S. Han, "Proxylessnas: Direct neural architecture search on target task and hardware," in *Proc. Int. Conf. Learn. Representations*, 2019.
- [8] F. P. Casale, J. Gordon, and N. Fusi, "Probabilistic neural architecture search," 2019, arXiv: 1902.05116.
- [9] X. Chang, Y. Yu, Y. Yang, and E. P. Xing, "Semantic pooling for complex event analysis in untrimmed videos," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 8, pp. 1617–1632, Aug. 2017.

- [10] S. Chen, Y. Liu, X. Gao, and Z. Han, "MobileFaceNets: Efficient CNNs for accurate real-time face verification on mobile devices," in *Proc. Chinese Conf. Biometric Recognit.*, 2018, pp. 428–438.
- [11] X. Chen, L. Xie, J. Wu, and Q. Tian, "Progressive differentiable architecture search: Bridging the depth gap between search and evaluation," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 1294–1303.
- [12] Y. Chen et al., "RENAS: Reinforced evolutionary neural architecture search," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2019, pp. 4782–4791.
- [13] X. Cheng et al., "Hierarchical neural architecture search for deep stereo matching," in Proc. Annu. Conf. Neural Inf. Process. Syst., 2020.
- [14] X. Chu, B. Zhang, R. Xu, and J. Li, "Fairnas: Rethinking evaluation fairness of weight sharing neural architecture search," 2019, arXiv: 1907.01845.
- [15] X. Dong and Y. Yang, "One-shot neural architecture search via self-evaluated template network," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 3680–3689.
- [16] X. Dong and Y. Yang, "Searching for a robust neural architecture in four GPU hours," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 1761–1770.
- [17] X. Dong and Y. Yang, "Nas-bench-201: Extending the scope of reproducible neural architecture search," in *Proc. Int. Conf. Learn. Representations*, 2020.
- [18] T. Elsken, Jan H. Metzen, and F. Hutter, "Neural architecture search: A survey," J. Mach. Learn. Rese., vol. 20, no. 55, pp. 1–21, 2019.
- [19] M. Guo, Z. Zhong, W. Wu, D. Lin, and J. Yan, "IRLAS: Inverse reinforcement learning for architecture search," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 9013–9021.
- [20] Z. Guo *et al.*, "Single path one-shot neural architecture search with uniform sampling," in *Proc. Eur. Conf. Comput. Vis.*, Springer, 2020, pp. 544–560.
- [21] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [22] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," in Proc. NIPS Deep Learn. Representation Learn. Workshop, 2015.
- [23] A. G. Howard *et al.*, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," 2017, *arXiv*: 1704.04861.
- [24] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with gumbel-softmax," in Proc. Int. Conf. Learn. Representations, 2017.
- [25] M. G. Kendall, "The treatment of ties in ranking problems," Biometrika, vol. 33, no. 3, pp. 239–251, 1945.
- [26] J. Kirkpatrick et al., "Overcoming catastrophic forgetting in neural networks," Proc. Nat. Acad. Sci. United States America, vol. 114. no. 13, pp. 3521–3526, 2017.
- no. 13, pp. 3521–3526, 2017.
 [27] A. Klein and F. Hutter, "Tabular benchmarks for joint architecture and hyperparameter optimization," 2019, *arXiv*: 1905.04970.
- [28] S.-W. Lee, J.-H. Kim, J. Jun, J.-W. Ha, and B.-T. Zhang, "Overcoming catastrophic forgetting by incremental moment matching," in *Proc.* 31st Int. Conf. Neural Inf. Process. Syst., 2017, pp. 4655–4665.
- [29] C. Li *et al.*, "Blockwisely supervised neural architecture search with knowledge distillation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 1986–1995.
 [30] L. Li and A. Talwalkar, "Random search and reproducibility for
- [30] L. Li and A. Talwalkar, "Random search and reproducibility for neural architecture search," in *Proc. Assoc. Uncertainty Artif. Intell.*, 2019, pp. 367–377.
- [31] X. Li et al., "Improving one-shot NAS by suppressing the posterior fading," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2020, pp. 13836–13845.
- [32] X. Li, Y. Zhou, T. Wu, R. Socher, and C. Xiong, "Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 3925–3934.
- [33] Z. Li and D. Hoiem, "Learning without forgetting," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 12, pp. 2935–2947, Dec. 2018.
- [34] Z. Li, L. Yao, X. Chang, K. Zhan, J. Sun, and H. Zhang, "Zero-shot event detection via event-adaptive concept relevance mining," *Pattern Recognit.*, vol. 88, pp. 595–603, 2019.
- [35] D. Lian et al., "Towards fast adaptation of neural architectures with meta learning," in Proc. Int. Conf. Learn. Representations, 2020.
- [36] C. Liu *et al.*, "Progressive neural architecture search," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 19–35.
 [37] H. Liu, K. Simonyan, O. Vinyals, C. Fernando, and K. Kavukcuo-
- [37] H. Liu, K. Simonyan, O. Vinyals, C. Fernando, and K. Kavukcuoglu, "Hierarchical representations for efficient architecture search," in *Proc. Int. Conf. Learn. Representations*, 2018.
- [38] H. Liu, K. Simonyan, and Y. Yang, "Darts: Differentiable architecture search," in Proc. Int. Conf. Learn. Representations, 2019.

- [39] R. Luo, F. Tian, T. Qin, E. Chen, and T.-Y. Liu, "Neural architecture optimization," in *Proc. Advances Neural Inf. Process. Syst.*, 2018, pp. 7816–7827.
- [40] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun, "Shufflenet v2: Practical guidelines for efficient CNN architecture design," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 116–131.
- [41] C. J. Maddison, A. Mnih, and Y. W. Teh, "The concrete distribution: A continuous relaxation of discrete random variables," in *Proc. Int. Conf. Learn. Representations*, 2017.
- [42] J. Mei et al., "Atomnas: Fine-grained end-to-end neural architecture search," in Proc. Int. Conf. Learn. Representations, 2020.
- [43] N. Nayman, A. Noy, T. Ridnik, I. Friedman, R. Jin, and L. Zelnik-Manor, "XNAS: Neural architecture search with expert advice," in *Proc. Conf. Neural Inf. Process. Syst.*, 2019, pp. 1977–1987.
- [44] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," *Neural Netw.*, vol. 113, pp. 54–71, 2019.
- [45] H. Pham, M. Guan, B. Zoph, Q. Le, and J. Dean, "Efficient neural architecture search via parameter sharing," in *Proc. 35th Int. Conf. Mach. Learn.*, 2018, pp. 4095–4104.
- [46] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, "Regularized evolution for image classifier architecture search," in *Proc. AAAI Conf. Artif. Intell.*, 2019, pp. 4780–4789.
- [47] P. Řen *et al.*, "A comprehensive survey of neural architecture search: Challenges and solutions," 2020, *arXiv*: 2006.02903
- [48] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2018, pp. 4510–4520.
- [49] K. A. Sankararaman, S. De, Z. Xu, W. R. Huang, and T. Goldstein, "The impact of neural network overparameterization on gradient confusion and stochastic gradient descent," in *Proc. Int. Conf. Mach. Learn.*, 2020.
- [50] Y. Shu, W. Wang, and S. Cai, "Understanding architectures learnt by cell-based neural architecture search," in *Proc. Int. Conf. Learn. Representations*, 2020.
- [51] P. Singh, T. Jacobs, S. Nicolas, and M. Schmidt, "A study of the learning progress in neural architecture search techniques," 2019, arXiv: 1906.07590.
- [52] C. Szegedy et al., "Going deeper with convolutions," in Proc. IEEE/ CVF Conf. Comput. Vis. Pattern Recognit., 2015, pp. 1–9.
- [53] W. Wang, D. Tran, and M. Feiszli, "What makes training multimodal classification networks hard?," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 12695–12705.
- [54] C. White, W. Neiswanger, and Y. Savani, "Bananas: Bayesian optimization with neural architectures for neural architecture search," 2019, arXiv: 1910.11858.
- [55] S. Xie, H. Zheng, C. Liu, and L. Lin, "SNAS: Stochastic neural architecture search," in Proc. Int. Conf. Learn. Representations, 2019.
- [56] Y. Xu et al., "PC-DARTS: Partial channel connections for memoryefficient architecture search," in Proc. Int. Conf. Learn. Representations, 2020.
- [57] A. Yang, Pedro M Esperança, and F. M. Carlucci, "NAS evaluation is frustratingly hard," in Proc. Int. Conf. Learn. Representations, 2020.
- [58] Z. Yang et al., "Cars: Continuous evolution for efficient neural architecture search," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2020, pp. 1826–1835.
- [59] Q. Yao, J. Xu, W.-W. Tu, and Z. Zhu, "Efficient neural architecture search via proximal iterations," in Proc. AAAI Conf. Artif. Intell., 2020.
- [60] C. Ying, A. Klein, E. Christiansen, E. Real, K. Murphy, and F. Hutter, "NAS-bench-101: Towards reproducible neural architecture search," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 7105–7114.
- [61] K. Yu, C. Sciuto, M. Jaggi, C. Musat, and M. Salzmann, "Evaluating the search phase of neural architecture search," in *Proc. Int. Conf. Learn. Representations*, 2020.
- [62] A. Zela, J. Siems, and F. Hutter, "NAS-bench-1shot1: Benchmarking and dissecting one-shot neural architecture search," in *Proc. Int. Conf. Learn. Representations*, 2020.
- [63] C. Zhang, M. Ren, and R. Urtasun, "Graph hypernetworks for neural architecture search," in Proc. Int. Conf. Learn. Representations, 2019.
- [64] M. Zhang, H. Li, S. Pan, X. Chang, Z. Ge, and S. Su, "Differentiable neural architecture search in equivalent space with exploration enhancement," in *Proc. 34th Conf. Neural Inf. Process.*, 2020.
- [65] M. Zhang, H. Li, S. Pan, X. Chang, and S. Su, "Overcoming multimodel forgetting in one-shot NAS with diversity maximization," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 7806–7815.

- [66] M. Zhang, H. Li, S. Pan, T. Liu, and S. Su, "One-shot neural architecture search via novelty driven sampling," in *Proc. Int. Joint Conf. Artif. Intell.*, 2020.
- [67] X. Zhang, Z. Huang, N. Wang, S. Xiang, and C. Pan, "You only search once: Single shot neural architecture search via direct sparse optimization," *IEEE Trans. Pattern Anal. Mach. Intell.*, early access, 2020, doi: 10.1109/TPAMI.2020.3020300.
- [68] X. Zhang, X. Zhou, M. Lin, and J. Sun, "Shufflenet: An extremely efficient convolutional neural network for mobile devices," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.* 2018, pp. 6848–6856.
- [69] X. Zheng, R. Ji, L. Tang, B. Zhang, J. Liu, and Q. Tian, "Multinomial distribution learning for effective neural architecture search," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 1304–1313.
- Proc. IEEE/CVF Int. Conf. Comput. Vis., 2019, pp. 1304–1313.
 [70] H. Zhou, M. Yang, J. Wang, and W. Pan, "Bayesnas: A Bayesian approach for neural architecture search," in Proc. Int. Conf. Mach. Learn., 2019, pp. 7603–7613.
- [71] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," in Proc. Int. Conf. Learn. Representations, 2017.
- [72] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 8697–8710.



Miao Zhang received the PhD degree from the Beijing Institute of Technology (BIT), China. He is also working toward the dual PhD degree at the University of Technology Sydney (UTS). His major research interests include AutoML, neural architecture search, Bayesian optimization, continual learning, and deep learning.



Chuan Zhou received the PhD degree from the Chinese Academy of Sciences, in 2013. He won the outstanding doctoral dissertation of Chinese Academy of Sciences, in 2014, the best paper award of ICCS-14, and the Best Student Paper Award of IJCNN-17. Currently, he is an associate professor with the Academy of Mathematics and Systems Science, Chinese Academy of Sciences. His research interests include socail network analysis and graph mining. To date, he has published more than 70 papers, including the *IEEE Transactions on Knowledge and Data Engineering*, ICDM, AAAI, CIKM, IJCAI, and WWW.



Zongyuan Ge is a full-time senior research fellow employed by Monash University vice chancellor and provost office with specific interest and expertise in Medical AI development. He has a strong background in statistical analysis, machine learning and computer vision research. So far, he has published more than 40 peer-reviewed publications and patents, which are first/senior author. He was selected as one of the 200 Most Qualified Young Researchers in Computer and Mathematics by the Scientific Committee of the Heidelberg Laureate Forum Foundation in 2017 and received Monash Exceptional Research Award 2019.



Huiqi Li (Senior Member, IEEE) received the PhD degree from Nanyang Technological University, Singapore, in 2003. She is currently a professor with the Beijing Institute of Technology. Her research interests include image processing and computer-aided diagnosis.



Steven Su (Senior Member, IEEE) received the PhD degree in control engineering from RSISE the Australian National University (ANU). He is currently an associate professor with the University of Technology Sydney (UTS). His research interests include system modeling and control, machine learning, wearable health monitoring, and rehabilitation engineering.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.



Shirui Pan received the PhD degree in computer science from the University of Technology Sydney (UTS), Ultimo, NSW, Australia. He is currently a lecturer with the Faculty of Information Technology, Monash University, Australia. Prior to this, he was a lecturer with the School of Software, University of Technology Sydney. His research interests include data mining and machine learning. To date, he has published more than 80 research papers in top-tier journals and conferences.



Xiaojun Chang is currently a faculty member with the Faculty of Information Technology, Monash University, Clayton, VIC, Australia. He is also a distinguished adjunct professor with the Faculty of Computing and Information Technology, King Abdulaziz University. He is an ARC Discovery Early Career Researcher Award (DECRA) Fellow from 2019 to 2021. He has achieved top performance in various international competitions, such as TREC-VID MED, TRECVID SIN, and TRECVID AVS.