

Towards Self-supervised and Weight-preserving Neural Architecture Search

Zhuowei Li^{*1}, Yibo Gao^{*†2}, Zhenzhou Zha^{†3}, Zhiqiang Hu⁴, Qing Xia⁴, and Shaoting Zhang⁴ Dimitris N. Metaxas¹

¹ Rutgers Univeristy, USA, z1502@cs.rutgers.edu

² University of Electronic Science and Technology of China, CN

³ Zhejiang University, CN

⁴ SenseTime Research, CN

Abstract. Neural architecture search (NAS) algorithms save tremendous labor from human experts. Recent advancements further reduce the computational overhead to an affordable level. However, it is still cumbersome to deploy the NAS techniques in real-world applications due to the fussy procedures and the supervised learning paradigm. In this work, we propose the self-supervised and weight-preserving neural architecture search (SSWP-NAS) as an extension of the current NAS framework by allowing the self-supervision and retaining the concomitant weights discovered during the search stage. As such, we simplify the workflow of NAS to a one-stage and proxy-free procedure. Experiments show that the architectures searched by the proposed framework achieve state-of-the-art accuracy on CIFAR-10, CIFAR-100, and ImageNet datasets without using manual labels. Moreover, we show that employing the concomitant weights as initialization consistently outperforms the random initialization and the two-stage weight pre-training method by a clear margin under semi-supervised learning scenarios. Codes are publicly available at <https://github.com/LzVv123456/SSWP-NAS>.

Keywords: Self-supervised Learning, Neural Architecture Search, Pre-training, Image Classification

1 Introduction

The development of NAS algorithms save considerable time and efforts of human experts through automating the neural architecture design process. It has achieved state-of-the-art performances in a series of vision tasks including image recognition [49,50,32], semantic segmentation [19,46] and object detection [39,9]. Recent advances on weight-sharing NAS [30] and differentiable NAS [23,41] further reduce the searching cost from thousands of GPU-days to a couple.

Despite the significant computational reduction made by current NAS methods, it is still cumbersome to deploy the NAS techniques in real-world applications due to the fussy procedures. As shown in Fig. 1, a typical workflow of NAS

* Equal contributions.

† This work was done during the internship at SenseTime.

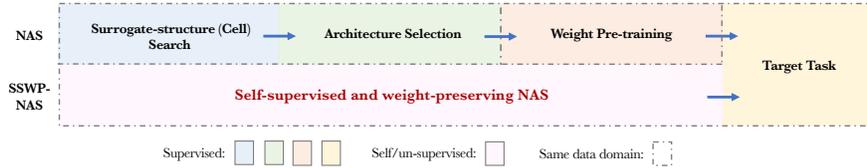


Fig. 1: Overview for the regime of general NAS and the proposed SSWP-NAS.

consists of the surrogate-structure search and the architecture selection two steps to acquire the architecture. Then a standalone procedure of weight pre-training need to be taken before transferring the architecture to downstream tasks. It is non-trivial to pre-train a network, taking even more time than the searching process. Besides, existing NAS workflows largely rely on manual annotations, making the domain-specific NAS even more unwieldy.

Driven by the inconvenience of the current NAS paradigm, we propose a new framework, namely self-supervised and weight-preserving neural architecture search (SSWP-NAS), as an extension of the current NAS methodology with the following two prominent properties: (1) SSWP-NAS is self-supervised so that it does not rely on manual signals to perform optimization. This property also removes the dependency on proxy-datasets (e.g. ImageNet [34]). (2) SSWP-NAS has the weight-preserving property, which means the concomitant weights generated during the search process can be retained and serve as initialization to benefit transfer learning. This property simplifies the current NAS workflow from the two-stage fashion to the one-stage. To achieve weight-preserving, we align the dimensionality of the network used during the search and train stage and leverage stochastic operation sampling strategy to reduce the memory footprint. To remove the dependency on manual labels, we probe how the designed searching process copes with the self-supervised learning objective. We also observe a persistent optimization challenge dubbed *network inflation issue* caused by the inconsistent optimization targets and the stochastic strategy. To overcome this challenge, we propose the forward progressive prune (FPP) operation that gradually bridges the gap between the optimization targets and reduces the extend of stochastic operations.

In experiments, the searched architecture using SSWP-NAS achieve state-of-the-art accuracy on CIFAR-10 (2.41% error rate), CIFAR-100 (16.47% error rate), and ImageNet (24.3% top-1 error rate under restricted resources) without using manual labels. Besides, concomitant weights as initialization consistently outperforms the random initialization and the two-stage weight pre-training method by a clear margin under semi-supervised learning scenarios. Moreover, we show that self-supervised learning objective consistently outperforms the supervised counterpart in our framework and FPP is beneficial under both supervised and self/un-supervised learning objectives. Comprehensive ablation studies have also been conducted towards the proposed designs. Our main contributions can be summarized in three-fold:

- We propose the SSWP-NAS that enjoys the self-supervised learning and weight-preserving property. It simplifies the current NAS workflow from the two-stage manner to the one stage. As a by-product, SSWP-NAS is also proxy-free, which means it relies on neither the surrogate structure nor the proxy-dataset.
- We propose the FPP to address the network inflation issue that occurs during the designed weight-preserving search process. We empirically show that FPP is beneficial under both supervised and self/un-supervised signals.
- SSWP-NAS searches the architecture and generate the pre-train weights concurrently while achieving state-of-the-art performance regarding the quality of both the architecture and the pre-train weights.

2 Related Work

Neural Architecture Search. Early works for NAS mainly leverage on reinforcement learning (RL) [49,50] or evolutionary algorithms (EA) [32,22] to optimize a controller that samples a sequence of discrete operations to form the architecture. This straightforward implementation consumes tremendous computational resources. As a remedy, the following works rely on the weight-sharing [30,23] methods and surrogate structures [50,23] to reduce the computational overhead. DARTS [23] further simplify the NAS framework by relaxing the search space from discrete to continuous. This *One-shot* searching framework then becomes popular in the NAS domain due to its simplicity and efficiency. The proposed SSWP-NAS also inheres to this line of research. While the original DARTS design is observed to suffer from performance degeneration [4] and mode collapse [18] issue. To this end, DARTS+ [18] introduces the early stop to suppress the over-characterized non-parameterized operations. P-DARTS [4] tries to alleviate the performance drop by gradually increasing the depth of the surrogate structure. ProxylessNAS [2] first achieves the differentiable architecture search without a surrogate structure. Despite the advancements made above methods, they still rely on supervised signals to perform optimization, and none of them are able to preserve the concomitant weights.

Self-supervised Learning. Not until recently, self-supervised learning largely relies on heuristic pretext tasks [8,27,44,48] to form the supervision signal, and their performance lags behind a lot compared with the supervised counterpart. Emerging of the contrastive-based self-supervised learning largely close the gap between the self-supervised and supervised weight pre-training [28,38,11,3,17]. It encourages a pulling force within a positive pair and pushing away negative pairs through minimizing a discriminative loss function. BYOL [10] and SimSiam [5] also demonstrate that the negative pairs are not necessary. In this work, we investigate how state-of-the-art self-supervised learning methods cope with the weight-preserving network search.

Self/un-supervised NAS. Most recently, some other works also explore the self/un-supervised learning objective under the NAS framework. UnNAS [20] explores how different pretexts tasks can replace the supervised discrimination

task. It shows that labels are not necessary for NAS, and metrics used in pre-text tasks can be a good proxy for the structure selection. RLNAS [45] shifts from the performance-based evaluation metric to the convergence-based metric, and it uses the random labels to generate supervision signals. Among self/un-supervised NAS works, SSNAS [15] and CSNAS [26] are most similar to our SSWP-NAS as they also explore the contrastive learning under the NAS framework. Nevertheless, they only consider the picture from an architecture optimization perspective, and their frameworks still fall into proxy-based searching. Dissimilarly, we are searching for the architecture and concomitant weights as integrity.

3 Methodology

In this section, we first introduce the prior knowledge about differential NAS which serves as the foundation of our framework. Then we detail how to extend the differential NAS towards the weight-preserving search and self-supervised optimization. Afterward, we investigate the network inflation issue and propose our solution. Finally, we demonstrate how to search a network (architecture plus weights) using SSWP-NAS.

3.1 Preliminary: differentiable NAS

Neural architecture search (NAS) task is generally formulated as a bi-level optimization task [1,6] where the upper-level variable α refers to the architecture parameters and lower-level variable w represents the operation parameters:

$$\min_{\alpha} \mathcal{L}_{val}(w^*(\alpha), \alpha) \quad (1)$$

$$s.t. \quad w^*(\alpha) = \operatorname{argmin}_w \mathcal{L}_{train}(w, \alpha) \quad (2)$$

In practice, two sets of parameters are optimized in an alternative manner that temporarily reduces the bi-level optimization to a single-level optimization. Differential NAS (DARTS) [23] further includes the architecture parameters directly to the computational graph through relaxing the search space from discrete to continuous. Thus it can effectively evolve both architecture structure and operation weights leveraging the stochastic gradient descent [33] techniques.

Inspired by the success of manually-designed structural motifs, NAS methods also shift the searching target from the intact architecture to a cell structure [50,21,32] which is then used to stack the final architecture repetitively. DARTS [23] also adopts this strategy by first searching a cell unit using a light-weight proxy architecture followed by constructing the final architecture with the searched cell structure. And the DARTS search space [23] represented by a cell can be interpreted as a directed acyclic graph (DAG) which consists of 7 nodes and 14 edges. For each edge, it is associated with a collection of candidate

operations \mathcal{O} weighted by a real-valued vector $\alpha^{(i,j)}$ of size $M = |\mathcal{O}|$. And the information flow $f_{i,j}$ between node x_i and x_j is defined as:

$$f_{i,j}(x_i) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})} o(x_i) \quad (3)$$

where $i < j$ and an intermediate node is computed based on all its predecessors: $x_j = \sum_{i < j} f_{i,j}(x_i)$. The final output of a cell is the concatenation of 4 intermediate nodes (except 2 input nodes and 1 output node) over the channel dimension. Here we refer the DARTS [23] paper for more details.

3.2 Towards Weight-preserving

One of the most prominent challenges in NAS is the surge of the memory footprint and the computational overhead. As a remedy, the decomposition from the whole architecture to the cell structure and the surrogate architecture searching is proposed to function as a trade-off between accuracy and efficiency. While reducing the computational overhead, the prevalence of such proxy strategy excludes the possibility of the weight-preserving property from the very beginning. Due to the non-identical structures used during the search and the train stage, concomitant weights discovered in the search phase are abandoned, and only the searched cell structure is reserved for the downstream application.

To retain the concomitant weights, we first build our target architecture using the DARTS [23] search space, which is well-established and contains abundant sub-graphs [43]. However, different from the DARTS and following advances [4,18], we align the dimensionality (widths and depths) of the architecture utilized during the search and the train stage. We further allow different cell structures at each level of the architecture. That is, instead of searching a single cell structure as the building unit, we search a set of cells to form the final architecture directly. As clued in section 3.1, this over-parameterized formation will consume $M = |\mathcal{O}|$ times memory footprint compared with a compact architecture. Disregarding the proxy strategy that contradicts the weight-preserving property in nature, we resort to stochastic algorithms [41,2,42] to cut down memory usage. Here, we adopt the path-binarization strategy, which is first proposed in ProxylessNAS [2] to overcome the accuracy degeneration issue caused by the depth-gap between the surrogate and the final architecture. Specifically, only a single operation on an edge is sampled and activated stochastically according to a learned distribution at each iteration. As such, the memory footprint during the search is reduced to the same magnitude as a compact architecture.

Through renouncing the proxy strategy and adopting the stochastic technique, we meet the indispensable conditions towards the weight-preserving. While in practice, we observe a consistently undesirable edge gained by the skip-connection operation. Similar phenomenons have also been observed in previous supervised and proxy-based search [18,4,35]. Such over-ratings for non-parameterized operations not only hinder the general quality of the structure, but also result in insufficient updates of parameterized operations during the

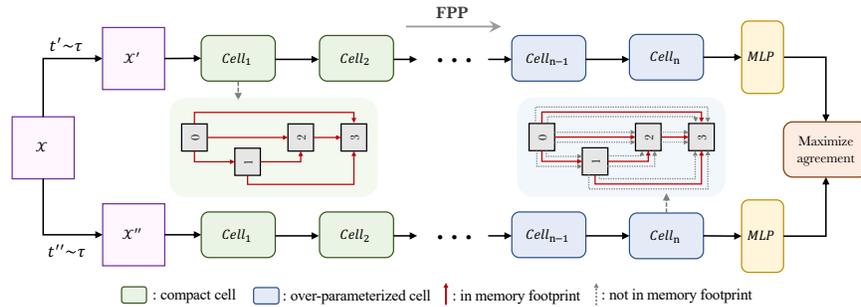


Fig. 2: The proposed SSWP-NAS framework. τ is a collection of data augmentations, t' and t'' are two transformations sampled from τ that transfer image x to two different views x' and x'' , respectively. Cell-structure plotted here is simplified for the interpretation purpose.

lower-level optimization. To this end, we introduce a non-parameterized operation dropout ($p = 0.2$ across our experiments) during the lower-level optimization. By doing so, we implicitly regularize the importance of non-parameterized operation and increase the sampling odds for parameterized operation without direct interference with the learned upper-level parameter distribution.

3.3 Towards Self-supervised Learning

It is conceptually straightforward to remove the dependency of the manual annotation by replacing the supervised signal with a self/un-supervised counterpart. Following the common practice of NAS, we are seeking a favorable architecture through modeling the conditional probability (discriminative model) rather than learning the joint distribution (generative model). As such, we subscribe our exploration for the learning objective to the discriminative sub-field. And among the bag of self/un-supervised discriminative pretext tasks, contrastive learning based methodologies [28,38,11,3,10] demonstrate an outstanding representation learning ability over the peers. So we further focus our attention on probing how the contrastive learning copes with the designed weight-preserving search process.

Given that the concept of the positive and negative pair lies at core of the contrastive learning, we study two representative methods, SimCLR [3] and BYOL [10] which are formulated with and without negative pairs, respectively. According to our pilot trials, BYOL only achieves 3.05% top-1 error rate on cifar-10 while SimCLR achieving 2.56%. BYOL lags behind a lot. We observe that the failure of the BYOL framework originates from the incompatible patterns between the stochastic operation sampling and exponential moving average (EMA). BYOL essentially forms a teacher-student pair using the current snapshot of the operation weights θ_t and its corresponding EMA counterpart $\theta'_t = \tau\theta'_{t-1} + (1 - \tau)\theta_t$ where τ is a smoothing factor. Then BYOL employs a

symmetric consistency loss between the pair as the supervision signal. However, stochastic operation sampling turns the parameter distribution θ_t into ξ_t where ξ and θ are different operation distributions. As a result, this inconsistency causes the mode collapse of the EMA, and thus fails the BYOL framework. As such, we take the strategy from SimCLR [3] where positive pairs are formed by two transformations of the same image and negative pairs are constructed using different images sampled from a mini-batch. The intact search architecture is then optimized in an end-to-end manner through minimizing the infoNCE loss [36,40] with mini-batch size N :

$$\mathcal{L}_{batch} = \sum_{i=1}^{2N} \mathcal{L}_i \quad (4)$$

$$\mathcal{L}_i = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_{j(i)})/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)} \quad (5)$$

Here, $\mathbf{z}_i = P(E(\tilde{\mathbf{x}}_i))$. $E(\cdot)$ denotes the searching encoder architecture and $P(\cdot)$ is a projection neck (multi-layer perceptron) added at end of the searching structure. The Subscript i and $j(i)$ refers to two different views of the same image. And $\text{sim}(\cdot)$ measures the cosine similarity between two given vectors. Fig. 2 displays the construction of our framework. This self-supervision pipeline is totally orthogonal to the search space and search process design. It is worth noting that $E_t(\cdot) \neq E_{t-1}(\cdot)$ at different time steps and $P(\cdot)$ is only employed during the search stage and being replaced by a linear classification layer for supervised train stage. By leveraging this self-guided signal generation method, SSWP-NAS gets rid of the dependency on manual labels.

3.4 Network Inflation Challenge

Despite a neat couple made by the weight-preserving search and the self-supervised optimization, it is challenging to optimize such a framework in practice. We observe that the difficulty comes from the combination of the over-parameterized structure and the stochastic operation sampling strategy, which in together we referred to as *network inflation* issue.

From the macro perspective, we are optimizing two sets of over-parameterized variable distributions α and w alternatively through the whole searching phase. At the end of search process, differentiable NAS [23] relies on the non-linear prune operation to approximate the compact variable $\tilde{\alpha}^*$ and \tilde{w}^* using the over-parameterized variable α^* and w^* . This leaping relaxation essentially relies on the good generalization ability of the hierarchical convolutional structures. From a micro view, a hierarchical convolutional structure can be viewed as a sequential model. Given a random tie-breaking index i in a N -layer structures at time step t , we are maximizing posterior probability $\arg \max_{w_{i \sim N}} p(w_{i \sim N} | w_{1 \sim i-1}^t, D)$ where D is the data distribution and subscript $a \sim b$ denotes layer index from a to b . Due to the stochastic sampling strategy, conditions $w_{1 \sim i-1}^{t-1} \neq w_{1 \sim i-1}^t$ while $w_{i \sim N}$ can be deemed as unchanged. Even though stochastic algorithm theoretically guarantees the same global convergence, if it exists, this non-stationary condition

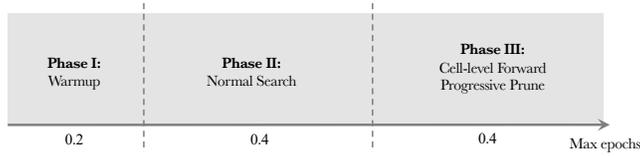


Fig. 3: Different phases during the search stage.

increases the difficulty of optimization at each time step. Here we take the lower-level variable w as an example, and the upper-level variable follows the similar formulation.

To handle the network inflation challenge, we propose a simple and effective solution, namely *forward progressive prune (FPP)*. Different from the general differentiable NAS pipeline where an architecture prune process is conducted for the target cell structure at the end of the search stage. We impose a cell-level progressive prune in the forward propagation direction during the search stage. At each prune step, we prune all edges contained in a cell, and for each edge, we only keep the operation with highest learned credit. By doing so, we reformulate the optimization target as $\arg \max_{w_{i \sim N}} p(w_{i \sim N} | \tilde{w}_{1 \sim i-1}^*, D)$ where the condition $\tilde{w}_{1 \sim i-1}^*$ is fixed and we transfer the optimization target from $p(w_{1 \sim N})$ to $p(w_{i \sim N} | \tilde{w}_{1 \sim i-1}^*)$ which is closer to the final goal $p(\tilde{w}_{1 \sim N}^*)$. As such, FPP gradually aligns the searching optimization target with the final objective. And it fixes the non-stationary conditions at some point during the search and allows the following layers to adjust according to the preceding layers. The forward propagation direction design follows the common acknowledgment that shallow layers of a CNN architecture capture the low-level features which are easier to learn and the deeper layers grab more semantics which rely on low-level features to compose.

3.5 Searching with SSWP-NAS

SSWP-NAS is a proxy-free search, so one can search the target architecture, not surrogate architecture, directly on target data domain without using labels. As depicted in Fig. 3, the search stage of SSWP-NAS is divided into three phases. At the warm-up phase, we only update operation parameters to allow a better initialization of parameterized operations. Then we update both operation parameters and architecture parameters alternatively as the standard differentiable NAS. Finally, we start the FPP phase, in which we progressively perform cell-level prune. The only extra hyper-parameter we introduce is the time span of different phases. We empirically suggest that the split $[0.2, 0.4, 0.4]$ works well in general. Given the proportion for phase III, the time step for pruning of two adjacent cells is calculated as $T_{step} = \frac{\text{max epoch} \times \text{FPP ratio}}{\text{cell num}}$.

4 Experiments

Experimental section is organized as following. We first provide the general settings used across our experiments. In the second part, we isolate the concomitant weights and benchmark the architecture quality searched using SSWP-NAS. Then we conduct comprehensive ablation studies towards our designs. Finally, we thoroughly study the effectiveness of concomitant weights. Detailed searched architecture are attached in *Appendix C*.

4.1 Experimental Settings

Following the well-established benchmark, we conduct our experiments on CIFAR-10/100 [16] and ImageNet [34] three datasets. We search an architecture consisting 20 cells with 36 initial channels for 300 epochs using mini-batch size 96 on CIFAR-10/100. For ImageNet, we search a structure that contains 14 cells with 48 initial channels for 100 epochs using mini-batch size 256. The magnitudes of searched architectures are kept the same as the final architecture trained in DARTS [23]. All hyper-parameters used during the search and train stage, except our proposed designs, are inherited from DARTS [23] without extra tricks. We attach detailed hyper-parameters in *Appendix A*. And 1_{st} order optimization is employed across all settings. Experiments related to CIFAR-10/100 [16] are conducted on a single Nvidia A100 40GB and it scales up to 4 for ImageNet [34].

4.2 Benchmarking SSWP-NAS

CIFAR-10/100. We search SSWP-NAS for 300 ($e = 300$) and 500 ($e = 500$) epochs on CIFAR dataset. Then we train the searched architecture from scratch in a supervised manner. As shown in Table 1, SSWP-NAS $_{e=300}$ achieves state-of-the-art performance and outperforms DARTS (1_{st} order) significantly by searching the same epochs. When we extend the searching duration from 300 epochs to 500 epochs, SSWP-NAS surpasses existing methods by a clear margin.

ImageNet. We search 50 and 100 epochs on ImageNet [34] and then train the searched structure for 250 epochs. As displayed in Table 2, SSWP-NAS also achieves state-of-the-art accuracy on ImageNet under limited budgets. One potential drawback is that SSWP-NAS takes a relative longer time to search when considering the architecture solely as we are directly searching for the final structure instead of a surrogate structure. However, the weight-preserving property offsets this computational overhead as all other methods need an extra non-trivial pre-train step. By achieving state-of-the-art performance on both CIFAR and ImageNet datasets, it also manifests the generality of the proposed SSWP-NAS framework.

4.3 Ablation Study

For simplicity, we abbreviate self-supervised learning and supervised learning as SSL and SL, respectively, in the following sections.

Table 1: Comparisons between SSWP-NAS and state-of-the-art methods on CIFAR-10/100. Here the search cost only counts the time used for the search. proxy-based methods commonly need another architecture selection procedure which generally costs another 1 GPU day [23]. While this procedure is not needed in our framework.

Architecture	Test Error (%)		Params (M)	Search Cost (GPU days)	Search Type	Search Method
	CIFAR-10	CIFAR-100				
DenseNet-BC[14]	3.46	17.18	25.6	-	supervised	manual
NASNet-A[50]	2.65	-	3.3	1800	supervised	RL
AmoebaNet-A[32]	3.34 ± 0.06	-	3.2	3150	supervised	evolution
AmoebaNet-B[32]	2.55 ± 0.05	-	2.8	3150	supervised	evolution
PNAS[21]	3.41 ± 0.09	-	3.2	225	supervised	SMBO
Hierarchical Evolution [22]	3.75 ± 0.12	-	15.7	300	supervised	evolution
ENAS[30]	2.89	-	4.6	0.5	supervised	RL
NAONet[24]	3.18	15.67	10.6	200	supervised	NAO
DARTS (1st order)[23]	3.0 ± 0.14	17.76	3.3	1.5	supervised	gradient
DARTS (2nd order)[23]	2.76 ± 0.09	17.54	3.3	4.0	supervised	gradient
SNAS (moderate)[41]	2.85 ± 0.02	-	2.8	1.5	supervised	gradient
ProxylessNAS-G[2]	2.08	-	5.7	4.0	supervised	gradient
P-DARTS[4]	2.50	16.55	3.4	0.3	supervised	gradient
PC-DARTS[42]	2.57 ± 0.07	-	3.6	0.1	supervised	gradient
BayesNAS[47]	2.81 ± 0.04	-	3.4	0.2	supervised	gradient
CSNAS _{N=5} [26]	2.66 ± 0.07	-	3.4	1.0	self-supervised	SMBO-TPE
SSNAS[15]	2.61	16.64	-	-	self-supervised	gradient
SSWP-NAS_{e=300} [†]	2.56 ± 0.07	17.27	4.0	1.0	self-supervised	gradient
SSWP-NAS_{e=500} [†]	2.41 ± 0.07	16.47	3.8	1.8	self-supervised	gradient

[†]: run 5 times with different seeds.

Table 2: Comparison with state-of-the-art architectures on ImageNet (restricted resources)

Architecture	Test Error (%)		Params (M)	×+ (M)	Search Cost (GPU days)	Search Type	Search Method
	Top-1	Top-5					
Inception-v2 [37]	25.2	7.8	11.2	-	-	-	manual
MobileNet-v3 (Large 1.0) [13]	24.8	-	5.4	219	-	-	manual
ShuffleNet(2×)-v2 [25]	25.1	-	≈ 5	591	-	-	manual
NASNet-A [50]	26.0	8.4	5.3	564	2000	supervised	RL
AmoebaNet-A [32]	25.5	8.0	5.1	555	3150	supervised	evolution
AmoebaNet-B [32]	26.0	8.5	5.3	555	3150	supervised	evolution
PNAS [21]	25.8	8.1	5.1	588	255	supervised	SMBO
DARTS (2nd order) [23]	26.7	8.7	4.7	574	4	supervised	gradient
P-DARTS [4]	24.1	7.3	5.4	597	2.0	supervised	gradient
PC-DARTS [42]	24.2	7.3	5.3	597	3.8	supervised	gradient
ProxylessNAS [2]	24.9	7.5	7.1	465	8.3	supervised	gradient
SSNAS [15]	27.8	9.6	5.2	-	-	self-supervised	gradient
CSNAS _{N=5} [26]	25.8	8.3	5.1	590	2.5	self-supervised	SMBO-TPE
SSWP-NAS_{e=50}	24.8	7.8	5.0	597	3.5	self-supervised	gradient
SSWP-NAS_{e=100}	24.3	7.5	4.9	595	7	self-supervised	gradient

Table 3: Ablation studies on learning objective, FPP, skip-connection dropout and proxy-free search.

(a) **SSL vs. SL and FPP.** Ablations for learning objectives and FPP.

Name	Test Error (%)	Dataset	Search Epoch
SL	2.92	CIFAR-10	300
SL+BPP	2.93	CIFAR-10	300
SL+FPP	2.77	CIFAR-10	300
SSL	2.65	CIFAR-10	300
SSL+BPP	2.73	CIFAR-10	300
SSL+FPP	2.56	CIFAR-10	300

(b) **FPP time span.** Ablations for different ratios of searching phases.

Name	Test Error (%)	Ratio	Dataset	Search Epoch
shorter	2.65	[0.2, 0.6, 0.2]	CIFAR-10	300
longer	2.68	[0.2, 0.2, 0.6]	CIFAR-10	300
default	2.56	[0.2, 0.4, 0.4]	CIFAR-10	300

(c) **Dropout for skip-connection.** Ablations on different dropout rate for skip connections.

Name	Test Error (%)	Dataset	Search Epoch
no drop	2.63	CIFAR-10	300
0.5 drop	2.68	CIFAR-10	300
0.2 drop	2.56	CIFAR-10	300

(d) **Proxy-free search.** Ablations for searching w/o proxy dataset.

Name	Accuracy (%)			Search Epoch	Search Dataset	Train Dataset
	Top-1	Top-5	Epoch			
transfer	69.6	88.45	600	CIFAR-10	ImageNet-tiny	
proxy-free	70.53	88.54	300	ImageNet-tiny	ImageNet-tiny	

SSL vs. SL. In order to substantiate the effectiveness of SSL, we compare it with the SL (cross-entropy loss function). By keeping all other settings the same, we switch between SL and SSL objectives. As shown in Table 3a, SSL outperforms SL significantly in our framework. This result not only exhibits the inessentiality of human annotations for architecture search, but also suggests that the manual interference may function adversely by limiting the optimization manifold. SSL objective may support learning a more generic structure for feature extraction. Beyond this frank improvement, we further show that the SSL searched architecture boosts the self-supervised weight pre-training in *Appendix B*

Effectiveness of FPP. To demonstrate the effectiveness of the FPP module, we implement SSWP-NAS with and without FPP. When disregarding FPP, we remove phase III during the search and adjust the ratio of phases I and II to 0.2 and 0.8 accordingly. Then the architecture prune process is conducted once for all cells after the search stage. Besides the original FPP, we also implement a reversed version, namely backward progressive pruning (BPP), that starts pruning from the last cell and propagates backwardly. According to Table 3a, FPP consistently improves the qualities of the searched architectures under both SSL and SL scenarios. BPP, on the contrary, even degenerates the performances. This result coincides with our knowledge that shallow layers extract low-level features, which are easy to learn. In contrast, deeper layers capture higher-level semantic features and leverage the low-level features to compose.

Time span of FPP. In this subsection, we study the role of time span in the FPP process and recommend the default setting. We allow a longer and a shorter FPP time span by adjusting the ratio to [0.2, 0.2, 0.6] and [0.2, 0.6, 0.2], respectively. Table 3b compares the result of different time spans. It is shown that a balanced split between phase II and phase III strikes for a better result. Shorter

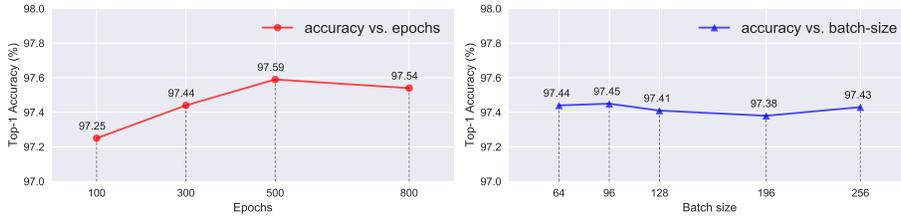


Fig. 4: For ablations of searching epochs, we fix the batch-size to 96 and increasing the searching epochs from 100 to 800. As for ablations of the batch-size, we fix the searching epoch as 300 and scale batch-size from 64 and 256.

FPP may result in the under-optimization of a single cell during the pruning phase. Longer FPP may squeeze the space of the normal bi-level optimization process, leading to potentially sub-optimal results. And a longer searching duration will relieve pressures of both phases by allowing more epochs at each phase and thus leads to a better result, as shown in Fig. 4.

Dropout for skip-connection. Here we add a dropout rate for skip-connection during the lower-level optimization process. As shown in Table 3c, using a $p = 0.2$ (result from a coarse grid search) improves the overall quality of the architecture. However, the architecture quality is relatively sensitive to the dropout rate as the skip-connection is a crucial component to prevent gradient vanishing issue. Over-suppressing of the skip-connection can hamper the performance.

Proxy-free search. It is well-established [12,31,3] that the performance of transfer learning drops when the gap between the target domain and the pre-trained domain is large. Given these priors, we focus on how proxy-free search can benefit the architecture quality. We use ImageNet-tiny [34] as the target domain and treat CIFAR-10 [16] as the proxy domain. Then we carry out transfer learning and proxy-free search. We doubled the search epochs on CIFAR-10 to keep the same iterations used during the search (ImageNet-tiny contains 10^5 images with size 64×64 , and CIFAR-10 contains 5×10^4 images with size 32×32 in the training dataset). So the only difference between the two settings is the domain gap (image dimension, context, etc.) itself. As shown in Table 3d, proxy-free search results in a better architecture quality. This result suggests that given the existence of labeled proxy datasets, it is still favorable to search on the target domain directly without the label.

Impact of epochs and batch-size. As verified in self-supervised weight pre-training [3,11,10], network weights benefit from longer training and larger mini-batch size. Therefore, we examine how these two hyper-parameters affect the architecture searching process. As displayed in Fig. 4, in consent with SSL weight pre-training, the quality of searched architecture also improves as the searching duration increases. It saturates around 500 epochs on CIFAR-10/100. This result encourages a longer searching epoch and manifests that our design does not suffer from the mode collapse issue [18,2]. And even searching for only 100 epochs (cost 0.4 GPU-day) on CIFAR-10, SSWP-NAS still surpasses DARTS

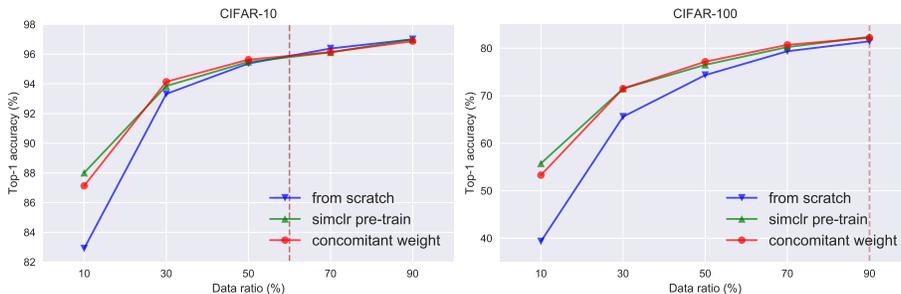


Fig. 5: Performances of the random initialization, simclr weight pre-train and concomitant weights on CIFAR-10/100 under different training data proportions.

(1_{st} order) [23]. Dissimilar to weight pre-training, the searched architecture does not take advantage of a larger mini-batch size in our approach referencing Fig. 4. As a result, one can use SSWP-NAS reliably under the limited GPU memory without worrying about the degeneration of architecture quality. However, this conclusion is only effective under the scenario where the magnitude of mini-batch size is a few hundred. We do not verify the impact of mini-batch size in thousands as used in several self-supervised weight-pretraining works [3,11,10,5].

4.4 Weight-preserving Benefits Semi-supervised Learning

In this section, we use both the searched architecture and concomitant weights from SSWP-NAS and probe their performances under semi-supervised scenarios. In particular, we first search using SSWP-NAS for 500 epochs with batch-size 96 on CIFAR dataset (train data). Then we gradually reduce the labeled training data from 90% to 10% with step-size 20% to mimic semi-supervised learning scenarios. Two baselines are included to demonstrate the effectiveness of concomitant weights. The first baseline noted as *random initialization* only uses the architecture. It is then trained from scratch using the given data ratio with a random initialization. The second baseline, *simclr pre-train*, in which we take the architecture searched by SSWP-NAS and pre-train it using SimCLR [3] for another 500 epochs using the same batch-size. The second baseline corresponds to a common two-stage framework without weight-preserving property. All three methods share the same architecture and the only difference is how they are initialized. We run the above settings on both CIFAR-10 and CIFAR-100.

As shown in Fig. 5, concomitant weights transfer clear positive information by outperforming the random initialization significantly. The gap between the concomitant weights and random initialization is bridged when using around 60% train data on CIFAR-10. However, this benefit diminishes not until around 90% of the train data on CIFAR-100. This result agrees with our intuition that self-supervised weight pre-training contributes more when labels are relatively scarce for each category, and the task is more challenging. More importantly, con-

Table 4: Ablation studies regarding concomitant weights.

(a) Comparisons of the two-stage pre-trained weights with the concomitant weights under different learning rates.

Learning rate	Accuracy (%)		Train data ratio (%)	Dataset
	SimCLR	Concomitant		
	Pre-train	Weight		
0.01	75.82	73.18	50%	CIFAR-100
0.025	76.1	74.96	50%	CIFAR-100
0.1	76.0	77.2	50%	CIFAR-100
0.2	75.99	77.54	50%	CIFAR-100
0.5	74.18	75.7	50%	CIFAR-100

(b) Impact of dropout and FPP on concomitant weights.

FPP	Dropout	Relative gain	Train data ratio	Dataset
✗	✗	-	50%	CIFAR-100
✗	✓	0.23	50%	CIFAR-100
✓	✗	0.01	50%	CIFAR-100
✓	✓	0.19	50%	CIFAR-100

comitant weights also surpass the two-stage framework pre-trained using SimCLR [3] except for the extremely scarce data setting (with only 10% data). On the one hand, this result substantiates our contribution by merging the common two-stage pipeline into one-stage; on the other hand, it further suggests an interesting potential that evolving both architecture and weights simultaneously may serve a better paradigm than the current isolated manner.

To better understand the differences between two-stage pre-trained weights and our concomitant weights, we further verify their reactions to different learning rates. Here we use 50% data from CIFAR-100 as a proxy-task. Table 4a shows that, unlike the typical two-stage pre-trained weights, concomitant weights consistently enjoy a larger learning rate. This result implies different statistical distributions between the concomitant weights and two-stage pre-trained weights.

Finally, we also investigate the impact of proposed modules on the quality of concomitant weights. Since architecture and concomitant weights are evolved simultaneously, we can not drive to the conclusion by direct comparing the accuracy. To this end, we use *relative gain* to isolate the impact of the target operation towards the concomitant weights. $relative\ gain = (acc_f^w - acc_s^w) - (acc_f^o - acc_s^o)$ where superscript *w/o* denotes with or without target operation. Subscript *f* and *s* represents fine-tuning and train from scratch, respectively. By doing so, we offset the influence of the architecture structure and focus on concomitant weights. As exhibited in Table 4b, dropout improves the quality of concomitant weights, this result agree with our assumption that insufficient update of parameterized operation may hinder the quality of the concomitant weights. And according to experiments, FPP does not have a clear impact on overall quality of concomitant weights.

5 Conclusion

In this work, Instead of trying to further reduce the computational overhead of the search process, the proposed SSWP-NAS strikes for a simplified workflow of NAS. It enjoys both self-supervising and weight-preserving two properties. Experiments show that self-supervised learning consistently benefits SSWP-NAS,

and the concomitant weights successfully merge the two-stage framework into the one stage. Comprehensive ablation studies substantiate the effectiveness of our proposed designs. For future work, it is important to probe how architecture and concomitant weights can bootstrap each other. And it is also compelling to design the new self-supervised paradigm specifically for joint-optimization of architecture and concomitant weights. From the practical consideration, enabling multi-objective and hardware-aware learning would be useful.

References

1. Anandalingam, G., Friesz, T.L.: Hierarchical optimization: An introduction. *Annals of Operations Research* **34**(1), 1–11 (dec 1992). <https://doi.org/10.1007/BF02098169>
2. Cai, H., Zhu, L., Han, S.: ProxylessNAS: Direct neural architecture search on target task and hardware. In: International Conference on Learning Representations (2019), <https://arxiv.org/pdf/1812.00332.pdf>
3. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. arXiv preprint arXiv:2002.05709 (2020)
4. Chen, X., Xie, L., Wu, J., Tian, Q.: Progressive differentiable architecture search: Bridging the depth gap between search and evaluation. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1294–1303 (2019)
5. Chen, X., He, K.: Exploring simple siamese representation learning. 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 15745–15753 (2021)
6. Colson, B., Marcotte, P., Savard, G.: An overview of bilevel optimization (2007)
7. Devries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. *CoRR* **abs/1708.04552** (2017), <http://arxiv.org/abs/1708.04552>
8. Doersch, C., Gupta, A.K., Efros, A.A.: Unsupervised visual representation learning by context prediction. 2015 IEEE International Conference on Computer Vision (ICCV) pp. 1422–1430 (2015)
9. Ghiasi, G., Lin, T.Y., Pang, R., Le, Q.V.: Nas-fpn: Learning scalable feature pyramid architecture for object detection. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 7029–7038 (2019)
10. Grill, J.B., Strub, F., Altché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., Piot, B., kavukcuoglu, k., Munos, R., Valko, M.: Bootstrap your own latent - a new approach to self-supervised learning. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M.F., Lin, H. (eds.) *Advances in Neural Information Processing Systems*. vol. 33, pp. 21271–21284. Curran Associates, Inc. (2020)
11. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.B.: Momentum contrast for unsupervised visual representation learning. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 9726–9735 (2020)
12. He, K., Girshick, R.B., Dollár, P.: Rethinking imagenet pre-training. 2019 IEEE/CVF International Conference on Computer Vision (ICCV) pp. 4917–4926 (2019)
13. Howard, A.G., Sandler, M., Chu, G., Chen, L.C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., Le, Q.V., Adam, H.: Searching for mobilenetv3. 2019 IEEE/CVF International Conference on Computer Vision (ICCV) pp. 1314–1324 (2019)

14. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4700–4708 (2017)
15. Kaplan, S., Giryes, R.: Self-supervised neural architecture search. CoRR **abs/2007.01500** (2020)
16. Krizhevsky, A.: Learning multiple layers of features from tiny images. Tech. rep. (2009)
17. Li, J., Zhou, P., Xiong, C., Hoi, S.: Prototypical contrastive learning of unsupervised representations. In: International Conference on Learning Representations (2021)
18. Liang, H., Zhang, S., Sun, J., He, X., Huang, W., Zhuang, K., Li, Z.: DARTS+: improved differentiable architecture search with early stopping. CoRR **abs/1909.06035** (2019)
19. Liu, C., Chen, L.C., Schroff, F., Adam, H., Hua, W., Yuille, A.L., Fei-Fei, L.: Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 82–92 (2019)
20. Liu, C., Doll'ar, P., He, K., Girshick, R.B., Yuille, A.L., Xie, S.: Are labels necessary for neural architecture search? In: ECCV (2020)
21. Liu, C., Zoph, B., Neumann, M., Shlens, J., Hua, W., Li, L.J., Fei-Fei, L., Yuille, A., Huang, J., Murphy, K.: Progressive neural architecture search. In: Proceedings of the European conference on computer vision (ECCV). pp. 19–34 (2018)
22. Liu, H., Simonyan, K., Vinyals, O., Fernando, C., Kavukcuoglu, K.: Hierarchical representations for efficient architecture search. In: International Conference on Learning Representations (2018), <https://openreview.net/forum?id=BJQRKzbA->
23. Liu, H., Simonyan, K., Yang, Y.: DARTS: Differentiable architecture search. In: International Conference on Learning Representations (2019)
24. Luo, R., Tian, F., Qin, T., Chen, E., Liu, T.Y.: Neural architecture optimization. In: Proceedings of the 32nd International Conference on Neural Information Processing Systems. p. 7827–7838 (2018)
25. Ma, N., Zhang, X., Zheng, H.T., Sun, J.: Shufflenet v2: Practical guidelines for efficient cnn architecture design. In: Proceedings of the European Conference on Computer Vision (ECCV) (September 2018)
26. Nguyen, N., Chang, J.M.: Contrastive self-supervised neural architecture search. CoRR **abs/2102.10557** (2021)
27. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: ECCV (2016)
28. van den Oord, A., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. ArXiv **abs/1807.03748** (2018)
29. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Advances in Neural Information Processing Systems 32, pp. 8024–8035. Curran Associates, Inc. (2019)
30. Pham, H., Guan, M., Zoph, B., Le, Q., Dean, J.: Efficient neural architecture search via parameters sharing. In: Dy, J., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 80, pp. 4095–4104. PMLR (10–15 Jul 2018)

31. Raghu, M., Zhang, C., Kleinberg, J., Bengio, S.: Transfusion: Understanding transfer learning for medical imaging. In: *Advances in Neural Information Processing Systems*. vol. 32. Curran Associates, Inc. (2019)
32. Real, E., Aggarwal, A., Huang, Y., Le, Q.V.: Regularized evolution for image classifier architecture search. *Proceedings of the AAAI Conference on Artificial Intelligence* **33**(01), 4780–4789 (Jul 2019)
33. Ruder, S.: An overview of gradient descent optimization algorithms. *CoRR* **abs/1609.04747** (2016)
34. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* **115**(3), 211–252 (2015). <https://doi.org/10.1007/s11263-015-0816-y>
35. Shu, Y., Wang, W., Cai, S.: Understanding architectures learnt by cell-based neural architecture search. In: *International Conference on Learning Representations* (2020)
36. Sohn, K.: Improved deep metric learning with multi-class n-pair loss objective. In: Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*. vol. 29. Curran Associates, Inc. (2016), <https://proceedings.neurips.cc/paper/2016/file/6b180037abbeba991d8b1232f8a8ca9-Paper.pdf>
37. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2016)
38. Tian, Y., Krishnan, D., Isola, P.: Contrastive multiview coding (2020)
39. Wang, N., Gao, Y., Chen, H., Wang, P., Tian, Z., Shen, C., Zhang, Y.: Nas-fcos: Fast neural architecture search for object detection. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2020)
40. Wu, Z., Xiong, Y., Yu, S.X., Lin, D.: Unsupervised feature learning via non-parametric instance discrimination. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition* pp. 3733–3742 (2018)
41. Xie, S., Zheng, H., Liu, C., Lin, L.: SNAS: stochastic neural architecture search. In: *International Conference on Learning Representations* (2019)
42. Xu, Y., Xie, L., Zhang, X., Chen, X., Qi, G.J., Tian, Q., Xiong, H.: Pc-darts: Partial channel connections for memory-efficient architecture search. In: *International Conference on Learning Representations* (2020)
43. Ying, C., Klein, A., Christiansen, E., Real, E., 0002, K.M., Hutter, F.: NAS-Bench-101: Towards Reproducible Neural Architecture Search. In: *Proceedings of the 36th International Conference on Machine Learning*. pp. 7105–7114. PMLR (2019)
44. Zhang, R., Isola, P., Efros, A.A.: Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 645–654 (2017)
45. Zhang, X., Hou, P., Zhang, X., Sun, J.: Neural architecture search with random labels. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 10902–10911 (2021)
46. Zhang, Y., Qiu, Z., Liu, J., Yao, T., Liu, D., Mei, T.: Customizable architecture search for semantic segmentation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 11633–11642 (2019)
47. Zhou, H., Yang, M., 0012, J.W., Pan, W.: BayesNAS: A Bayesian Approach for Neural Architecture Search. In: *Proceedings of the 36th International Conference on Machine Learning*. pp. 7603–7613. PMLR (2019)

48. Zhuang, C., Zhai, A., Yamins, D.: Local aggregation for unsupervised learning of visual embeddings. 2019 IEEE/CVF International Conference on Computer Vision (ICCV) pp. 6001–6011 (2019)
49. Zoph, B., Le, Q.V.: Neural architecture search with reinforcement learning. ArXiv [abs/1611.01578](https://arxiv.org/abs/1611.01578) (2017)
50. Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V.: Learning transferable architectures for scalable image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8697–8710 (2018)

6 Supplementary Materials

We abbreviate self-supervised learning and supervised learning as SSL and SL for simplicity.

6.1 Detailed Hyper-parameters

Detailed hyper-parameters used for both search and train stages are listed in Table 5. We follow the same suggestions as described in DARTS [23] and ProxylessNAS [2]. Variable w refers to the lower-level parameter (operation weights) and α denotes upper-level parameter (architecture weights). To speed up the training process for ImageNet, we also employ the distributed data-parallel, the automatic mixed precision, and the synchronized batch normalization techniques implemented in pytorch [29] framework.

Table 5: Detailed hyper-parameters used across the experiments.

Hyper-parameters	CIFAR-10/100		ImageNet	
	Search	Train	Search	Train
batch size	96	128	256	1024
learning rate (w)	0.025	0.025	0.025	0.4
minimum learning rate (w)	0	0	0	0
optimizer (w)	sgd	sgd	sgd	sgd
scheduler (w)	CosineAnnealing	CosineAnnealing	CosineAnnealing	CosineAnnealing
momentum (w)	0.9	0.9	0.9	0.9
weight decay (w)	4×10^{-5}	3×10^{-4}	4×10^{-5}	3×10^{-4}
learning rate (α)	0.001	-	0.001	-
optimizer (α)	adam	-	adam	-
adam β_1	0	-	0	-
adam β_2	0.99	-	0.99	-
auxiliary weight	-	0.4	-	0.4
cutout [7] length	-	16	-	-
drop-path rate	-	0.3	-	-

6.2 Self-supervised architecture search benefits self-supervised weight pre-training

In this section, we show that architecture searched by self-supervised learning objective also benefits the typical self-supervised weight pre-training method. We first search two architectures using supervised and self-supervised learning objectives, respectively. Then we pre-train the two searched architectures using the SimCLR [3] framework. Finally, we conduct the standard linear probe experiments [3,11,10,5] on pre-trained two networks with the same setting. As shown in Table 6, SSL searched architecture also benefits the downstream SSL based weight pre-training in a typical two-stage workflow (architecture search and weight pre-training are treated separately). This result suggests that the performance of the architecture and its corresponding weights are correlated under the SSL framework. It is recommended to use SSL-based NAS instead of SL-based NAS when considering using a self-supervised weight pre-training method.

Table 6: Linear probe results for architectures searched with SSL and SL.

Name	Accuracy (%)	Search epoch	Pre-train epoch	Search dataset	Pre-train dataset
SL	64.4	300	150	CIFAR-10	CIFAR-10
SSL	65.0	300	150	CIFAR-10	CIFAR-10

6.3 Architecture Structure Searched by SSWP-NAS

Detailed architecture structures searched by SSWP-NAS are attached as Fig. 6.

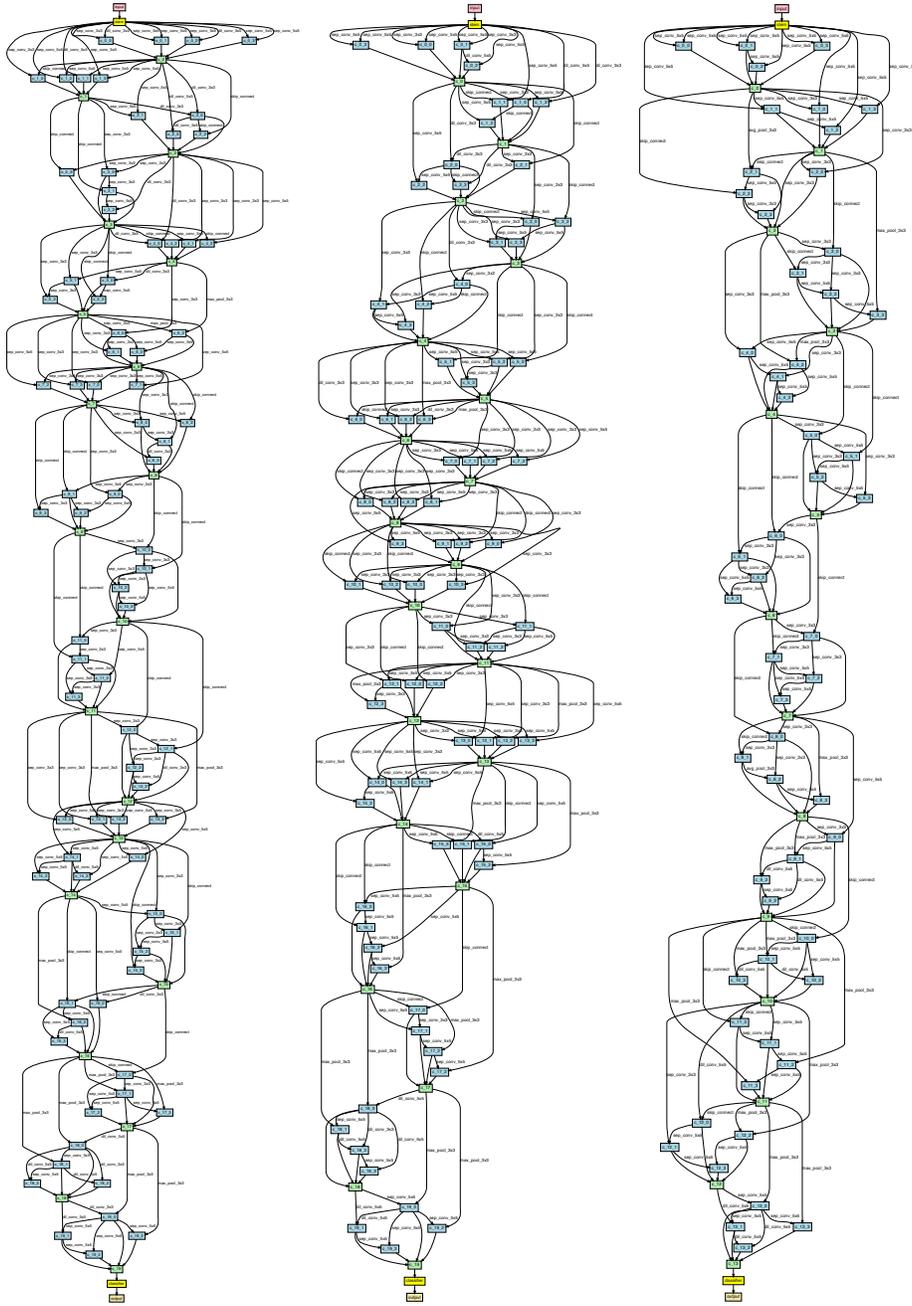


Fig. 6: Detailed architecture structures searched by SSWP-NAS. Architectures from left to right correspond to searching 300 epochs on CIFAR-10/100, 500 epochs on CIFAR-10/100, and 100 epochs on ImageNet accordingly.