

SGAS: Sequential Greedy Architecture Search

<https://www.deepgcns.org/auto/sgas>

Guohao Li^{*1}, Guocheng Qian^{*1}, Itzel C. Delgadillo^{*1}, Matthias Müller², Ali Thabet¹, Bernard Ghanem¹

¹King Abdullah University of Science and Technology (KAUST), Saudi Arabia

²Intelligent Systems Lab, Intel Labs, Germany

Abstract

Architecture design has become a crucial component of successful deep learning. Recent progress in automatic neural architecture search (NAS) shows a lot of promise. However, discovered architectures often fail to generalize in the final evaluation. Architectures with a higher validation accuracy during the search phase may perform worse in the evaluation (see Figure 1). Aiming to alleviate this common issue, we introduce sequential greedy architecture search (SGAS), an efficient method for neural architecture search. By dividing the search procedure into sub-problems, SGAS chooses and prunes candidate operations in a greedy fashion. We apply SGAS to search architectures for Convolutional Neural Networks (CNN) and Graph Convolutional Networks (GCN). Extensive experiments show that SGAS is able to find state-of-the-art architectures for tasks such as image classification, point cloud classification and node classification in protein-protein interaction graphs with minimal computational cost.

1. Introduction

Deep learning has revolutionized computer vision by learning features directly from data. As a result deep neural networks have achieved state-of-the-art results on many difficult tasks such as image classification [13], object detection [32], object tracking [39], semantic segmentation [11], depth estimation [16] and activity understanding [7], to name just a few examples. While there was a big emphasis on feature engineering before deep learning, the focus has now shifted to architecture engineering. In particular many novel architectures have been proposed, such as LeCun [28], AlexNet [27], VGG [47], GoogLeNet [49], ResNet [19], DenseNet [22], ResNeXt [57] and SENet [21]. Results on each of the above mentioned tasks keep improving every year by innovations in architecture design. In essence, the community has shifted from feature engineering to architecture engineering.

^{*}equal contribution

Search Ranking	Evaluation Ranking			
	DARTS (1st Order)	DARTS (2nd Order)	SGAS (Criterion 1)	SGAS (Criterion 2)
1	4	8	2	4
2	2	5	3	6
3	7	9	1	1
4	8	6	6	2
5	1	4	8	3
6	5	3	7	10
7	10	7	4	5
8	9	1	5	8
9	6	10	10	9
10	3	2	9	7
Kendall τ	0.16	-0.29	0.56	0.42
Avg. Acc.	97.15	97.18	97.34	97.33

Figure 1. **Comparison of search-evaluation Kendall τ coefficients.** We show Kendall τ correlations for architecture rankings between the search and the evaluation phase of DARTS and SGAS. Architectures are obtained from 10 independent search runs.

In recent years, many efforts have been made to reduce the manual intervention required to obtain better models for a particular task. As a matter of fact, a new area of research, commonly referred to as meta-learning, has emerged in order to tackle such problems. The idea of meta-learning is to leverage prior experience in order to quickly find good algorithm configurations, network architectures and any required parameters for a new learning task. Examples of recent meta-learning approaches include automatic hyper-parameter search [15], data-augmentation [12], finding novel optimizers [2] and architecture search [65]. In particular, architecture search has sparked a lot of interest in the community. In this task, the search space is huge and manual search is prohibitive.

Early work by Zoph *et al.* [65], based on Reinforcement Learning, has shown very promising results. However, its high computational cost has prevented widespread adoption. Recently, differentiable architecture search (DARTS) [35] has been proposed as an alternative which makes architecture search differentiable and much more efficient. This has opened up a path towards computationally feasible ar-

chitecture search. However, despite their success, current approaches still have a lot of limitations. During the search phase, network architectures are usually constructed from basic building blocks and evaluated on a validation set. Due to computational cost, the size of considered architectures is limited. In the evaluation phase, the best building blocks are used to construct larger architectures and they are evaluated on the test set. **As a result there is a large discrepancy between the validation accuracy during search and the test accuracy during evaluation.** In this work, we propose a novel greedy architecture search algorithm, SGAS, which addresses this discrepancy and searches very efficiently.

Contributions. Our contributions can be summarized as the following: (1) **We propose SGAS, a greedy approach for neural architecture search with high correlation between the validation accuracy during the search phase and the final evaluation accuracy.** (2) Our method discovers top-performing architectures with much less search cost than previous state-of-the-art methods such as DARTS. (3) **Our proposed method is able to search architectures for both CNNs and GCNs across various datasets and tasks.**

2. Related Work

In the past, considerable success was achieved with hand-crafted architectures. One of the earliest successful architectures was LeNet [28], a very simple convolutional neural network for optical character recognition. Other prominent networks include AlexNet [27], VGG [47] and GoogLeNet [49] which revolutionized computer vision by outperforming all previous approaches in the ImageNet [13] challenge by a large margin. ResNet [19] and DenseNet [22] were further milestones in architecture design. They showed the importance of residual and dense connections for designing very deep networks, an insight that influences modern architecture design to this day. Until recently, architecture innovations were a result of human insight and experimentation. The first successful attempts for architecture search were using reinforcement learning [65] and evolutionary algorithms [43]. These works were extended with NASNet [66] where a new cell-based search space and regularization technique were proposed. Another extension, ENAS [40], represents the entire search space as a single directed acyclic graph. A controller discovers architectures by searching for subgraphs that maximize the expected reward on the validation set. This setup allows for parameter sharing between child models making search very efficient. Further, PNAS [33] introduced a sequential model-based optimization (SMBO) strategy in order to search for structures of increasing complexity. PNAS needs to evaluate 5 times less models and reduces the computational cost by a factor of 8 compared to NASNet. Yet, PNAS still requires thousands of GPU hours. One shot approaches

[6, 5, 8] further reduce the search time by training a single over-parameterized network with inherited/shared weights. In order to search in a continuous domain [44, 1, 46, 53], DARTS [35] proposes a continuous relaxation of the architecture representation, making architecture search differentiable and hence much more efficient. As a result, DARTS is able to find good convolutional architectures at a fraction of the computational cost making NAS broadly accessible. Owing to the large success of DARTS, several extensions have been proposed recently. SNAS [58] optimizes parameters of a joint distribution for the search space in a cell. The authors propose a search gradient which optimizes the same objective as RL-based NAS, but leads to more efficient structural decisions. P-DARTS [9] attempts to overcome the depth gap issue between search and evaluation. This is accomplished by increasing the depth of searched architectures gradually during the training procedure. PC-DARTS [61] leverages the redundancy in network space and only samples a subset of channels in super-net during search to reduce computation.

3. Methodology

3.1. Preliminary - DARTS

By reducing the search problem to searching for the best cell structure, cell-based NAS methods [66, 33, 43] are able to learn scalable and transferable architectures. The networks are composed of layers with identical cell structure but different weights. A cell is usually represented as a directed acyclic graph (DAG) with N nodes including two input nodes, several intermediate nodes and a single output node. Each node is a latent representation denoted as $x^{(i)}$, where i is its topological order in the DAG. Each directed edge (i, j) in the DAG is associated with an operation $o^{(i,j)}$ that transfers the information from node $x^{(i)}$ to node $x^{(j)}$. In *Differentiable Architecture Search* (DARTS) [35] and its variants [58, 9, 61, 18], the optimal architecture is derived from a discrete search space by relaxing the selection of operations to a continuous optimization problem. During the search phase, the operation of each edge is parameterized by architectural parameters $\alpha^{(i,j)}$ as a softmax mixture over all the possible operations within the operation space \mathcal{O} , *i.e.* $\bar{o}^{(i,j)}(x^{(i)}) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})} o(x^{(i)})$. The input nodes are represented by the outputs from the previous two cells. Each intermediate node aggregates information flows from all of its predecessors, $x^{(j)} = \sum_{i < j} \bar{o}^{(i,j)}(x^{(i)})$. The output node is defined as a concatenation of a fixed number of its predecessors. The learning procedure of architectural parameters involves a bi-level optimization problem:

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

$$\text{s.t. } \mathcal{W}^*(\mathcal{A}) = \operatorname{argmin}_{\mathcal{W}} \mathcal{L}_{train}(\mathcal{W}, \mathcal{A}) \quad (2)$$

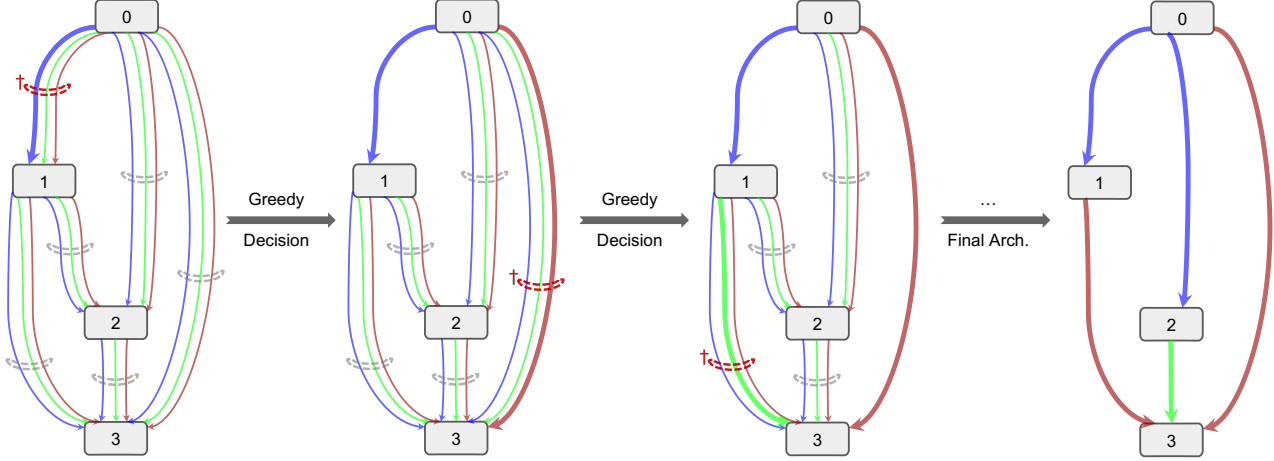


Figure 2. **Illustration of Sequential Greedy Architecture Search.** At each greedy decision epoch, an edge (i^\dagger, j^\dagger) is selected based on the selection criterion. A greedy decision will be made for the edge (i^\dagger, j^\dagger) by replacing $\bar{o}^{(i^\dagger, j^\dagger)}$ with $o^{(i^\dagger, j^\dagger)} = \operatorname{argmax}_{o \in \mathcal{O}} \alpha_o^{(i^\dagger, j^\dagger)}$. The corresponding architectural parameter $\alpha^{(i^\dagger, j^\dagger)}$ will be removed from the bi-level optimization. Operations which were not chosen in a mixture operation will be pruned. At the end of the search phase, a stand-alone architecture without weight sharing will be obtained.

\mathcal{L}_{train} and \mathcal{L}_{val} denote the training and validation loss respectively. Owing to the continuous relaxation, the search is realized by optimizing a supernet. \mathcal{W} is the set of weights of the supernet and \mathcal{A} is the set of the architectural parameters. DARTS [35] proposed to solve this bi-level problem by a first/second order approximation. At the end of the search, the final architecture is derived by selecting the operation with highest weight for every mixture operation, $o^{(i,j)} = \operatorname{argmax}_{o \in \mathcal{O}} \alpha_o^{(i,j)}$.

3.2. Search-Evaluation Correlation

A popular pipeline of existing NAS algorithms [66, 35] includes two stages: a search phase and an evaluation phase. In order to reduce computational overhead, previous works [66, 35] first search over a pre-defined search space with a lightweight proxy model on a small proxy dataset. After the best architecture cell/encoding is obtained, the final architecture is built and trained from scratch on the target dataset. This requires that the true performance during evaluation can be inferred during the search phase. However, this assumption usually does not hold due to the discrepancy in dataset, hyper-parameters and network architectures between the search and evaluation phases. The best ranking derived from the search phase does not imply the actual ranking in the final evaluation. In practice, the correlation between the performances of derived architectures during the search and evaluation phases is usually low. In this paper, we refer to this issue as **degenerate search-evaluation correlation**. Recent work by Sciuto *et al.* [45] also analyzes this issue and suggests that the Kendall τ metric [24] could be used to evaluate the search phase. They show that the widely used *weight sharing* technique actually decreases the correlation. The Kendall τ metric [24] is a common

measurement of the correlation between two rankings. The Kendall τ coefficient can be computed as $\tau = \frac{N_c - N_d}{\frac{1}{2}n(n-1)}$, where N_c and N_d are the number of concordant pairs and the number of discordant pairs respectively. It is a number in the range from -1 to 1 where -1 corresponds to a perfect negative correlation and 1 to a perfect positive correlation. If the Kendall τ coefficient is 0 , the rankings are completely independent. An ideal NAS method should have a high search-evaluation Kendall τ coefficient. We take DARTS [35] as an example and show its Kendall τ in Figure 1. It is calculated between the rankings during search phase and evaluation phase. The rankings are determined according to the validation accuracy and the final evaluation accuracy after 10 different runs on the CIFAR-10 dataset. The Kendall τ coefficients for DARTS are only 0.16 and -0.29 for the 1st-order and 2nd-order versions respectively. Therefore, it is impossible to make reliable predictions regarding the final test accuracy based on the search phase.

3.3. Sequential Greedy Architecture Search

In order to alleviate the **degenerate search-evaluation correlation** problem, the core aspects are to reduce (1) the discrepancy between the search and evaluation phases and (2) the negative effect of weight sharing. We propose to solve the bi-level optimization (Equation 1, 2) in a sequential greedy fashion to reduce the model discrepancy and the weight sharing progressively. As mentioned in Section 3.1, DARTS-based methods [35, 9, 61] solve the relaxed problem fully and obtain all the selected operations at the end. Instead of solving the complete problem directly, we divide it into sub-problems and solve them sequentially with a greedy algorithm. The sub-problems are defined based on the directed edges in the DAG. We pick the operation

Algorithm 1: SGAS – Sequential Greedy Architecture Search

Create architectural parameters $\mathcal{A} = \{\alpha^{(i,j)}\}$ and supernet weights \mathcal{W}

Create a mixed operation $\bar{o}^{(i,j)}$ parameterized by $\alpha^{(i,j)}$ for each edge (i, j)

while not terminated **do**

1. Update undetermined architecture parameters \mathcal{A} by descending $\nabla_{\mathcal{A}} \mathcal{L}_{val}(\mathcal{W}, \mathcal{A})$
2. Update weights \mathcal{W} by descending $\nabla_{\mathcal{W}} \mathcal{L}_{train}(\mathcal{W}, \mathcal{A})$
(since the weights of unchosen operations are pruned, only the remaining weights need to be updated)
3. If a decision epoch, select an edge (i^\dagger, j^\dagger) based on the greedy *Selection Criterion*
Determine the operation by replacing $\bar{o}^{(i^\dagger, j^\dagger)}$ with $o^{(i^\dagger, j^\dagger)} = \arg\max_{o \in \mathcal{O}} \alpha_o^{(i^\dagger, j^\dagger)}$
Prune unchosen weights from \mathcal{W} , Remove $\alpha^{(i^\dagger, j^\dagger)}$ from \mathcal{A}

Derive the final architecture based on chosen operations

for edges greedily in a sequential manner and solve the remaining sub-problem iteratively. The iterative procedure is shown in Algorithm 1. **At each decision epoch, we choose one edge (i^\dagger, j^\dagger) according to a pre-defined selection criterion.** A greedy optimal choice is made for the selected edge by replacing the corresponding mixture operation $\bar{o}^{(i^\dagger, j^\dagger)}$ with $o^{(i^\dagger, j^\dagger)} = \arg\max_{o \in \mathcal{O}} \alpha_o^{(i^\dagger, j^\dagger)}$. The architectural parameters $\alpha^{(i^\dagger, j^\dagger)}$ and the weights of the remaining paths within the mixture operations are no longer needed; we prune and exclude them from the latter optimization. As a side benefit, the efficiency improves as parameters in \mathcal{A} and \mathcal{W} are pruned gradually in the optimization loop. The search procedure of the remaining \mathcal{A} and \mathcal{W} forms a new sub-problem which will be solved iteratively. **At the end of the search phase, a stand-alone network without weight sharing is obtained,** as illustrated in Figure 2. **Therefore, the model discrepancy is minimized and the validation accuracy during the search phase reflects the final evaluation accuracy much better.** To maintain the optimality, the design of the *selection criterion* is crucial. We consider three aspects of edges which are the *edge importance*, the *selection certainty* and the *selection stability*.

Edge Importance. Similar to DARTS [35], a *zero* operation is included in the search space to indicate a lack of connection. **Edges that are important should have a low weight in the *zero* operation.** Thus, the edge importance is defined as the summation of weights over non-zero operations:

$$S_{EI}^{(i,j)} = \sum_{o \in \mathcal{O}, o \neq \text{zero}} \frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})} \quad (3)$$

Selection Certainty. Entropy is a common measurement of uncertainty of a distribution. The normalized softmax weights of non-zero operations can be regarded as a distribution, $p_o^{(i,j)} = \frac{\exp(\alpha_o^{(i,j)})}{S_{EI}^{(i,j)} \sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})}$, $o \in \mathcal{O}, o \neq \text{zero}$. We define the *selection certainty* as the complement of the

normalized entropy of the operation distribution:

$$S_{SC}^{(i,j)} = 1 - \frac{-\sum_{o \in \mathcal{O}, o \neq \text{zero}} p_o^{(i,j)} \log(p_o^{(i,j)})}{\log(|\mathcal{O}| - 1)} \quad (4)$$

Selection Stability. In order to incorporate the history information, we measure the movement of the operation distribution. KullbackLeibler divergence and histogram intersection [48] are two popular methods to detect changes in distribution. For simplicity, we choose the latter one. The average selection stability at step T with a history window size K is computed as follows:

$$S_{SS}^{(i,j)} = \frac{1}{K} \sum_{t=T-K}^{T-1} \sum_{o_t \in \mathcal{O}, o_t \neq \text{zero}} \min(p_{o_t}^{(i,j)}, p_{o_T}^{(i,j)}) \quad (5)$$

In our experiments, we consider the following two criteria:

Criterion 1. An edge (i^\dagger, j^\dagger) with a high *edge importance* $S_{EI}^{(i,j)}$ and a high *selection certainty* $S_{SC}^{(i,j)}$ will be selected. We normalize $S_{EI}^{(i,j)}$ and $S_{SC}^{(i,j)}$, compute the final score and pick the edge with the highest score:

$$S_1^{(i,j)} = \text{normalize}(S_{EI}^{(i,j)}) * \text{normalize}(S_{SC}^{(i,j)}) \quad (6)$$

Criterion 2. In addition to *Criterion 1*, we also consider that the selected edge (i^\dagger, j^\dagger) should have a high *selection stability*. The final score is defined as follows:

$$S_2^{(i,j)} = S_1^{(i,j)} * \text{normalize}(S_{SS}^{(i,j)}) \quad (7)$$

where *normalize*(\cdot) denotes a standard Min-Max scaling normalization. For a fair comparison with existing works [65, 43, 35], two incoming edges are preserved for every intermediate node in the DAG. Once a node has two determined incoming edges, its other incoming edges will be pruned. We refer to our method as **Sequential Greedy Architecture Search (SGAS)**. Figure 1 shows that SGAS with Criterion 1 and 2 improves the Kendall τ correlation coefficients to 0.56 and 0.42 respectively. As expected from the much higher search-evaluation correlation SGAS outperform DARTS in terms of average accuracy significantly.

4. Experiments

We use our SGAS to automatically find architectures for both CNNs and GCNs. The CNN architectures discovered by SGAS outperform the state-of-the-art (SOTA) in image classification on CIFAR-10 [26] and ImageNet [13]. Similarly, the discovered GCN architectures outperform the state-of-the-art methods for point cloud classification on ModelNet [56] and node classification in biological graphs using the PPI [64] dataset.

4.1. Searching CNN architectures with SGAS

4.1.1 Architecture Search on CIFAR-10

As is common practice, we first search for normal cells and reduction cells with a small network for image classification on CIFAR-10. CIFAR-10 is a small popular dataset containing 50K training images and 10K testing images. Then, a larger network is constructed by making necessary changes in channel size and stacking the searched cells multiple times. The larger network is retrained on CIFAR-10 to compare its performance with other state-of-the-art methods. Finally, we show the transferability of our SGAS by stacking even more cells and evaluating on ImageNet. We show that SGAS consistently achieves the top performance.

Search Space. We keep our search space the same as DARTS, which has 8 candidate operations: *skip-connect*, *max-pool-3×3*, *avg-pool-3×3*, *sep-conv-3×3*, *sep-conv-5×5*, *dil-conv-3×3*, *dil-conv-5×5*, *zero*. **During the search phase, we stack 6 normal cells and 2 reduction cells to form a network. Two reduction cells are inserted at a network depth of 1/3 and 2/3 respectively.** The stride of each convolution in normal cells is 1, so the spatial size of an input feature map does not change. In reduction cells, convolutions with stride 2 are used to reduce the spatial resolution of feature maps. **There are 7 nodes with 4 intermediate nodes and 14 edges in each cell during search.** The first and second input nodes of the cell are set equal to the outputs of the two previous cells respectively. The output node of a cell is the depth-wise concatenation of all the intermediate nodes.

Training Settings. We keep the training setting the same as in DARTS [35]. A small network consisting of 6 normal cells and 2 reduction cells with an initial channel size 16 is trained on CIFAR-10. We perform architecture search for 50 epochs with a batch size of 64. SGD is used to optimize the model weights \mathcal{W} with an initial learning rate 0.025, momentum 0.9 and weight decay 3×10^{-4} . For architecture parameters \mathcal{A} , the Adam optimizer with an initial learning rate 3×10^{-4} , momentum (0.5, 0.999) and weight decay 10^{-3} is used. **Instead of training the entire super-net throughout the search phase, SGAS makes decisions sequentially in a greedy fashion. After warming up for 9 epochs, SGAS begins to select one operation for one se-**

lected edge every 5 epochs using *Criterion 1* or *Criterion 2* as the selection criterion. For *Criterion 2*, we set the history window size K to 4. The batch size is increased by 8 after each greedy decision, which further boosts the searching efficiency of SGAS. We provide a thorough discussion and ablation study on the choices of hyper-parameters in the **supplement material**. The search takes only 0.25 day (6 hours) on a single NVIDIA GTX 1080Ti.

4.1.2 Architecture Evaluation on CIFAR-10

We run 10 *independent searches* to get 10 architectures with *Criterion 1* or *Criterion 2*, as shown in Figure 1. To highlight the stability of the search method, we evaluate the discovered architectures on CIFAR-10 and report the mean and standard deviation of the test accuracy across those 10 models and the performance of the best model in Table 1. It is important to mention that other related works in Table 1 only report the mean and standard deviation for their *best architecture* with different runs on evaluation.

Training Settings. We train a large network of 20 cells with a initial channel size 36. The SGD optimizer is used during 600 epochs with a batch size of 96. The other hyper-parameters remain the same as the search phase. Cutout with length 16, auxiliary towers with weight 0.4 and path dropout with probability 0.3 are used as in DARTS [35].

Evaluation Results and Analysis. We compare our results with other methods in Table 1 and report the average and best performance for both *Criterion 1* and *Criterion 2*. We outperform our baseline DARTS by a significant margin with test errors of 2.39% and 2.44% respectively while only using 0.25 day (6 hours) on a single NVIDIA GTX 1080Ti.

4.1.3 Architecture Evaluation on ImageNet

The architecture evaluation on ImageNet uses the cell architectures that we obtained after searching on CIFAR-10.

Training Settings. We choose the 3 best performing cell architectures on CIFAR-10 for each *Criterion* and train them on ImageNet. For this evaluation, we build a large network with 14 cells and 48 initial channels and train for 250 epochs with a batch size of 1024. The SGD optimizer with an initial learning rate of 0.5, a momentum of 0.9 and a weight decay of 3×10^{-5} is used. **We run these experiments on 8 Nvidia Tesla V100 GPUs for three days.**

Evaluation Results and Analysis. In Table 2 we compare our models with SOTA hand-crafted architectures (manual) and models obtained through other search methods. **We apply the mobile setting for ImageNet, which has an image size of 224×224 and restricts the number of multi-add operations to 600M.** Our best performing models SGAS (Cri.1 best) and SGAS (Cri.2 best) outperform all the other methods with top-1 errors 24.2% and 24.1% respectively

Architecture	Test Err. (%)	Params (M)	Search Cost (GPU-days)	Search Method
DenseNet-BC [22]	3.46	25.6	-	manual
NASNet-A [66]	2.65	3.3	1800	RL
AmoebaNet-A [43]	3.34±0.06	3.2	3150	evolution
AmoebaNet-B [43]	2.55±0.05	2.8	3150	evolution
Hier-Evolution [34]	3.75±0.12	15.7	300	evolution
PNAS [33]	3.41±0.09	3.2	225	SMBO
ENAS [40]	2.89	4.6	0.5	RL
NAONet-WS [37]	3.53	3.1	0.4	NAO
DARTS (1 st order) [35]	3.00±0.14	3.3	0.4	gradient
DARTS (2 nd order) [35]	2.76±0.09	3.3	1	gradient
SNAS (mild) [58]	2.98	2.9	1.5	gradient
ProxylessNAS [8]	2.08	-	4	gradient
P-DARTS [9]	2.5	3.4	0.3	gradient
BayesNAS [63]	2.81±0.04	3.4	0.2	gradient
PC-DARTS [61]	2.57±0.07	3.6	0.1	gradient
SGAS (Cri.1 avg.)	2.66±0.24*	3.7	0.25	gradient
SGAS (Cri.1 best)	2.39	3.8	0.25	gradient
SGAS (Cri.2 avg.)	2.67±0.21*	3.9	0.25	gradient
SGAS (Cri.2 best)	2.44	4.1	0.25	gradient

Table 1. **Performance comparison with state-of-the-art image classifiers on CIFAR-10.** We report the average and best performance of SGAS (Cri.1) and SGAS (Cri.2). *Criterion 1* and *Criterion 2* are used in the search respectively. *Note that mean and standard derivation are computed across 10 independently searched architectures.

while using only a search cost of 0.25 GPU day on one NVIDIA GTX 1080Ti. SGAS (Cri.2) outperforms SGAS (Cri.1) showing the effectiveness of integrating *selection stability* into the selection criterion. The best performing cells of SGAS (Cri.2 best) are illustrated in Figure 3.

4.2. Searching GCN architectures with SGAS

Recently, GCNs have achieved impressive performance on point cloud segmentation [30], biological graph node classification [29] and video recognition [60] by training DeepGCNs [30, 29]. However, this hand-crafted architecture design requires adequate effort by an human expert. The main component of DeepGCNs is the GCN backbone. We explore an automatic way to design the GCN backbone using SGAS. Our backbone network is formed by stacking the graph convolutional cell discovered by SGAS. Our GCN cell consists of 6 nodes. We use fixed 1×1 convolutions in the first two nodes, and set the input to them equal to the output from the previous two layers. Our experiments on GCNs have two stages. First, we apply SGAS to search for the graph convolutional cell using a small dataset and obtain 10 architectures from 10 runs. Then, 10 larger networks are constructed by stacking each discovered cell multiple times. The larger networks are trained on the same dataset or a larger one to evaluate their performance. We report the best and average performance of these 10 architectures. We show the effectiveness of SGAS in GCN architecture search by comparisons with SOTA hand-crafted methods and Random Search.

4.2.1 Architecture Search on ModelNet10

ModelNet [56] is a dataset for 3D object classification with two variants, ModelNet10 and ModelNet40 containing objects from 10 and 40 classes respectively. We conduct GCN architecture search on ModelNet10 and then evaluate the final performance on ModelNet40.

Search Space. Our graph convolutional cell has 10 candidate operations: *conv-1×1*, *MRConv* [30], *EdgeConv* [55], *GAT* [52], *SemiGCN* [25], *GIN* [59], *SAGE* [17], *RelSAGE*, *skip-connect*, and *zero* operation. Please refer to our **supplement material** for more details of these GCN operators. We use *k nearest neighbor* in the first operation of each cell to construct edges (we use $k = 9$ by default unless it is specified). These edges are then shared in the following operations inside the cell. Dilated graph convolutions with the same linearly increasing dilation rate schedule as proposed in DeepGCNs [30] are applied to the cells.

Training Settings. We sample 1024 points from the 3D models in ModelNet10. We use 2 cells with 32 initial channels and search the architectures for 50 epochs with batch size 28. SGD is used to optimize the model weights with initial learning rate 0.005, momentum 0.9 and weight decay 3×10^{-4} . The Adam optimizer with the same parameters as in the search for CNNs is used to optimize architecture parameters. After warming up for 9 epochs, SGAS begins to select one operation for a selected edge every 7 epochs. We experimented with both selection criteria, *Criterion 1* and *Criterion 2*. We use a history window of 4 for *Cri.2*. The

Architecture	Test Err. (%)		Params (M)	$\times +$ (M)	Search Cost (GPU-days)	Search Method
	top-1	top-5				
Inception-v1 [49]	30.2	10.1	6.6	1448	-	manual
MobileNet [20]	29.4	10.5	4.2	569	-	manual
ShuffleNet 2x (v1) [62]	26.4	10.2	~ 5	524	-	manual
ShuffleNet 2x (v2) [38]	25.1	-	~ 5	591	-	manual
NASNet-A [66]	26	8.4	5.3	564	1800	RL
NASNet-B [66]	27.2	8.7	5.3	488	1800	RL
NASNet-C [66]	27.5	9	4.9	558	1800	RL
AmoebaNet-A [43]	25.5	8	5.1	555	3150	evolution
AmoebaNet-B [43]	26	8.5	5.3	555	3150	evolution
AmoebaNet-C [43]	24.3	7.6	6.4	570	3150	evolution
FairNAS-A [10]	24.7	7.6	4.6	388	12	evolution
PNAS [33]	25.8	8.1	5.1	588	225	SMBO
MnasNet-92 [50]	25.2	8	4.4	388	-	RL
DARTS (2 nd order) [35]	26.7	8.7	4.7	574	4.0	gradient
SNAS (mild) [58]	27.3	9.2	4.3	522	1.5	gradient
ProxylessNAS [8]	24.9	7.5	7.1	465	8.3	gradient
P-DARTS [9]	24.4	7.4	4.9	557	0.3	gradient
BayesNAS [63]	26.5	8.9	3.9	-	0.2	gradient
PC-DARTS [61]	25.1	7.8	5.3	586	0.1	gradient
SGAS (Cri.1 avg.)	24.41 \pm 0.16	7.29 \pm 0.09	5.3	579	0.25	gradient
SGAS (Cri.1 best)	24.2	7.2	5.3	585	0.25	gradient
SGAS (Cri.2 avg.)	24.38 \pm 0.22	7.39 \pm 0.07	5.4	597	0.25	gradient
SGAS (Cri.2 best)	24.1	7.3	5.4	598	0.25	gradient

Table 2. **Comparison with state-of-the-art classifiers on ImageNet.** We transfer the top 3 performing architectures on CIFAR-10 to ImageNet in the mobile setting. $\times +$ denote multiply-add operations. The average and best performance of SGAS are reported.

Architecture	OA (%)	Params (M)	Search Cost (GPU-days)
3DmFV-Net [4]	91.6	45.77	manual
SpecGCN [54]	91.5	2.05	manual
PointNet++ [42]	90.7	1.48	manual
PCNN [3]	92.3	8.2	manual
PointCNN [31]	92.2	0.6	manual
DGCNN [55]	92.2	1.84	manual
KPConv [51]	92.9	14.3	manual
Random Search	92.65 \pm 0.33	8.77	random
SGAS (Cri.1 avg.)	92.69 \pm 0.20	8.78	0.19
SGAS (Cri.1 best)	92.87	8.63	0.19
SGAS (Cri.2 avg.)	92.93 \pm 0.19	8.87	0.19
SGAS (Cri.2 best)	93.23	8.49	0.19
SGAS (Cri.2 small best)	93.07	3.86	0.19

Table 3. **Comparison with state-of-the-art architectures for 3D object classification on ModelNet40.** 10 architectures are derived for both SGAS and Random Search within the same search space.

batch size increases by 4 after each decision. The search takes around 0.19 GPU day on one NVIDIA GTX 1080Ti.

4.2.2 Architecture Evaluation on ModelNet40

After searching for 10 architectures on ModelNet10, we form a large backbone network for each and train them on ModelNet40. The performance of 3D point cloud classification is evaluated with the overall accuracy (OA). We also

apply Random Search to the same search space to obtain 10 architectures as our random search baseline.

Training Settings. We stack the searched cell 9 times with channel size 128. We also form small networks by stacking the cell 3 times with the same channel size. We use $k = 20$ for all the large networks and $k = 9$ for the small ones. Adam is used to optimize the weights with initial learning rate 0.001 and weight decay 1×10^{-4} . We sample 1024 points as input. Our architectures are all trained for 400 epochs with batch size of 32. We report the mean and standard deviation of the accuracy on the test dataset of the 10 discovered architectures; we also report the accuracy of the best performing model of the big and the small networks.

Evaluation Results and Analysis. We compare the performance of our discovered architectures with SOTA hand-crafted methods and architectures obtained by Random Search for 3D point clouds classification on ModelNet40. Table 3 shows that SGAS (Cri.2 best), the best architecture discovered by our SGAS with *Criterion 2*, outperforms all the other models. The smaller network SGAS (Cri.2 small best) discovered by SGAS with *Criterion 2* also outperforms all the hand-crafted architectures. Owing to a well-designed search space, Random Search is a strong baseline. The performance of SGAS surpasses the hand-crafted architectures and Random Search, demonstrating the effectiveness of SGAS for GCN architecture search. The best architecture for this task can be found in Figure 4 (a).

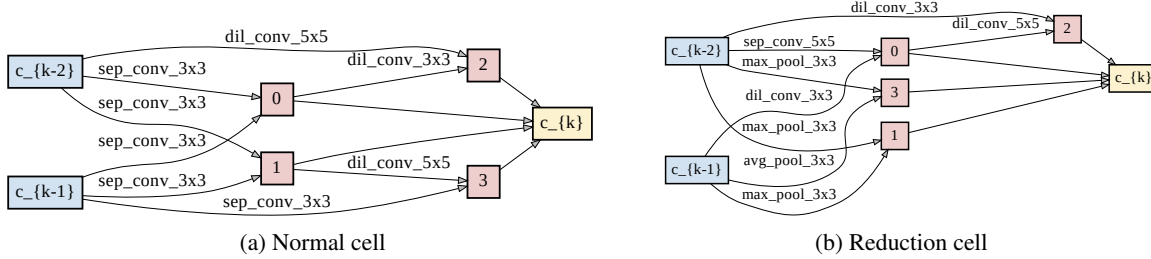


Figure 3. Best cell architecture on Imagenet with SGAS Crit. 2

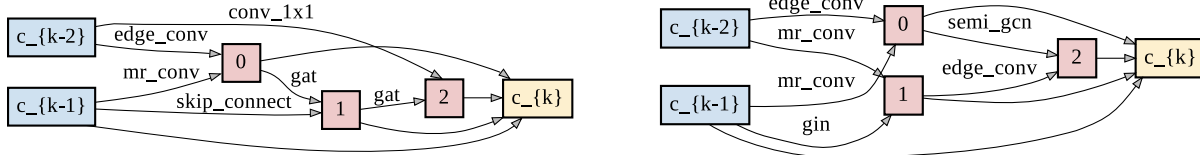


Figure 4. Best cell architectures on ModelNet40 and PPI

4.2.3 Architecture Search on PPI

PPI is a popular biological graph dataset in the data mining domain. We search for GCN architectures on the PPI dataset for the task of node classification.

Training Settings. We use 1 cell with 32 channels. We train and search the architectures for 50 epochs with a batch size of 6 on PPI. We do not increase the batch size after making decisions since PPI is small and only contains 20 batches. The other parameters are the same as when searching on ModelNet10. **The search takes around 0.003 day (4 minutes) on a Nvidia Tesla V100 GPU (16GB).**

4.2.4 Architecture Evaluation on PPI

We evaluate architectures on the PPI test set. We report the mean, standard derivation and the best accuracy and compare them with the SOTA methods and Random Search. We also conduct an ablation study on number of cells and channel size which we include in the **supplement material**.

Training Settings. We stack the discovered cell 5 times with channel size 512. Adam is used to optimize the model weights with initial learning rate 0.002. We use a cosine annealing to schedule the learning rate. Our architectures are trained for 2000 epochs with batch size of 1 as suggested in DeepGCNs [29]. We find the best model on the validation dataset and obtain the micro-F1 score on the test dataset.

Evaluation Results and Analysis. We compare the average and best performance of SGAS to other state-of-the-arts methods and Random Search on node classification on the PPI dataset. Table 4 shows the best architecture discovered by our SGAS outperforms the state-of-the-art DenseMRGCN-14 [29] by $\sim 0.03\%$ with ~ 30.24 M less parameters. The average performance of SGAS also surpasses the Random Search baseline consistently. In addition, SGAS (Cri.2 avg.) outperforms SGAS (Cri.1 avg.) in terms of both mean and standard deviation. This indicates

Architecture	micro-F1 (%)	Params (M)	Search Cost (GPU-days)
GraphSAGE (LSTM) [17]	61.2	0.26	manual
GeniePath [36]	97.9	1.81	manual
GAT [52]	97.3 \pm 0.2	3.64	manual
DenseMRGCN-14 [29]	99.43	53.42	manual
ResMRGCN-28 [29]	99.41	14.76	manual
Random Search	99.36 \pm 0.04	23.70	random
SGAS (Cri.1 avg.)	99.38 \pm 0.17	25.01	0.003
SGAS (Cri.1 best)	99.46	23.18	0.003
SGAS (Cri.2 avg.)	99.40 \pm 0.09	25.93	0.003
SGAS (Cri.2 best)	99.46	29.73	0.003
SGAS (small)	98.89	0.40	0.003

Table 4. Comparison with state-of-the-art architectures for node classification on PPI. SGAS (small) is the small network stacking the cell searched by SGAS (Cri.1).

that *Criterion 2* provides more stable results. We visualize the architecture with the best performance in Figure 4 (b).

5. Conclusion

In this work, we propose the Sequential Greedy Architectural Search (SGAS) algorithm to design architectures automatically for CNNs and GCNs. The bi-level optimization problem in NAS is solved in a greedy fashion using heuristic criteria which take the edge importance, the selection certainty and the selection stability into consideration. Such an approach alleviates the effect of the **degenerate search-evaluation correlation** problem and reflects the true ranking of architectures. As a result, architectures discovered by SGAS achieve state-of-the-art performance on CIFAR-10, ImageNet, ModelNet and PPI datasets.

Acknowledgments. This work was supported by the King Abdullah University of Science and Technology (KAUST) Office of Sponsored Research (OSR) through VCC funding. The authors thank the KAUST IBEX team for helping to optimize the training workloads on ImageNet.

References

- [1] Karim Ahmed and Lorenzo Torresani. Connectivity learning in multi-branch networks. *arXiv preprint arXiv:1709.09582*, 2017. [2](#)
- [2] Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient descent. In *Advances in neural information processing systems*, pages 3981–3989, 2016. [1](#)
- [3] Matan Atzmon, Haggai Maron, and Yaron Lipman. Point convolutional neural networks by extension operators. *ACM Trans. Graph.*, 37(4):71:1–71:12, July 2018. [7](#)
- [4] Yizhak Ben-Shabat, Michael Lindenbaum, and Anath Fischer. 3dmfv: Three-dimensional point cloud classification in real-time using convolutional neural networks. *IEEE Robotics and Automation Letters*, 3(4):3145–3152, 2018. [7](#)
- [5] Gabriel Bender, Pieter-Jan Kindermans, Barret Zoph, Vijay Vasudevan, and Quoc Le. Understanding and simplifying one-shot architecture search. In *International Conference on Machine Learning*, pages 549–558, 2018. [2](#)
- [6] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Smash: one-shot model architecture search through hypernetworks. *arXiv preprint arXiv:1708.05344*, 2017. [2](#)
- [7] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–970, 2015. [1](#)
- [8] Han Cai, Ligeng Zhu, and Song Han. Proxylessnas: Direct neural architecture search on target task and hardware. *arXiv preprint arXiv:1812.00332*, 2018. [2](#), [6](#), [7](#)
- [9] Xin Chen, Lingxi Xie, Jun Wu, and Qi Tian. Progressive differentiable architecture search: Bridging the depth gap between search and evaluation. *arXiv preprint arXiv:1904.12760*, 2019. [2](#), [3](#), [6](#), [7](#), [12](#)
- [10] Xiangxiang Chu, Bo Zhang, Ruijun Xu, and Jixiang Li. Fairnas: Rethinking evaluation fairness of weight sharing neural architecture search, 2019. [7](#)
- [11] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. [1](#)
- [12] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018. [1](#)
- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009. [1](#), [2](#), [5](#)
- [14] Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric. *arXiv preprint arXiv:1903.02428*, 2019. [13](#)
- [15] Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. Bilevel programming for hyperparameter optimization and meta-learning. *arXiv preprint arXiv:1806.04910*, 2018. [1](#)
- [16] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. [1](#)
- [17] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pages 1024–1034, 2017. [6](#), [8](#), [12](#), [13](#)
- [18] Chaoyang He, Haishan Ye, Li Shen, and Tong Zhang. Milenas: Efficient neural architecture search via mixed-level reformulation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#)
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [1](#), [2](#)
- [20] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. [7](#)
- [21] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018. [1](#)
- [22] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. [1](#), [2](#), [6](#)
- [23] Siyu Huang, Xi Li, Zhi-Qi Cheng, Zhongfei Zhang, and Alexander Hauptmann. Gnas: A greedy neural architecture search method for multi-attribute learning. In *2018 ACM Multimedia Conference on Multimedia Conference*, pages 2049–2057. ACM, 2018. [12](#)
- [24] Maurice G Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938. [3](#)
- [25] Thomas Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *ArXiv*, abs/1609.02907, 2016. [6](#), [12](#)
- [26] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009. [5](#)
- [27] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. [1](#), [2](#)
- [28] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. [1](#), [2](#)
- [29] Guohao Li, Matthias Müller, Guocheng Qian, Itzel C. Delgadillo, Abdullellah Abualshour, Ali K. Thabet, and Bernard Ghanem. Deepgcns: Making gcns go as deep as cnns. *ArXiv*, abs/1910.06849, 2019. [6](#), [8](#)
- [30] Guohao Li, Matthias Müller, Ali Thabet, and Bernard Ghanem. Deepgcns: Can gcns go as deep as cnns? In *The*

- IEEE International Conference on Computer Vision (ICCV)*, 2019. 6, 12, 13
- [31] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. In *NeurIPS*, 2018. 7
 - [32] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1
 - [33] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 19–34, 2018. 2, 6, 7, 12
 - [34] Hanxiao Liu, Karen Simonyan, Oriol Vinyals, Chrisantha Fernando, and Koray Kavukcuoglu. Hierarchical representations for efficient architecture search, 2017. 6
 - [35] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018. 1, 2, 3, 4, 5, 6, 7
 - [36] Ziqi Liu, Chaochao Chen, Longfei Li, Jun Zhou, Xiaolong Li, Le Song, and Yuan Qi. Geniepath: Graph neural networks with adaptive receptive paths. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4424–4431, 2019. 8
 - [37] Renqian Luo, Fei Tian, Tao Qin, Enhong Chen, and Tie-Yan Liu. Neural architecture optimization. In *Advances in neural information processing systems*, pages 7816–7827, 2018. 6
 - [38] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 116–131, 2018. 7
 - [39] Matthias Muller, Adel Bibi, Silvio Giancola, Salman Alsubaihi, and Bernard Ghanem. Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 300–317, 2018. 1
 - [40] Hieu Pham, Melody Y Guan, Barret Zoph, Quoc V Le, and Jeff Dean. Efficient neural architecture search via parameter sharing. *arXiv preprint arXiv:1802.03268*, 2018. 2, 6
 - [41] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 77–85, 2016. 12
 - [42] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in neural information processing systems*, pages 5099–5108, 2017. 7
 - [43] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4780–4789, 2019. 2, 4, 6, 7
 - [44] Shreyas Saxena and Jakob Verbeek. Convolutional neural fabrics. In *Advances in Neural Information Processing Systems*, pages 4053–4061, 2016. 2
 - [45] Christian Sciuto, Kaicheng Yu, Martin Jaggi, Claudiu Musat, and Mathieu Salzmann. Evaluating the search phase of neural architecture search. *arXiv preprint arXiv:1902.08142*, 2019. 3
 - [46] Richard Shin, Charles Packer, and Dawn Song. Differentiable neural network architecture search. 2018. 2
 - [47] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 1, 2
 - [48] Michael J Swain and Dana H Ballard. Color indexing. *International journal of computer vision*, 7(1):11–32, 1991. 4
 - [49] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Computer Vision and Pattern Recognition (CVPR)*, 2015. 1, 2, 7
 - [50] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2820–2828, 2019. 7
 - [51] Hugues Thomas, Charles R. Qi, Jean-Emmanuel Deschaud, Beatriz Marcotequi, François Goulette, and Leonidas J. Guibas. Kpconv: Flexible and deformable convolution for point clouds. *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 7
 - [52] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017. 6, 8, 12
 - [53] Tom Veniat and Ludovic Denoyer. Learning time/memory-efficient deep architectures with budgeted super networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3492–3500, 2018. 2
 - [54] Chu Wang, Babak Samari, and Kaleem Siddiqi. Local spectral graph convolution for point set feature learning. In *The European Conference on Computer Vision (ECCV)*, September 2018. 7
 - [55] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon. Dynamic graph cnn for learning on point clouds. *ACM Transactions on Graphics (TOG)*, 2019. 6, 7, 12
 - [56] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1912–1920, 2014. 5, 6
 - [57] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 1
 - [58] Sirui Xie, Hehui Zheng, Chunxiao Liu, and Liang Lin. Snas: stochastic neural architecture search. *arXiv preprint arXiv:1812.09926*, 2018. 2, 6, 7
 - [59] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *ArXiv*, abs/1810.00826, 2018. 6, 12

- [60] Mengmeng Xu, Chen Zhao, David S. Rojas, Ali Thabet, and Bernard Ghanem. G-tad: Sub-graph localization for temporal action detection, 2019. 6
- [61] Yuhui Xu, Lingxi Xie, Xiaopeng Zhang, Xin Chen, Guo-Jun Qi, Qi Tian, and Hongkai Xiong. Pc-darts: Partial channel connections for memory-efficient differentiable architecture search. *arXiv preprint arXiv:1907.05737*, 2019. 2, 3, 6, 7
- [62] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6848–6856, 2018. 7
- [63] Hongpeng Zhou, Minghao Yang, Jun Wang, and Wei Pan. Bayesnas: A bayesian approach for neural architecture search, 2019. 6, 7
- [64] Marinka Zitnik and Jure Leskovec. Predicting multicellular function through multi-layer tissue networks. In *Bioinformatics*, 2017. 5
- [65] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016. 1, 2, 4
- [66] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018. 2, 3, 6, 7

SGAS: Sequential Greedy Architecture Search

– Supplementary Material –

A. Discussion

A.1. Greedy Methods in NAS

The idea of incorporating greedy algorithms into NAS has been explored in several works. PNAS [33] proposes a sequential model-based optimization (SMBO) approach to accelerate the search for CNN architectures. They start from a simple search space and learn a predictor function. Then they greedily grow the search space by predicting scores of candidate cells using the learned predictor function. GNAS [23] learns a global tree-like architecture for multi-attribute learning by iteratively updating layer-wise local architectures in a greedy manner. P-DARTS [9] can also be regarded as a greedy approach, in which they bridge the depth gap between search and evaluation by gradually increasing the depth of the search networks while shrinking the number of candidate operations.

A.2. Selection Criteria and Hyper-parameters

Edge Importance and Selection Certainty. Edge importance and selection certainty are combined into a single criterion, since the algorithm will be agnostic to the selection distribution of an edge, if we only consider edge importance. In this case, an edge may be selected with a sub-optimal operation at early epochs. On the other hand, we need to select 8 out of 14 edges in a DAG with 4 intermediate nodes for a fair comparison with DARTS. Only considering selection certainty may fail to select the optimal 8 edges, since an edge with a high selection certainty may have a high weight on *Zero* operation (low edge importance).

Choices of Hyper-parameters. Three extra parameters are introduced in SGAS: (1) *length of warm-up phase* (2) *interval of greedy decisions* (3) *history window size for Cri.2*. We provide a discussion on the default choices of them:

(1) Since the softmax weights of operations are initialized under a uniform distribution, choosing an operation for an edge after a short period of warming up leads to stable results. We simply set the *length of the warm-up phase* to 9 epochs so that the first greedy decision will be made at the 10th epoch. (2) For CNN experiments, the *interval between greedy decisions* is chosen to be 5. Since designing a normal cell with 4 intermediate nodes needs to select 8 out of 14 edges (8 decisions to be made). For a fair comparison to our baseline DARTS, we want the search phase to last up to 50 epochs, which is the length of search epochs in DARTS. For GCN architectures, in order to learn a compact

network, we search a normal cell with 3 intermediate nodes. Thus, we have 6 decisions to make (6 out of 9 edges). Similarly, to keep the length of the search phase less than 50 epochs, we set the *interval between greedy decisions* to be 7. (3) The *history window size for Cri.2* is always set as 4, which is simply chosen to be slightly smaller than the *interval between greedy decisions*.

Ablation Study on Hyper-parameters. In order to better understand the effects of the choices of hyper-parameters, we conduct ablation studies on *interval of greedy decisions* T and *history window size* K for *Cri.2* on CIFAR-10 in Table 5. The default values of T and K are 5 and 4 respectively. We find that larger T and K stabilize the search and produce standard deviations in the test error. The test error only has a standard deviation as 0.08 when $T = 7$. When $T = 3$, the average test error increases significantly from 2.67% to 2.86%. We also find K is less sensitive than T .

T	K	Avg.		Best	
		Params (M)	Test Err.(%)	Params (M)	Test Err.(%)
5	4	3.91±0.22	2.67±0.21	4.09	2.44
3	4	4.09±0.24	2.86±0.12	4.39	2.69
7	4	3.66±0.16	2.65±0.08	3.68	2.54
5	2	3.87±0.20	2.73±0.16	3.94	2.51
5	6	3.93±0.26	2.67±0.17	3.70	2.47

Table 5. Ablation study on *interval of greedy decisions* T and *history window size* K for *Cri.2* on CIFAR-10. We use SGAS (Cri.2) as our search method. We report the average and best performance of searched architectures.

B. Experimental Details

B.1. GCN Experiments

GCN operators. Similar as the search for CNN, SGAS selects one operation from a candidate operation search space for each edge in the DAG. We choose the following 10 operations as our candidate operations: *conv-1×1*, *MRConv* [30], *EdgeConv* [55], *GAT* [52], *SemiGCN* [25], *GIN* [59], *SAGE* [17], *RelSAGE*, *skip-connect*, and *zero* operation. *conv-1×1* is a basic convolution operation without aggregating information from neighbors, which is similar to PointNet [41]. *MRConv* [30], *EdgeConv* [55], *GAT* [52], *SemiGCN* [25], *GIN* [59] and *SAGE* [17] are widely used GCN operators in the graph learning domain and the 3D computer vision domain. *RelSAGE* is a modified Graph-

SAGE (SAGE) [17] which combines the ideas from *MR-Conv* [30] and GraphSAGE [17]. Instead of aggregating the node features with its neighbor features directly, we aggregate the node features with the difference between the node features and its neighbor features:

$$\mathbf{h}_v^{(k)} = \sigma \left(\mathbf{W}^{(k)} \cdot f_k \left(\mathbf{h}_v^{(k-1)}, \left\{ \mathbf{h}_u^{(k-1)} - \mathbf{h}_v^{(k-1)}, \forall u \in \mathcal{N}(v) \right\} \right) \right)$$

where $\mathbf{h}_v^{(k)}$ is the feature of the center node v in k -th layer. $\mathcal{N}(v)$ denotes the neighbors of node v . f_k is a max aggregation function and σ is a ReLU activation function as GraphSAGE [17]. The GCNs operators are implemented using Pytorch Geometric [14]. We also add *skip-connect* (similar as residual graph connections [30]) and *zero* operation in our search space.

Ablation Study on GCN Cells. We conduct an ablation study on the parameter size of the best cell searched on PPI by SGAS (Cri.1 best). Table 6 shows the trade-off between the parameter size and the final performance. To derive a compact model, we can use a smaller number of cells or less channels in the architecture searched by SGAS.

Number of Cells	Channel Size	Params (M)	micro-F1(%)
5	64	0.40	98.894
5	128	1.52	99.369
5	256	5.89	99.429
5	512	23.18	99.462
1	256	1.22	99.157
3	256	3.52	99.418
5	256	5.89	99.429
7	256	8.25	99.433

Table 6. **Ablation study on channel size and number of cells on node classification on PPI.** We use SGAS (Cri.1 best), the best architecture we discovered by using *Criterion 1* to conduct experiments.

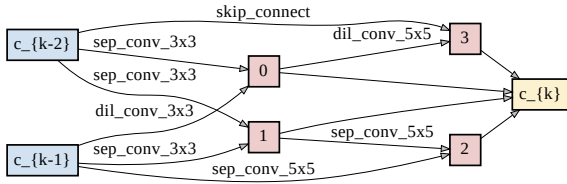
B.2. More Details

Cell Visualizations. We visualize the best cells discovered by SGAS with different criteria (*Criterion 1* and *Criterion 2*) mentioned in the experiment section. Figure 5 shows the best cells for CNNs on CIFAR-10 and ImageNet. Figure 6 shows the best cells for GCNs on ModelNet-40 and PPI.

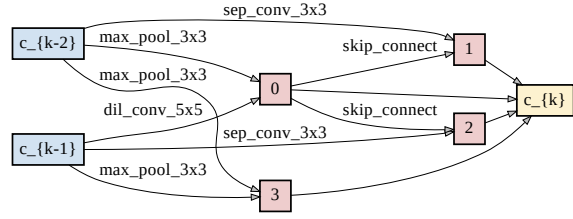
Detailed results. Here we provide the detailed results mentioned in the experimental section of the paper. In the CNN experiments, we compare SGAS with DARTS on CIFAR-10 and ImageNet. We execute the search phase 10 times for both SGAS (Cri.1 and Cri.2) and DARTS (1st and 2nd order) to obtain 10 different architectures per method. For each resulting architecture, we run the evaluation phase and assign a ranking based on the evaluation accuracy. To measure the discrepancy between the search and evaluation, we calculate the Kendall τ correlation between the ranking of the search phase and the evaluation phase. We show these

results in Table 7 and Table 8 for SGAS, and Table 9 and Table 10 for DARTS. For ImageNet, we evaluate the top three architectures found on CIFAR-10. We show the results in Table 11 and Table 12 for both *Criterion 1* and *Criterion 2*.

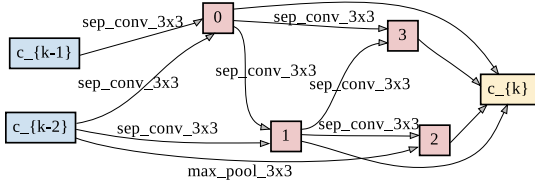
In the GCN experiments, we compare SGAS (Cri.1 and Cri.2) with a random search baseline on ModelNet and PPI. Similar as in experiments for CNNs, we conduct the search phase 10 times for each method. For experiments on ModelNet, we search cells on ModelNet10 and then evaluate the searched cells on ModelNet40. The results are shown for *Criterion 1*, *Criterion 2* and random search in Table 13, Table 14 and Table 15 respectively. The results on PPI are presented in Table 16 and Table 17 for each *Criterion* and in Table 18 for random search.



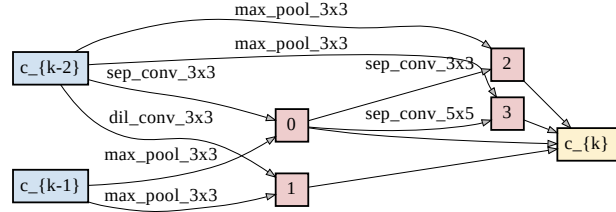
(a) Normal cell of the best model with SGAS *Cri. 1* on CIFAR-10 and ImageNet



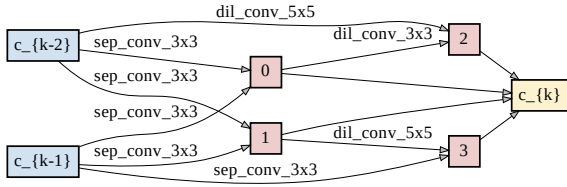
(b) Reduction cell of the best model with SGAS *Cri. 1* on CIFAR-10 and ImageNet



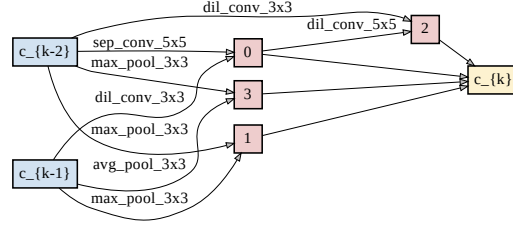
(c) Normal cell of the best model with SGAS *Cri. 2* on CIFAR-10



(d) Reduction cell of the best model with SGAS *Cri. 2* on CIFAR-10

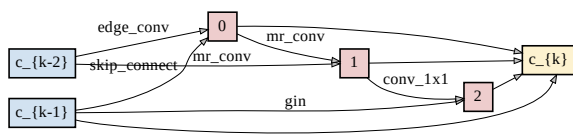


(e) Normal cell of the best model with SGAS *Cri. 2* on ImageNet

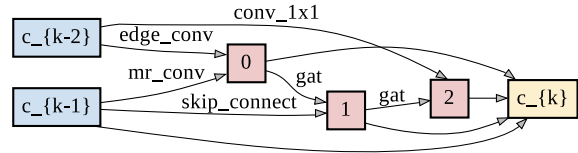


(f) Reduction cell of the best model with SGAS *Cri. 2* on ImageNet

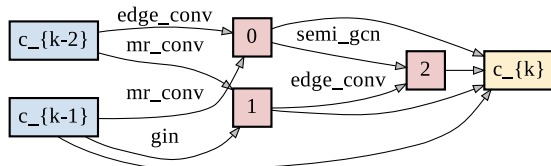
Figure 5. Best cell architecture for image classification tasks



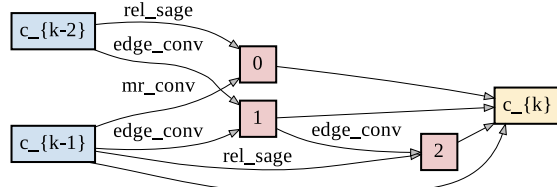
(a) Normal cell of the best model with SGAS *Cri. 1* on ModelNet40



(b) Normal cell of the best model with SGAS *Cri. 2* on ModelNet40



(c) Normal cell of the best model with SGAS *Cri. 1* on PPI



(d) Normal cell of the best model with SGAS *Cri. 2* on PPI

Figure 6. Best cell architectures on ModelNet40 and PPI with each *Criterion*

Experiment	Validation error (%)	Params (M)	Test error (%)	Evaluation ranking
Cri.1_CIFAR_1	16.94	3.75	2.44	2
Cri.1_CIFAR_2	17.33	3.73	2.50	3
Cri.1_CIFAR_3	17.90	3.80	2.39	1
Cri.1_CIFAR_4	17.90	3.32	2.63	6
Cri.1_CIFAR_5	17.99	3.45	2.78	8
Cri.1_CIFAR_6	18.43	3.47	2.68	7
Cri.1_CIFAR_7	18.72	3.83	2.51	4
Cri.1_CIFAR_8	19.82	3.66	2.61	5
Cri.1_CIFAR_9	19.93	3.98	3.18	10
Cri.1_CIFAR_10	21.53	3.61	2.87	9
Average	18.65±1.4	3.66±0.2	2.66±0.24	Kendall τ
Best Model	17.90	3.80	2.39	0.56

Table 7. Results of SGAS *Criterion 1* on CIFAR-10

Experiment	Validation error (%)	Params (M)	Test error (%)	Evaluation ranking
Cri.2_CIFAR_1	16.48	4.14	2.57	4
Cri.2_CIFAR_2	17.26	3.88	2.60	6
Cri.2_CIFAR_3	17.31	4.09	2.44	1
Cri.2_CIFAR_4	17.47	3.91	2.49	2
Cri.2_CIFAR_5	17.53	3.69	2.52	3
Cri.2_CIFAR_6	17.98	3.95	3.12	10
Cri.2_CIFAR_7	18.28	3.69	2.58	5
Cri.2_CIFAR_8	18.28	4.33	2.85	8
Cri.2_CIFAR_9	19.48	3.73	2.85	9
Cri.2_CIFAR_10	19.98	3.68	2.66	7
Average	18.00±1.06	3.91±0.22	2.67±0.21	Kendall τ
Best Model	17.31	4.09	2.44	0.42

Table 8. Results of SGAS *Criterion 2* on CIFAR-10

Experiment	Validation error (%)	Params (M)	Test error (%)	Evaluation ranking
DARTS_1st_CIFAR_1	11.37	3.27	2.83	4
DARTS_1st_CIFAR_2	11.45	3.65	2.57	2
DARTS_1st_CIFAR_3	11.47	2.29	2.94	7
DARTS_1st_CIFAR_4	11.48	2.65	2.96	8
DARTS_1st_CIFAR_5	11.65	3.09	2.50	1
DARTS_1st_CIFAR_6	11.75	2.86	2.84	5
DARTS_1st_CIFAR_7	11.77	2.09	3.06	10
DARTS_1st_CIFAR_8	11.81	2.52	3.01	9
DARTS_1st_CIFAR_9	11.82	2.65	2.94	6
DARTS_1st_CIFAR_10	11.94	3.27	2.82	3
Average	11.65±0.19	2.84±0.49	2.85±0.18	Kendall τ
Best Model	11.65	3.09	2.50	0.16

Table 9. Results of DARTS 1st order on CIFAR-10

Experiment	Validation error (%)	Params (M)	Test error (%)	Evaluation ranking
DARTS_2nd_CIFAR_1	11.35	2.91	2.96	8
DARTS_2nd_CIFAR_2	11.51	2.93	2.73	5
DARTS_2nd_CIFAR_3	11.68	2.20	3.01	9
DARTS_2nd_CIFAR_4	11.76	2.66	2.75	6
DARTS_2nd_CIFAR_5	11.80	3.09	2.72	4
DARTS_2nd_CIFAR_6	11.82	3.40	2.62	3
DARTS_2nd_CIFAR_7	11.83	2.91	2.82	7
DARTS_2nd_CIFAR_8	11.93	3.20	2.51	1
DARTS_2nd_CIFAR_9	11.95	2.14	3.48	10
DARTS_2nd_CIFAR_10	12.03	2.55	2.62	2
Average	11.77 \pm 0.21	2.8 \pm 0.41	2.82 \pm 0.28	Kendall τ
Best Model	11.93	3.20	2.51	-0.29

Table 10. Results of DARTS 2nd order on CIFAR-10

Experiment	Test error top-1 (%)	Test error top-5 (%)	Params (M)	$\times +$
Cri.1_ImageNet_1	24.47	7.23	5.25	578
Cri.1_ImageNet_2	24.53	7.40	5.23	574
Cri.1_ImageNet_3	24.22	7.25	5.29	585
Average	24.41 \pm 0.16	7.29 \pm 0.09	5.25 \pm 0.03	579
Best Model	24.22	7.25	5.29	585

Table 11. Results of SGAS with *Criterion 1* on ImageNet. Note that the chosen architectures are the three best ones from the results obtained on CIFAR-10.

Experiment	Test error top-1 (%)	Test error top-5 (%)	Params (M)	$\times +$
Cri.2_ImageNet_1	24.44	7.41	5.70	621
Cri.2_ImageNet_2	24.13	7.31	5.44	598
Cri.2_ImageNet_3	24.55	7.44	5.20	571
Average	24.38 \pm 0.22	7.39 \pm 0.07	5.44 \pm 0.25	597
Best Model	24.13	7.31	5.44	598

Table 12. Results of SGAS with *Criterion 2* on ImageNet. Note that the chosen the architectures are the three best ones from the results obtained on CIFAR-10.

Experiment	Params (M)	Test OA (%)
Cri.1_ModelNet_1	8.79	92.71
Cri.1_ModelNet_2	9.23	92.83
Cri.1_ModelNet_3	8.79	92.79
Cri.1_ModelNet_4	8.78	92.34
Cri.1_ModelNet_5	8.93	92.79
Cri.1_ModelNet_6	8.19	92.30
Cri.1_ModelNet_7	8.63	92.83
Cri.1_ModelNet_8	8.63	92.71
Cri.1_ModelNet_9	8.63	92.87
Cri.1_ModelNet_10	9.23	92.79
Average	8.78 \pm 0.30	92.69 \pm 0.20
Best Model	8.63	92.87

Table 13. Results of SGAS with *Criterion 1* on ModelNet40. Architectures are formed by stacking 9 cells with 128 channel size.

Experiment	Params (M)	Test OA (%)
Cri.2_ModelNet_1	8.78	92.91
Cri.2_ModelNet_2	8.78	92.67
Cri.2_ModelNet_3	9.08	92.79
Cri.2_ModelNet_4	8.49	93.23
Cri.2_ModelNet_5	9.08	93.03
Cri.2_ModelNet_6	9.08	93.07
Cri.2_ModelNet_7	8.78	93.11
Cri.2_ModelNet_8	8.63	92.67
Cri.2_ModelNet_9	8.78	92.83
Cri.2_ModelNet_10	9.23	92.95
Average	8.87 \pm 0.23	92.93 \pm 0.19
Best Model	8.49	93.23

Table 14. **Results of SGAS with *Criterion 2* on ModelNet40.** Architectures are formed by stacking 9 cells with 128 channel size.

Experiment	Params (M)	Test OA (%)
Random_ModelNet_1	9.22	92.79
Random_ModelNet_2	8.93	92.67
Random_ModelNet_3	9.08	92.71
Random_ModelNet_4	8.78	92.46
Random_ModelNet_5	8.19	92.79
Random_ModelNet_6	8.63	92.54
Random_ModelNet_7	8.93	91.94
Random_ModelNet_8	8.63	92.99
Random_ModelNet_9	8.79	93.15
Random_ModelNet_10	8.49	92.46
Average	8.77 \pm 0.30	92.65 \pm 0.33
Best Model	8.79	93.15

Table 15. **Results of random search on ModelNet40.** Architectures are formed by stacking 9 cells with 128 channel size.

Experiment	Params (M)	Test micro-F1 (%)
Cri.1_PPI_1	27.11	99.45
Cri.1_PPI_2	23.18	99.42
Cri.1_PPI_3	25.80	98.91
Cri.1_PPI_4	25.80	99.38
Cri.1_PPI_5	24.49	99.44
Cri.1_PPI_6	29.73	99.44
Cri.1_PPI_7	24.50	99.44
Cri.1_PPI_8	21.87	99.43
Cri.1_PPI_9	24.49	99.44
Cri.1_PPI_10	23.18	99.46
Average	25.01 \pm 2.24	99.38 \pm 0.17
Best Model	23.18	99.46

Table 16. **Results of SGAS with *Criterion 1* on PPI.**

Experiment	Params (M)	Test micro-F1 (%)
Cri.2_PPI_1	25.80	99.17
Cri.2_PPI_2	28.42	99.46
Cri.2_PPI_3	20.55	99.40
Cri.2_PPI_4	21.87	99.43
Cri.2_PPI_5	24.49	99.44
Cri.2_PPI_6	28.42	99.43
Cri.2_PPI_7	29.73	99.42
Cri.2_PPI_8	25.80	99.45
Cri.2_PPI_9	28.42	99.44
Cri.2_PPI_10	25.79	99.41
Average	25.93 \pm 2.99	99.40 \pm 0.09
Best Model	28.42	99.46

Table 17. Results of SGAS with *Criterion 2* on PPI.

Experiment	Params (M)	Test micro-F1 (%)
Random_PPI_1	20.57	99.27
Random_PPI_2	24.48	99.37
Random_PPI_3	24.49	99.40
Random_PPI_4	19.24	99.40
Random_PPI_5	24.48	99.37
Random_PPI_6	21.85	99.36
Random_PPI_7	27.11	99.32
Random_PPI_8	23.17	99.40
Random_PPI_9	24.48	99.39
Random_PPI_10	27.11	99.38
Average	23.7 \pm 2.56	99.36 \pm 0.04
Best Model	23.17	99.40

Table 18. Results of random search on PPI.