

Single-level Optimization For Differential Architecture Search

Pengfei Hou

houpengfei2020@126.com

Ying Jin

jiny18@mails.tsinghua.edu.cn

Abstract

In this paper, we point out that differential architecture search (DARTS) makes gradient of architecture parameters biased for network weights and architecture parameters are updated in different datasets alternatively in the bi-level optimization framework. The bias causes the architecture parameters of non-learnable operations to surpass that of learnable operations. Moreover, using softmax as architecture parameters' activation function and inappropriate learning rate would exacerbate the bias. As a result, it's frequently observed that non-learnable operations are dominated in the search phase. To reduce the bias, we propose to use single-level to replace bi-level optimization and non-competitive activation function like sigmoid to replace softmax. As a result, we could search high-performance architectures steadily. Experiments on NAS Benchmark 201 validate our hypothesis and stably find out nearly the optimal architecture. On DARTS space, we search the state-of-the-art architecture with 77.0% top1 accuracy (training setting follows PDARTS and without any additional module) on ImageNet-1K and steadily search architectures up to 76.5% top1 accuracy (but not select the best from the searched architectures) which is comparable with current reported best result.

1. Introduction

Neural architecture search (NAS) has helped to find more optimal architecture than manual design. Generally NAS is formulated as a bi-level optimization problem[2] as:

$$\begin{aligned} \alpha^* &= \arg \min_{\alpha \in \mathbf{A}} \mathcal{L}_{val}(\alpha, w_\alpha^*) \\ s.t. \ w_\alpha^* &= \arg \min_{w_\alpha} \mathcal{L}_{train}(\alpha, w) \end{aligned} \quad (1)$$

where α denote architecture and \mathbf{A} denote architecture search space, w_α denote the network weights bound with the architecture α , \mathcal{L}_{train} and \mathcal{L}_{val} denote optimization loss on training and validation dataset. Due to inner optimization on \mathcal{L}_{train} that any architecture has to be trained fully, therefore it costs huge computation sources to search the

optimal architecture. To avoid training each architecture from scratch, weight-sharing methods [28] are proposed to construct a super network where all architectures share the same weights. DARTS relaxes the search space to be continuous and approximates w_α^* by adapting w using only a single training step, without solving the inner optimization completely. The approximation scheme is as follows:

$$\nabla_\alpha \mathcal{L}_{val}(\alpha, w_\alpha^*) \approx \nabla_\alpha \mathcal{L}_{val}(\alpha, w - \xi \nabla_w \mathcal{L}_{train}(\alpha, w)) \quad (2)$$

It saves computation costs a lot and finds out competitive architectures [8, 35].

However, many papers [35, 27, 38, 22, 12, 15] have reported that DARTS easily converges to non-learnable operations including skip-connnet, pooling and zero, etc. and it could not search steadily high-performance architectures. Furthermore we find out that non-learnable operations are dominated in the very early stage of search phase. And once the situation happens, it's likely to last until the end of the search phase 1.

For the phenomenon, we propose the hypothesis that it's caused from two aspects and try to give theoretic explanation. For one thing, the approximation to bi-level optimization which has to train w and α on different datasets computes biased gradients for α . Even the gradients of w and α are computed on the same dataset but not on the same batch, the bias also exists. On the early stage of training process, learnable operations' α could not be learned normally. As a result, non-learnable operations' α would surpass over learnable operations'. For the other thing, using softmax as α 's activation function and inappropriate (in specific is too large) learning rate would make surpassing architecture parameters larger. It means under the bi-level optimization framework the non-learnable operations would dominate over learnable operations from the start to finish.

In view of the above two aspects, we conduct experimental analysis to verify them based on NAS-Benchmark-201[15]. Furthermore, we give improvement on DARTS. On one side, we propose to use single-level optimization instead of bi-level optimization and update w and α on

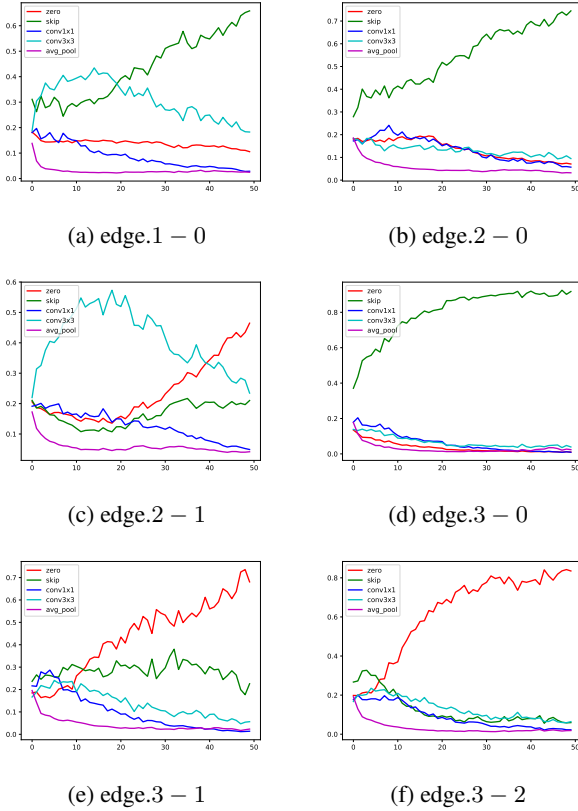


Figure 1: Values of softmax(α)(probability) on NAS-Benchmark-201. Use bi-level optimization(DARTS) to train for 50 epochs on CIFAR10. Non-learnable operations’ probability is far bigger than learnable operations’. It shows that non-learnable operations are likely dominated in the early stage and the situation would last until the end.

the same data batch. On another side, we propose to use uncompetitive activation function like sigmoid to replace softmax. And if there are too many non-operations in the searched result, we advice to decrease learning rate. As a result, we steadily search high-performance architectures. Experiments on NAS-Benchmark-201[15] search nearly the optimal architecture in the space with variance of 0. On DARTS space, we search the SOTA architecture with top1 accuracy of 77% on ImageNet-1K [21] for the first time. And we could search architectures steadily with top1 accuracy up to 76.5% which is comparable with reported best result 76.6% [19, 37]. The training phase is the same as PDARTS[27] and without any additional module.

In general, our contributions are as follows:

- We find out that non-learnable operations dominate learnable operations from the start to finish. According to the phenomenon, we propose a hypothesis that

bi-level optimization bias the gradients of architecture parameters. Using softmax as activation function and inappropriate learning rate would make the bias more serious. We give a theoretical explanation and conduct experiments to verify it.

- As a substitute, we propose to use single-level optimization which calculates and backward the gradients of w and α on the same data batch as well as at the same time:

$$\alpha^t, w^t \leftarrow \eta \nabla_{\alpha, w} \mathcal{L}_{train}(\alpha^{t-1}, w^{t-1}) \quad (3)$$

And we use uncompetitive activation function like sigmoid to replace softmax. Meanwhile, we do normalization for the initialization of sigmoid. More details are introduced bellow.

- Our improvement on DARTS could search high-performance architectures with little variance. On NAS-Benchmark-201 we find out nearly the optimal architecture. On DARTS space, we find the SOTA architecture which get 77.0% top1 accuracy on ImageNet. And we steadily find architectures (not select the best architecture of several tries) up to 76.5% top1 accuracy which is comparable with current reported best result.

2. Related Works

Neural Architecture Search. Neural architecture search(NAS) is an automatic method to design neural architecture instead of human design. Early NAS methods adopt reinforcement learning (RL) or evolutionary strategy [39, 2, 3, 32, 31, 40] to search among thousands of individually trained completely networks, which costs huge computation sources. Recent works focus on efficient weight-sharing methods, which falls into two categories: one-shot approaches [6, 4, 1, 7, 18, 34, 30] and gradient-based approaches [33, 28, 9, 8, 20, 12, 35, 24], both achieve state-of-the-art results on a series of tasks [10, 17, 25, 36, 16, 29] in various search space. They construct a super network/graph which share weights with all sub network/graphs. The former commonly does heuristic search and evaluation sampled architectures in the super network to get the optimal architecture. The latter relaxes the search space to be continuous and introduce differential and learnable architecture parameters. In the this paper, we mainly argue gradient-based approaches.

Differential Architecture Search. Gradient-based approaches are commonly formulated as an approximation to the bi-level optimization which updates network

weights and updates architecture parameters in the training and validation dataset alternatively[28]. Although it has searched competitive architectures, however, many papers pointed out DARTS-based methods don't work stably. [15] show DARTS perform badly on NAS-Benchmark-201 and performance drops fast during search. There are many works trying to give explanation and solve it. [4] show the relationship between DARTS searched architectures' performance and the domain eigenvalue of $\nabla_{\alpha}^2 \mathcal{L}_{valid}$. [12] show darts easily converge to skip-connect operation. [27, 35, 19, 22] show DARTS also easily over-converge to non-parametric operations. [12] introduces collaborative competition approach which offer each operation an independent architecture weight to avoid competition between operations. [19] observes the co-adaption problem and Matthew effect that operations with less parameters are trained maturely earlier. However, we indicate that Matthew effect is mainly caused by gradient bias in bi-level optimization but not co-adaption problem. In fact, in single-level optimization, the Matthew effect would be disappeared.

Single-level Optimization. Recently there are some works adopt single/one-level optimization to replace bi-level optimization and get much more stable results. Although in the original paper [28] it shows that single-level performs worse than bi-level optimization and they indicate single-level would cause overfitting, but recent works show the inverse results [22, 5, 19]. [22] advocates the overfitting phenomenon is caused by network's depth gap between search phase and evaluation phase and get much more stable results than bi-level-based DARTS. [5] indicates that super network parameters are trained much more effective than architecture parameters and do data augmentation in search phase. And based on single-level optimization, [19] combines operations dropout with single-level optimization to solve co-adaptation problem in the paper. However, on one side, these works focus on other technology to get satisfactory results but single-level optimization is as basis. Our experiments demonstrate that stand-alone single-level optimization is enough to search steadily high-performance architectures. On the other side, these works show single-level better than bi-level optimization mainly by empirical experiments but not give clear explanation on why bi-level doesn't work. In this paper, we propose that the reason is from biased gradients caused by calculating gradients on different data batch in the bi-level optimization framework.

3. Hypothesis And Methodology

In this section, we propose the hypothesis of the gradient bias of architecture parameters in the bi-level optimization framework and try to give theoretic explanation on the domination phenomenon of non-learnable operations. In out-

line, bi-level optimization has to compute gradients of network weights w and architecture parameters α separably on the training and validation dataset. It's the computation on different dataset that makes the gradients of learnable operations' α biased. Moreover, the bias leads non-learnable operations' α accumulating gradient faster than learnable operations'. With softmax as α 's activation function, the gap between them would be expanded. Finally, it results in the domination phenomenon.

3.1. Preliminary of Differential Architecture Search

Following [28] we search the DARTS space which stacks of several cells, where each cell is a directed acyclic graph. Each cell consists of sequential nodes where each node represent latent feature map x^i . The edge from node i to node j represent a connection and is bounded with one from candidate operations \mathbb{O} : convolution, pooling, zero, skip-connect, etc. Let $o^{i,j}$ denote operation from node i to j . Each intermediate node is computed based on all of its predecessors:

$$x_j = \sum_{i < j} o^{i,j}(x_i) \quad (4)$$

The problem of searching the best architecture is transformed to search the best operation $o_{k^*}^{i,j}$ among all cells and edges. DARTS make the search space continuous by relaxing the categorical choice of a particular operation to a softmax over all possible operations. Let $\alpha_{o_k}^{i,j}$ or $\alpha_k^{i,j}$ denote bounded parameter with operation o_k

$$f_{i,j}(x_i) = \sum_{o \in \mathbb{O}} \frac{\exp(\alpha_o^{i,j})}{\sum_{o' \in \mathbb{O}} \exp(\alpha_{o'}^{i,j})} o(x_i) \quad (5)$$

Then the intermediate node is computed based on all of its predecessors:

$$x_j = \sum_{i < j} f_{i,j}(x_i) \quad (6)$$

Without loss of generality, we consider a single edge. Loss function could be seen as

$$\mathcal{L} = f\left(\sum_i \frac{\exp(\alpha_i)}{\sum_j \exp(\alpha_j)} o_i(x); X, y\right) \quad (7)$$

3.2. Hypothesis

Prior works point out that DARTS easily converge to non-parametric operations like skip-connect, zero, pooling, etc. In fact, we find that in the early stage of search phase, non-learnable operations' architecture parameters have suppressed learnable operations'. Once the situation happens, it's hard to be retrieved 1. For the phenomenon, we have the following explanation.

Let p_i denote $\frac{\exp(\alpha_i)}{\sum_j \exp(\alpha_j)}$, assuming $o_i(x)$ could be expanded as $W_i x$, we have

$$\frac{\partial \mathcal{L}}{\partial p_i} = \mathbb{E}[\frac{\partial \mathcal{L}^T}{\partial \bar{x}}(W_i x); X, y] \quad (8)$$

$$\frac{\partial \mathcal{L}}{\partial W_i} = \mathbb{E}[p_i \frac{\partial \mathcal{L}}{\partial \bar{x}} x^T; X, y] \quad (9)$$

Considering $\frac{\partial \mathcal{L}}{\partial p_i}$, if $W_i x$ is negative correlated with $\frac{\partial \mathcal{L}}{\partial \bar{x}}$, p_i would be increased under gradient descent. Heuristically, if W_i fit x better than W_j , then $\frac{\partial \mathcal{L}}{\partial p_i}$ should be smaller than $\frac{\partial \mathcal{L}}{\partial p_j}$. For DARTS, p_i is updated on the validation dataset. On the iteration t and on validation dataset $\mathcal{D}_{\text{valid}} = \{X_{\text{val}}, y_{\text{val}}\}$, we have:

$$\frac{\partial \mathcal{L}}{\partial p_i^t} = \frac{1}{N} \sum_{j=1}^N \frac{\partial \mathcal{L}(x_{\text{val}}^j)}{\partial \bar{x}} (W_i^t x_{\text{val}}^j)^T \quad (10)$$

W_i^t is updated on the training dataset $\mathcal{D}_{\text{train}} = \{X_{\text{train}}, y_{\text{train}}\}$, we have:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_i^t} &= \frac{p_i}{M} \sum_{k=1}^M \frac{\partial \mathcal{L}(x_{\text{train}}^k)}{\partial \bar{x}} x_{\text{train}}^k{}^T \\ W_i^t &= W_i^{t-1} - \eta \frac{\partial \mathcal{L}}{\partial W_i^{t-1}} \\ &= W_i^{t-1} - \eta \frac{p_i^{t-1}}{M} \sum_{k=1}^M \frac{\partial \mathcal{L}(x_{\text{train}}^k)}{\partial \bar{x}} x_{\text{train}}^k{}^T \end{aligned} \quad (11)$$

As a result, the gradient of p_i is:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial p_i^t} &= \frac{1}{N} \sum_{j=1}^N \frac{\partial \mathcal{L}(x_{\text{val}}^j)}{\partial \bar{x}} (W_i^{t-1} x_{\text{val}}^j)^T \\ &\quad - \frac{\eta p_i^{t-1}}{NM} \sum_{j=1}^N \sum_{k=1}^M \left(\frac{\partial \mathcal{L}(x_{\text{val}}^j)}{\partial \bar{x}} \frac{\partial \mathcal{L}(x_{\text{train}}^k)}{\partial \bar{x}} \right) (x_{\text{train}}^k{}^T x_{\text{val}}^j) \end{aligned} \quad (12)$$

For bi-level optimization W and α are computed on the different batches, if samples are different $\frac{\partial \mathcal{L}(x_{\text{val}}^j)}{\partial \bar{x}}$ is independent $\frac{\partial \mathcal{L}(x_{\text{train}}^k)}{\partial \bar{x}}$ and x_{train}^k is independent of x_{val}^j . Thus we have

$$\sum_{j=1}^N \sum_{k=1}^M \left(\frac{\partial \mathcal{L}(x_{\text{val}}^j)}{\partial \bar{x}} \frac{\partial \mathcal{L}(x_{\text{train}}^k)}{\partial \bar{x}} \right) (x_{\text{train}}^k{}^T x_{\text{val}}^j) \approx 0 \quad (13)$$

Furthermore, on the early stage when iteration is not enough we have

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial p_i^t} &\approx \frac{1}{N} \sum_{j=1}^N \frac{\partial \mathcal{L}(x_{\text{val}}^j)}{\partial \bar{x}} (W_i^{t-1} x_{\text{val}}^j)^T \\ &\approx \frac{1}{N} \sum_{j=1}^N \frac{\partial \mathcal{L}(x_{\text{val}}^j)}{\partial \bar{x}} (W_i^0 x_{\text{val}}^j)^T \end{aligned} \quad (14)$$

Learnable operations are initialized randomly thus they are actually making noise on input x and fitting x worse than non-learnable operations like skip-connect and pooling. If W_i fit x better, $\frac{\partial \mathcal{L}}{\partial p_i^t}$ would be smaller and p_i would be increased. As a result, in the early stage, non-learnable operations' p_i would be increased faster than learnable operations. Furthermore, for the biggest α_i (assume α_0), it's more likely that the gradient of α_0 smaller than others compared under softmax activation. And large learning rate would expand the gap. It means DARTS is more easily converge to α_0 . We have the following theorems.

Theorem 3.1. For function like $l = f(\sum_i \frac{\exp(\alpha_i)}{\sum_j \exp(\alpha_j)} z_i)$, let $p_i = \frac{\exp(\alpha_i)}{\sum_j \exp(\alpha_j)}$, $z = \sum_i p_i z_i$, $i^* = \arg \max_i \frac{\partial l}{\partial z} z_i$, we define margin $\delta = \min_{i \neq i^*} \frac{\partial l}{\partial z} (z_{i^*} - z_i) \geq 0$, if $\alpha_{i^*} \geq \alpha_i$, it holds that for any $\epsilon > 0$, function gradient achieves $p_{i^*} > 1 - \epsilon$ in iterations

$$t \leq \frac{n \ln((1 - \epsilon)n)}{\eta \delta} \quad (15)$$

Please see proof in Appendix. However, for single-level optimization, we have

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial p_i^t} &\approx \frac{1}{N} \sum_{j=1}^N \frac{\partial \mathcal{L}(x_{\text{val}}^j)}{\partial \bar{x}} (W_i^{t-1} x_{\text{val}}^j)^T \\ &\quad - \frac{\eta p_i^{t-1}}{N^2} \sum_{j=1}^N \left(\frac{\partial \mathcal{L}(x_{\text{train}}^j)}{\partial \bar{x}} \frac{\partial \mathcal{L}(x_{\text{train}}^j)}{\partial \bar{x}} \right) (x_{\text{train}}^j{}^T x_{\text{train}}^j) \\ &\leq \frac{1}{N} \sum_{j=1}^N \frac{\partial \mathcal{L}(x_{\text{val}}^j)}{\partial \bar{x}} (W_i^{t-1} x_{\text{val}}^j)^T \end{aligned} \quad (16)$$

It means for learnable operations, $\frac{\partial \mathcal{L}}{\partial p_i}$ could be decreased as training process goes. And after adequate iterations, learnable operations $\frac{\partial \mathcal{L}}{\partial p_i}$ would be smaller than non-learnable operations. Under gradient descent, p_i would be increased faster than non-learnable operations until exceed them.

3.3. Methodology

Above explanation shows that:

- Bi-level optimization make gradients of architecture parameters biased for that training data makes no direct effort on the gradients of architecture parameters. In fact, even doing bi-level optimization on the same dataset but not on the same data batch, the gradients could yet be biased.
- Using softmax as α 's activation function and inappropriate learning rate would expand the gap between non-learnable and learnable operations.

Therefore, we propose:

- Use single-level optimization meanwhile calculate gradients of w and α at the same time to instead bi-level optimization.
- Use sigmoid or other uncompetitive activation function to replace softmax.
- Decrease the learning rate if there are too many non-learnable operations.

4. Experiments

In this section, we introduce our experiments on different datasets and search spaces and mainly divided into two parts. One is to do experiments on NAS-Bench-201 to verify our hypothesis. The other part is to compare our improved DARTS (single-level optimization with sigmoid activation function) with different algorithms. For the first part, specifically, we firstly analyze the gradient correlation [13](#) and the gradients of p_i in bi-level optimization and single-level optimization. Then we compare the effect of using softmax as activation function with other independent activation function. At last we compare different learning rates' effect.

For the second part, we do experiments in NAS-Benchmark-201 and DARTS space, and on Cifar10 and ImageNet-1K dataset. Experiments show that we could steadily search high-performance architectures. On NAS-Benchmark-201 space, we stably find out nearly the optimal architecture. On DARTS space, we search the state-of-the-art architecture with 77.0% top1 accuracy (training setting follows PDARTS and without any additional module) on ImageNet-1K and steadily search architectures up to 76.5% top1 accuracy which is comparable with current reported best result.

4.1. Verification On NAS Benchmark 201

4.1.1 Search space.

NAS-Bench-201 [\[15\]](#) builds a cell-based search space, where a cell could be seen as a directed acyclic graph consisting of 4 nodes and 6 edges. Each network is stacked by 15 cells. Each edge represents an operation selected from (1) zeroize, (2) skip connection, (3) 1-by-1 convolution, (4) 3-by-3 convolution, and (5) 3-by-3 average pooling. The search space has 15,625 neural cell candidates in total. And all the candidates are given training accuracy/valid accuracy/test accuracy on three datasets: CIFAR-10, CIFAR-100 and ImageNet-16-120. In order to facilitate the analysis, in this section we set optimizers of network weights and architecture parameters as SGD.

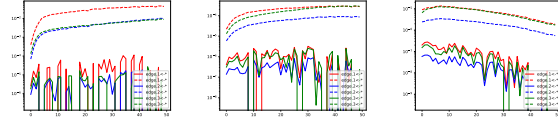


Figure 2: Comparison of gradient correlation [13](#) between bi-level and single-level optimization on Cifar10 dataset and in NAS-Benchmark-201 search space. It shows the gradient correlation on different nodes of first, middle and last cells. Solid lines represent bi-level and dashed lines represent single-level. Gradient correlation in bi-level optimization is nearly zero and far smaller than in single-level optimization.

4.1.2 Analysis on Gradient Correlation

In last section, we analyze the gradient correlation between bi-level optimization and single-level optimization. Fig [2](#) compares that the gradients correlation [13](#) in bi-level and single-level optimization in different cells. For bi-level optimization it's nearly zero and far less than that gradient-correlation in single-level. In fact, even from the same dataset but not the same batch, their gradients are still irrelevant. Table [1](#) compares the effect of whether backward propagations gradients of w and α in the same batch. The results show that even the gradients are from the same dataset but from different batches, the search progress collapses yet for the bias exists. The training phase lasts 50 epochs. w and α is optimized with SGD optimizer, learning rate is 0.005 (cosine scheduler). Original learning rate in [\[15\]](#) is 0.025, but it's too large and we show that it results badly in next section.

4.1.3 Analysis on Gradient Bias

We compare $\frac{\partial \mathcal{L}}{\partial p_i}$ in the last cell between bi-level and single-level optimization in Fig [3](#). For bi-level optimization, due to gradient bias, $\frac{\partial \mathcal{L}}{\partial p_i}$ of non-learnable operations are smaller than learnable operations and the gap is increased as the training process goes. Thus non-learnable operations' architecture parameter would be much larger than learnable operations'. For single-level optimization, $\frac{\partial \mathcal{L}}{\partial p_i}$ of learnable operations compete against non-learnable operations but finally get dominated. More comparison of different cells are shown in Appendix.

4.1.4 Analysis on Activation Function and Learning Rate

We compare the effects of learning rate and activation function on the searched results. Table [2](#) shows that larger learning rate would cause search phase collapsed, especially in

Table 1: Comparison the effect of whether gradients are calculated on the same data. Each setting is run 3 independent times to get average values. It shows that even on the same dataset but not the same mini batch, DARTS could not search good architectures for different batches are also low correlated. The experiments are under softmax as activation function. For learning rate damages search phase (it would be discussed bellow), the analysis is based on learning rate set as 0.005.

same dataset	same batch	CIFAR10	CIFAR100	ImageNet16-120
False	False	61.88±10.72	25.13±13.47	17.98±2.35
True	False	84.34±0.02	54.92±0.06	25.97±0.49
True	True	94.27±0.13	73.01±0.71	46.10±0.34

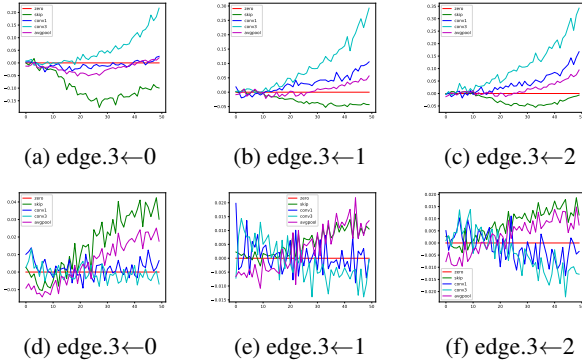


Figure 3: Comparison of gradient of p_i in the edges $3 \leftarrow *$ in the 15th cell between bi-level and single-level optimization. Figures (a)(b)(c) are on bi-level optimization and (d)(e)(f) are on single-level optimization.

the bi-level optimization framework. In the single-level optimization, the search phase performs more steadily. But too large learning rate also has bad effects. In addition, using sigmoid performs more steadily than softmax under different learning rate and optimization level. In this section, for the convenience of analysis, we use SGD as network weights' and architecture parameters' optimizer. In general NAS tasks, we advice to use SGD as network weights' optimizer and Adam as architecture parameters' optimizer. We could adjust network parameters' learning rate according to the searched non-parametric operations' ratio. If it's over a threshold such as 0.3, we should decrease it. And the process commonly doesn't cost too much time.

4.2. Comparison between algorithms

4.2.1 On NAS-201

At last, we compare our improvement with other algorithms. We use single-level to replace bi-level optimization, use sigmoid to replace softmax and α initialized as $-\ln(n-1)$. It shows single-level optimization could steadily get the SOTA and nearly the optimal architecture in NAS-201 space. On search phase, we train the super

Table 2: Comparison of effects of different activation function and learning rate. Each setting is run 3 independent times to get average values. Larger learning rate would cause search phase collapsed, especially in the bi-level optimization framework. In the single-level optimization, the search phase performs more steadily. And using sigmoid performs more steadily.

optim	activation	lr	CIFAR10	CIFAR100
bi	softmax	0.001	91.54±1.26	68.21±1.15
bi	softmax	0.005	61.88±10.72	25.13±13.47
bi	softmax	0.025	84.97±1.03	56.01±1.71
bi	sigmoid	0.001	92.53±1.62	69.62±2.05
bi	sigmoid	0.005	80.57±0.00	47.93±0.00
bi	sigmoid	0.025	74.27±14.12	41.86±18.56
single	softmax	0.001	94.36±0.00	73.51±0.00
single	softmax	0.005	94.27±0.13	73.01±0.71
single	softmax	0.025	86.58±0.00	58.33±0.00
single	sigmoid	0.001	94.36±0.00	73.51±0.00
single	sigmoid	0.005	94.36±0.00	73.51±0.00
single	sigmoid	0.025	93.10±0.00	69.24±0.00

model for 50 epochs, using SGD optimizer with a momentum of 0.9, cosine scheduler, batch size as 64, learning rate as 0.005, weight decay as 3×10^{-4} . For architecture parameters, we use Adam optimizer with a fixed learning rate of 3×10^{-4} , a momentum (0.5, 0.999). α is initialized as $-\ln(4)$ such that $\text{sigmoid}(\alpha) = 0.2$. And weight decay is set 0 to avoid extra gradients on zero-op. Results 3 show single-level optimization search steadily near the best architecture in NAS-201-Benchmark.

4.2.2 On DARTS

DARTS [15] is also a cell-based search space. Each cell contains 6 nodes and each node has to select 2 edges to connect with previous 2 nodes. Each edge has 8 operations: 3×3 and 5×5 separable convolution, 3×3 and 5×5 dilated

Table 3: Comparison of different NAS algorithms on NAS-Bench-201.

Method	CIFAR-10		CIFAR-100		ImageNet-16-120	
	validation	test	validation	test	validation	test
RSPS[23]	84.16±1.69	87.66±1.69	59.00±4.60	58.33±4.34	31.56±3.28	31.14±3.88
DARTS[28]	39.77±0.00	54.30±0.00	15.03±0.00	15.61±0.00	16.43±0.00	16.32±0.00
GDAS[14]	90.00±0.21	93.51±0.13	71.14±0.27	70.61±0.26	41.70±1.26	41.84±0.90
SETN [13]	82.25±5.17	86.19±4.63	56.86±7.59	56.87±7.77	32.54±3.63	31.90±4.07
ENAS [30]	39.77±0.00	54.30±0.00	15.03±0.00	15.61±0.00	16.43±0.00	16.32±0.00
CDARTS [37]	91.13±0.44	94.02±0.31	72.12±1.23	71.92 ±1.30	45.09±0.61	45.51±0.72
DARTS- [11]	91.03±0.44	93.80±0.40	71.36±1.51	71.53±1.51	44.87±1.46	45.12±0.82
single-level+sigmoid	91.55±0.00	94.36±0.00	73.49±0.00	73.51±0.00	44.87±0.00	46.34±0.00
optimal	91.61	94.37	73.49	73.51	46.77	47.31

separable convolution, 3×3 max-pooling, 3×3 average-pooling, skip-connect (identity), and zero (none). And the stacked networks have normal cells and reduction cells. The search space covers 10^{18} architectures which is quite large.

We directly do search process on the ImgeNet-1k dataset.

Table 4: Comparison with SOTA architectures on ImageNet in DARTS space. † denotes the average results of searched architectures (not the average results of retraining searched best architecture) and * denotes the best result. For single-level optimization and DARTS(ours), we search 3 times directly on the ImageNet-1K. Different from PDARTS training setting that re-train the searched architecture for 250 epochs, DropNAS train it for 600 epochs and DARTS+ 800 epochs.

Architecture	FLOPs	Params	Top-1.
	(M)	(M)	(%)
NASNet-A ([40])	564	5.3	74.0
AmoebaNet-C ([31])	570	6.4	75.7
PDARTS ([9])	557	4.9	75.6
PC-DARTS ([35])	597	5.3	75.8
DARTS ([28])	574	4.7	73.3
DARTS(ours)*	677	5.93	74.6
DARTS(ours)†	629	5.51	74.0
DARTS+([24])	591	5.1	76.3
CDARTS†([37])	732	6.1	76.3
CDARTS*([37])	704	6.3	76.6
DropNAS([19])	597	5.4	76.6
Single-level+softmax*	706	6.56	76.7
Single-level+softmax†	710	6.58	76.4
Single-level+sigmoid*	714	6.60	77.0
Single-level+sigmoid†	707	6.55	76.7

Specifically, input images are downsampled third times by

convolution layers at the beginning of the super model to reduce spatial resolution which follows [26]. On search phase, we train the super model for 50 epochs, using SGD optimizer with a momentum of 0.9, cosine scheduler, batch size as 360, learning rate as 0.025, weight decay as 3×10^{-4} . For architecture parameters, we use Adam optimizer with a fixed learning rate of 3×10^{-4} , a momentum (0.5, 0.999). For softmax as architecture parameters' activation function, α is initialized with $Gaussian(0, 10^{-3})$ and weight decay is 10^{-3} . For sigmoid, α is initialized as $-\ln(7)$ such that $sigmoid(\alpha) = 0.125$. And weight decay is set 0 to avoid extra gradients on zero-op. The search phase costs 4 GPUs for about 28 hours on NVIDIA GeForce RTX 2080ti. On retrain phase, we train the searched architecture following fully PDARTS setting and without any additional module. It lasts 250 epochs, using SGD optimizer with a momentum of 0.9, batch size as 768, an initial learning rate of 0.375 (decayed down to zero linearly), and a weight decay of 3×10^{-5} . Additional enhancements are adopted including label smoothing and an auxiliary loss tower during training as in PDARTS. Learning rate warm-up is applied for the first 5 epochs. The retrain phase costs 8 GPUs for about 3.5 days on NVIDIA GeForce RTX 2080ti. We search the state-of-the-art architecture with 77.0% top1 accuracy (training setting follows PDARTS and without any additional module) on ImageNet-1K and steadily search architectures up-to 76.5% top1 accuracy. More details are shown in Appendix.

5. conclusion

In this paper, we do analysis on the bi-level optimization in DARTS. Experiments show that bi-level optimization could cause gradients bias of architecture parameters and expand the gap between learnable and non-learnable operations. Furthermore, architecture parameters' activation function and large learning rate would aggregate the gap. We give a theoretical explanation and conduct experiments

to verify it. As a result, we propose to use single-level optimization as an instead. And we use uncompetitive activation function like sigmoid to replace softmax. Meanwhile, we do normalization for the initialization of sigmoid. Experiments on NAS-Benchmark-201 and DARTS space show that single-level optimization could steadily find out high-performance architectures.

References

- [1] Youhei Akimoto, Shinichi Shirakawa, Nozomu Yoshinari, Kento Uchida, Shota Saito, and Kouhei Nishida. Adaptive stochastic natural gradient method for one-shot neural architecture search. *arXiv preprint arXiv:1905.08537*, 2019. 2
- [2] Bowen Baker, Otkrist Gupta, Nikhil Naik, and Ramesh Raskar. Designing neural network architectures using reinforcement learning. *arXiv preprint arXiv:1611.02167*, 2016. 1, 2
- [3] Irwan Bello, Barret Zoph, Vijay Vasudevan, and Quoc V Le. Neural optimizer search with reinforcement learning. *arXiv preprint arXiv:1709.07417*, 2017. 2
- [4] 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, 3
- [5] Kaifeng Bi, Lingxi Xie, Xin Chen, Longhui Wei, and Qi Tian. Gold-nas: Gradual, one-level, differentiable. *arXiv preprint arXiv:2007.03331*, 2020. 3
- [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] Han Cai, Chuang Gan, and Song Han. Once for all: Train one network and specialize it for efficient deployment. *arXiv preprint arXiv:1908.09791*, 2019. 2
- [8] Han Cai, Ligeng Zhu, and Song Han. Proxylessnas: Direct neural architecture search on target task and hardware. *arXiv preprint arXiv:1812.00332*, 2018. 1, 2
- [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, 7
- [10] Yukang Chen, Tong Yang, Xiangyu Zhang, Gaofeng Meng, Chunhong Pan, and Jian Sun. Detnas: Neural architecture search on object detection. *arXiv preprint arXiv:1903.10979*, 2019. 2
- [11] Xiangxiang Chu, Xiaoxing Wang, Bo Zhang, Shun Lu, Xiaolin Wei, and Junchi Yan. Darts-: robustly stepping out of performance collapse without indicators. *arXiv preprint arXiv:2009.01027*, 2020. 7
- [12] Xiangxiang Chu, Tianbao Zhou, Bo Zhang, and Jixiang Li. Fair darts: Eliminating unfair advantages in differentiable architecture search. *arXiv preprint arXiv:1911.12126*, 2019. 1, 2, 3
- [13] Xuanyi Dong and Yi Yang. One-shot neural architecture search via self-evaluated template network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3681–3690, 2019. 7
- [14] Xuanyi Dong and Yi Yang. Searching for a robust neural architecture in four gpu hours. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1761–1770, 2019. 7
- [15] Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. *arXiv preprint arXiv:2001.00326*, 2020. 1, 2, 3, 5, 6
- [16] Yonggan Fu, Wuyang Chen, Haotao Wang, Haoran Li, Yingyan Lin, and Zhangyang Wang. Autogan-distiller: Searching to compress generative adversarial networks. *arXiv preprint arXiv:2006.08198*, 2020. 2
- [17] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Nas-fpn: Learning scalable feature pyramid architecture for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7036–7045, 2019. 2
- [18] Zichao Guo, Xiangyu Zhang, Haoyuan Mu, Wen Heng, Zechun Liu, Yichen Wei, and Jian Sun. Single path one-shot neural architecture search with uniform sampling. *arXiv preprint arXiv:1904.00420*, 2019. 2
- [19] Weijun Hong, Guilin Li, Weinan Zhang, Ruiming Tang, Yunhe Wang, Zhenguo Li, and Yong Yu. Dropnas: Grouped operation dropout for differentiable architecture search. In *International Joint Conference on Artificial Intelligence*, 2020. 2, 3, 7
- [20] Andrew Hundt, Varun Jain, and Gregory D Hager. sharp-darts: Faster and more accurate differentiable architecture search. *arXiv preprint arXiv:1903.09900*, 2019. 2
- [21] 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. 2
- [22] Guilin Li, Xing Zhang, Zitong Wang, Zhenguo Li, and Tong Zhang. Stacnas: Towards stable and consistent optimization for differentiable neural architecture search. *arXiv preprint arXiv:1909.11926*, 2019. 1, 3
- [23] Liam Li and Ameet Talwalkar. Random search and reproducibility for neural architecture search. In *Uncertainty in Artificial Intelligence*, pages 367–377. PMLR, 2020. 7
- [24] Hanwen Liang, Shifeng Zhang, Jiacheng Sun, Xingqiu He, Weiran Huang, Kechen Zhuang, and Zhenguo Li. Darts+: Improved differentiable architecture search with early stopping. *arXiv preprint arXiv:1909.06035*, 2019. 2, 7
- [25] Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, Wei Hua, Alan L Yuille, and Li Fei-Fei. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 82–92, 2019. 2
- [26] Chenxi Liu, Piotr Dollár, Kaiming He, Ross Girshick, Alan Yuille, and Saining Xie. Are labels necessary for neural architecture search? *arXiv preprint arXiv:2003.12056*, 2020. 7
- [27] 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. 1, 2, 3

- [28] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018. 1, 2, 3, 7
- [29] Marcelo Gennari do Nascimento, Theo W Costain, and Victor Adrian Prisacariu. Finding non-uniform quantization schemes using multi-task gaussian processes. *arXiv preprint arXiv:2007.07743*, 2020. 2
- [30] 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, 7
- [31] 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, 7
- [32] Esteban Real, Sherry Moore, Andrew Selle, Saurabh Saxena, Yutaka Leon Suematsu, Jie Tan, Quoc V Le, and Alexey Kurakin. Large-scale evolution of image classifiers. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2902–2911. JMLR. org, 2017. 2
- [33] Richard Shin, Charles Packer, and Dawn Song. Differentiable neural network architecture search. 2018. 2
- [34] Dimitrios Stamoulis, Ruizhou Ding, Di Wang, Dimitrios Lymberopoulos, Bodhi Priyantha, Jie Liu, and Diana Marculescu. Single-path nas: Device-aware efficient convnet design. *arXiv preprint arXiv:1905.04159*, 2019. 2
- [35] 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. 1, 2, 3, 7
- [36] Zhimin Xu, Si Zuo, Edmund Y Lam, Byoungcho Lee, and Ni Chen. Autosegnet: An automated neural network for image segmentation. *IEEE Access*, 8:92452–92461, 2020. 2
- [37] Hongyuan Yu and Houwen Peng. Cyclic differentiable architecture search. *arXiv preprint arXiv:2006.10724*, 2020. 2, 7
- [38] Arber Zela, Thomas Elsken, Tonmoy Saikia, Yassine Marrakchi, Thomas Brox, and Frank Hutter. Understanding and robustifying differentiable architecture search. *arXiv preprint arXiv:1909.09656*, 2019. 1
- [39] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016. 2
- [40] 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, 7

A. Proof of theorem 3.1

Proof. It’s easily to see that $p_i \geq p_j$ if $\alpha_i \geq \alpha_j$ and $p_{i^*} \geq \frac{1}{n}$ for $\sum_i p_i = 1$. Consider

$$\begin{aligned}
 \frac{\partial l}{\partial \alpha_i} &= \frac{\partial l}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \alpha_i} \\
 &= \frac{\partial l}{\partial \mathbf{z}} (p_i(1 - p_i)\mathbf{z}_i - p_i \sum_{j \neq i} p_j \mathbf{z}_j) \\
 &= \frac{\partial l}{\partial \mathbf{z}} \sum_j p_i p_j (\mathbf{z}_i - \mathbf{z}_j) \\
 &= p_i \sum_j p_j \frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_i - \mathbf{z}_j)
 \end{aligned} \tag{17}$$

For i^* and $\forall i \neq i^*$, because $\frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_{i^*} - \mathbf{z}_j) \geq \frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_i - \mathbf{z}_j)$ and $p_{i^*} \geq p_i \geq 0$, thus $\frac{\partial l}{\partial \alpha_{i^*}} \geq \frac{\partial l}{\partial \alpha_i}$.

As a result

$$\begin{aligned}
 \frac{\partial l}{\partial \alpha_{i^*}} - \frac{\partial l}{\partial \alpha_i} &= p_{i^*} \sum_j p_j \frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_{i^*} - \mathbf{z}_j) - p_i \sum_j p_j \frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_i - \mathbf{z}_j) \\
 &= (p_{i^*} - p_i) \sum_j p_j \frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_{i^*} - \mathbf{z}_j) + p_i \left(\sum_j p_j \frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_{i^*} - \mathbf{z}_j) \right. \\
 &\quad \left. - \sum_j p_j \frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_i - \mathbf{z}_j) \right) \\
 &= (p_{i^*} - p_i) \sum_j p_j \frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_{i^*} - \mathbf{z}_j) + p_i \sum_j p_j \frac{\partial l}{\partial \mathbf{z}} (\mathbf{z}_{i^*} - \mathbf{z}_i) \\
 &\geq (p_{i^*} - p_i) \delta + p_i \delta \\
 &= p_{i^*} \delta \\
 &\geq \frac{\delta}{n}
 \end{aligned} \tag{18}$$

Under gradient ascent, on update step t we have, for any $i \neq i^*$,

$$\begin{aligned}
 \alpha_{i^*}^t - \alpha_i^t &= \alpha_{i^*}^{t-1} - \alpha_i^{t-1} + \eta \left(\frac{dl_{t-1}}{d\alpha_{i^*}^{t-1}} - \frac{dl_{t-1}}{d\alpha_i^{t-1}} \right) \\
 &\geq \alpha_{i^*}^{t-1} - \alpha_i^{t-1} + \eta p_{i^*}^t \delta_t \\
 &\geq \alpha_{i^*}^0 - \alpha_i^0 + \eta \sum_t p_{i^*}^t \delta_t \\
 &\geq \frac{\eta t \delta}{n}
 \end{aligned} \tag{19}$$

When $t = \frac{n \ln((1-\epsilon)n)}{\eta \delta}$, we have

$$\alpha_{i^*} - \alpha_i \geq \ln((1 - \epsilon)n) \tag{20}$$

Thus

$$p_{i^*} = \frac{1}{\sum_i \exp(\alpha_i - \alpha_{i^*})} \geq \frac{1}{\sum_i \exp(-\ln((1-\epsilon)n))} = 1-\epsilon \quad (21)$$

Under gradient descent, for $i^* = \arg \min_i \frac{\partial l}{\partial z} z_i$, we have the same conclusion. \square

B. More comparison of gradient of p_i in section 4.2.3

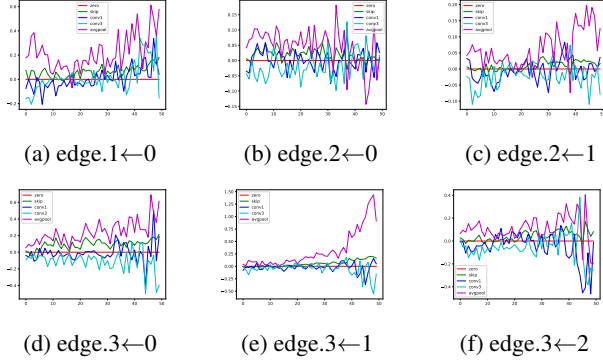


Figure 4: Bi-level optimization. On the 0th cell.

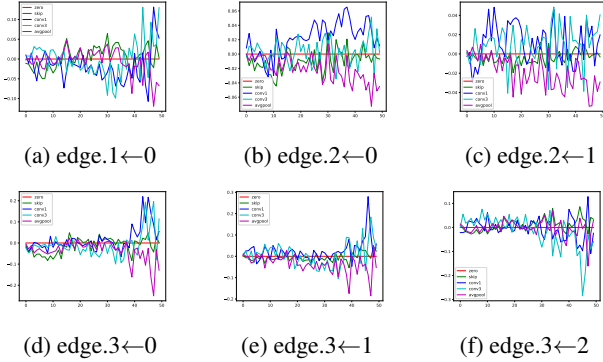


Figure 5: Bi-level optimization. On the 8th cell.

C. Searched results in DARTS

We use single-level optimization to directly search on the ImageNet-1K dataset in DARTS space. Initialize α_i as $-\ln(7)$ so that $\text{softmax}(\alpha_i) = 0.125$ would have more better results. It's first time that searched architecture in DARTS space get 77% top-1 accuracy on ImageNet-1K dataset. The training setting follows PDARTS and without any additional tricks. In addition, for single-level optimization, using half of dataset could search promising results as well.

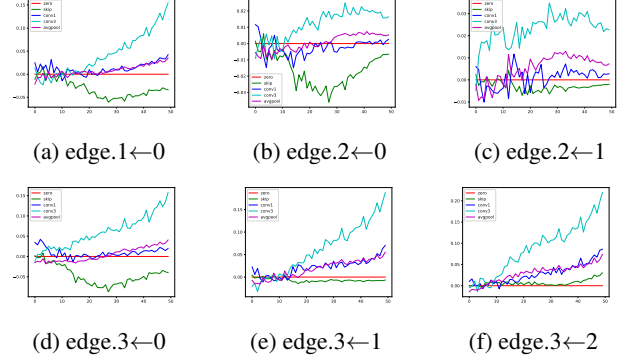


Figure 6: Bi-level optimization. On the 16th cell.

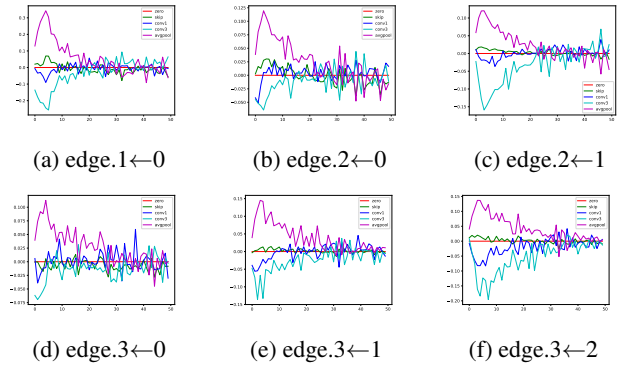


Figure 7: Single-level optimization. On the 0th cell.

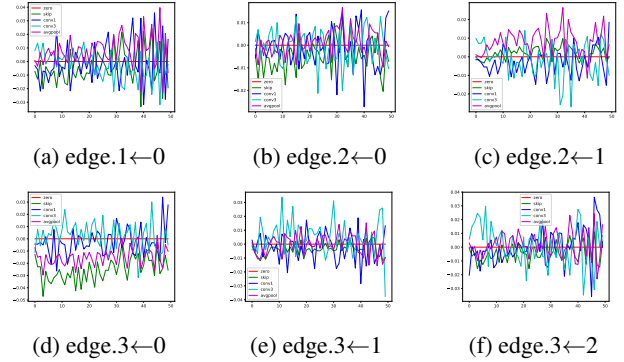


Figure 8: Single-level optimization. On the 8th cell.

D. Visualization of architectures

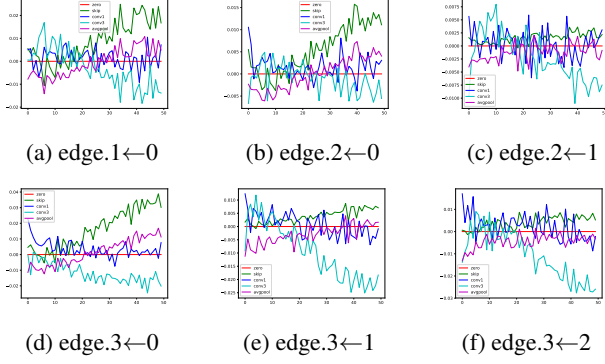


Figure 9: Single-level optimization. On the 16th cell.

Table 5: * denotes α is initialized as $-\ln(7)$. 'Data' means use full or half of data to do search.

Activation	Data (M)	Seed (M)	FLOPs (%)	Params	Top-1.
softmax	full	0	714.72	6.58	76.28
softmax	full	1	712.92	6.56	76.71
softmax	full	2	722.25	6.61	76.27
sigmoid	full	0	738.21	6.69	76.12
sigmoid	full	1	738.21	6.69	76.58
sigmoid	full	2	721.35	6.60	76.29
sigmoid*	full	0	707.89	6.50	76.26
sigmoid*	full	1	721.35	6.60	77.0
sigmoid*	full	2	700.61	6.40	76.95
softmax	half	0	712.01	6.55	76.51
softmax	half	1	692.18	6.36	76.31
softmax	half	2	709.04	6.45	76.51
sigmoid*	half	0	720.44	6.59	76.67
sigmoid*	half	1	692.18	6.36	76.78
sigmoid*	half	2	707.89	6.50	76.54

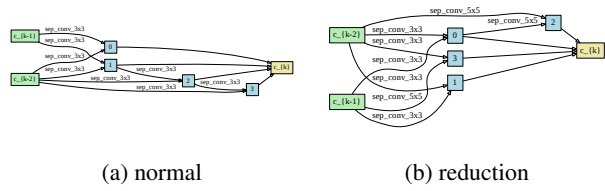


Figure 10: activation=softmax, data=full, seed=0

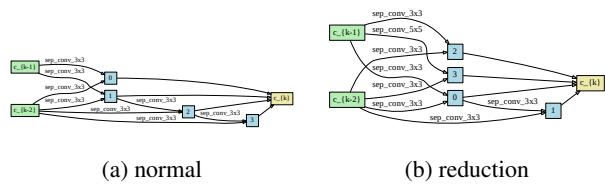


Figure 11: activation=softmax, data=full, seed=1

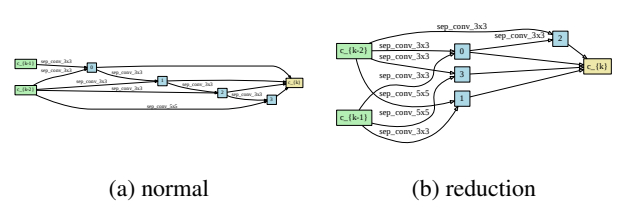


Figure 12: activation=softmax, data=full, seed=2

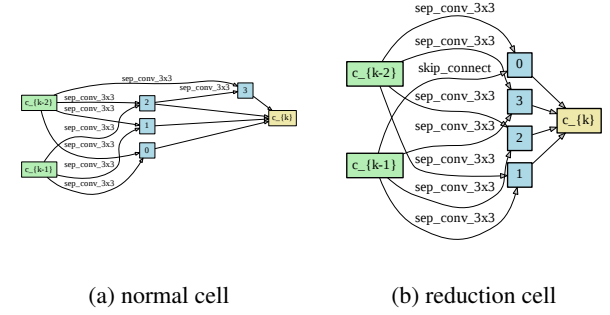


Figure 13: activation=sigmoid*, data=full, seed=0

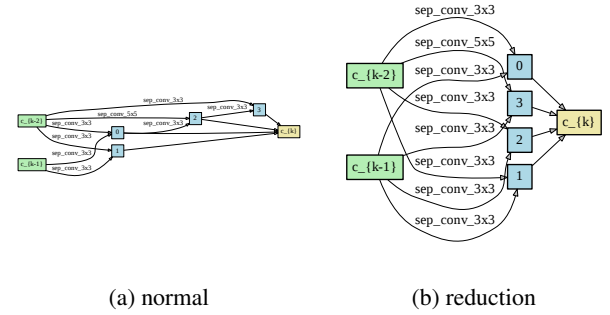


Figure 14: activation=sigmoid*, data=full, seed=1

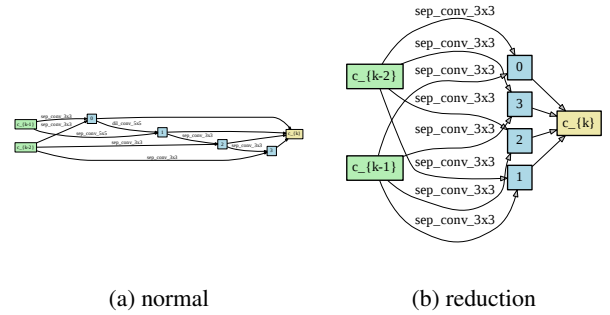


Figure 15: activation=sigmoid*, data=full, seed=2