Supplementary of "SCARLET-NAS: Bridging the Gap between Stability and Scalability in Weight-sharing Neural Architecture Search"

Xiangxiang Chu^{*} Bo Zhang Qingyuan Li Ruijun Xu Xudong Li

{chuxiangxiang,zhangboll,liqingyuan, xuruijun}@xiaomi.com, lixudong16@mails.ucas.edu.cn

1. Proof

Proof. First, we prove that Equation 3 (main text) holds for $\forall o \in \{0, 1, ..., n - 2\}$. In this case, it's sufficient to prove the output of the first convolution $Conv_{(c_l,m,k,k)}$ can be exactly matched by adding $Conv_{(c_l,c_{l+1},1,1)}$ before $Conv_{(c_{l+1},m,k,k)}$. Let $W_{c_l,c_{l+1},1,1}^1$ and $W_{c_l,m,k,k}^2$ be the weight tensors of $Conv_{(c_l,c_{l+1},1,1)}$ and $Conv_{(c_{l+1},m,k,1)}$ respectively. Let $W_{c_l,m,k,k}^3$ be the weight tensors of $Conv_{(c_l,m,k,1)}$. Let w be one element of the tensor. We have

$$y = Conv_{(c_l, c_{l+1}, 1, 1)}(x_l^{c_l}), z = Conv_{(c_{l+1}, m, k, 1)}(y)$$
(1)

$$y(i,j,c) = \sum_{p=1}^{c_l} w_{p,c,1,1}^1 x(i,j,p)$$
(2)

Also,

$$z(i,j,c) = \sum_{q=1}^{k} \sum_{p=1}^{c_{l+1}} w_{p,c,q,q}^{2} y(i+q,j+q,p)$$

$$= \sum_{q=1}^{k} \sum_{p=1}^{c_{l+1}} w_{p,c,q,q}^{2} (\sum_{u=1}^{c_{l}} w_{u,p,1,1}^{1} x(i+q,j+q,u))$$

$$= \sum_{q=1}^{k} \sum_{p=1}^{c_{l+1}} \sum_{u=1}^{c_{l}} w_{p,c,q,q}^{2} w_{u,p,1,1}^{1} x(i+q,j+q,u)$$

$$= \sum_{q=1}^{k} \sum_{u=1}^{c_{l}} w_{u,c,q,q}^{3} x(i+q,j+q,u)$$

(3)

Therefore, the first part is proved by setting

$$w_{u,c,q,q}^{3} = \sum_{p=1}^{c_{l+1}} w_{p,c,q,q}^{2} w_{u,p,1,1}^{1}.$$
 (4)

For o = n - 1, we replace a skip connection with an *ELS*. We can iteratively apply the first part of the proof till the end of searchable layers.

Algorithm 1 The constrained and weighted NAS pipeline.

Input: Supernet S, the number of generations N, population size n, validation dataset D, constraints C, objective weights w **Output:** A set of K individuals on the Pareto front. Train supernet S defined on the scalable search space. Uniformly generate the populations P_0 and Q_0 until each has n individuals satisfying C_{FLOPS} , C_{Accuracy} . for i = 0 to N - 1 do $R_i = P_i \cup Q_i$ $F = \text{non-dominated-sorting}(R_i)$ Pick *n* individuals to form P_{i+1} by ranks and the crowding distance weighted by w. $Q_{i+1} = \emptyset$ while $size(Q_{i+1}) < n$ do $M = \text{tournament-selection}(P_{i+1})$ $q_{i+1} = \operatorname{crossover}(M) \cup \operatorname{hierarchical-mutation}(M)$ {Check the FLOPS constraint at first (It takes < 1ms). if $FLOPS(q_{i+1}) > FLOPS_{max}$ then continue end if Evaluate model q_{i+1} with S on D {Check the accuracy constraint (It takes $\approx 60s$). if $Accuracy(q_{i+1}) > Acc_{min}$ then Add q_{i+1} to Q_{i+1} end if end while end for Select K equispaced models near Pareto-front from P_N

2. Algorithm

Our constrained and weighted NAS pipeline is listed in Algorithm 1 and Fig. 1.

^{*}This work was done when all the authors were at Xiaomi AI Lab.



Figure 1. Constrained and weighted NSGA-II Pipeline. It starts with a uniform initialization (top left) with some constraints (red) to generate the initial population. The trained scalable supernet serves as a fast evaluator to decide the relative performance of each model so that they can be grouped into several Fronts $(F_1, F_2, ...)$ by weighted non-dominated sorting (right). Only the top n of them make up the next generation P_{i+1} , based on which Q_{i+1} is produced with tournament selection, crossover and mutation (blue) under the same constraints (bottom left). The whole evolution loops until we reach Pareto-optimality.

3. Experiments

3.1. Search Space

For later experiments, we add skip connections to commonly used search space to construct S_1 and S_2 . They are described as follows,

Search Space S_1 . It is similar to ProxylessNAS [2], where MobileNetV2 [7] is adopted as its backbone. In particular, S_1 is represented as a block-level supernet with L = 19 layers of N = 7 choices each. Its total size is 7^{19} . The choices are,

- MobileNetV2's inverted bottleneck blocks [7] of two expansion rates (x) in (3,6), three kernel sizes (y) in (3,5,7), labelled as MBExKy¹,
- skip connection (the 6th choice²).

Search Space S_2 . On top of S_1 , we give each inverted bottleneck a squeeze-and-excitation [5] option (e.g., ExKy, $ExKy_SE$), similar to MnasNet [8]. Its total size thus becomes 13^{19} .

We have to notice that skip connections are commonly used [8, 6, 1], but meticulously neglected in recent singlepath one-shot methods [4, 3].

3.2. NSGA-II Hyperparameters

The hyperparameters for the weighted NSGA-II approach are given in Table 1.

Item	value	Item	value
Population N	70	Mutation Ratio	0.8
p_{rm}	0.2	p_{re}	0.65
p_{pr}	0.15	p_M	0.7
p_{K-M}	0.3		

Table 1. Hyperparameters for the weighted NSGA-II approach.

3.3. More Details about Scalable Supernet with ELS

Given an input of a chickadee³ image from ImageNet, we illustrate both high-level and low-level feature maps of the trained supernet with our proposed improvements in Figure 2. Pure skip connection easily interferes with the training process as it causes incongruence with other choice blocks. Note the channel size of feature map after Choice 6 in Figure 2 (a) is half of others because the previous channel size is 16, while other choice blocks output 32 channels. This effect is largely attenuated by ELS. As it goes deeper, we still observe consistent high-level features. Specifically, when ELS is not enforced, high-level features of deeper channels easily get blurred out, while the supernet with ELS enabled continues to learn useful features in deeper channels.

3.4. Search Space Evaluation

NAS results can benefit from good search space. To prove the validity of the proposed method, we show our search space has a wide range and is not particularly designed. We pick two extreme cases, one with all identity blocks (only the stem and the tail remains), the other with all K7E6s. They have the minimum and the maximum FLOPS respectively. We list their evaluation result in Table 2. The former has 24.1% top-1 accuracy on ImageNet, and the latter 76.8% at a cost of 557M FLOPs. Both are infeasible solutions as they violate either acc_{min} or $madds_{max}$. It's thus a challenging task to deal with such search space for ordinary search techniques.

3.5. Analysis of SCARLET Models

SCARLET-A makes full use of large kernels (five 5×5 and seven 7×7 kernels) to enlarge receptive field. Besides it activates many squeezing and excitation (12 out of 19) blocks to improve its classification performance. At the early stage, it appreciates either large kernels and small expansion ratios or small kernels and large expansion ratios to balance the trade-off between accuracy and FLOPs.

SCARLET-B chooses two identity operations. Compared with A, it shortens network depth at the last stages. Besides, it utilizes squeezing and excitation block extensively (14 out of 17). It places a large expansion block with large kernels at the tail stage.

¹The order of numbering o = (x - 3) + (y - 3)/2.

²zero-based numbering

³ImageNet ID: n01592084_7680

	20
Choice 1 Cho	-
Choice 2 Charles Charl	4
Choice 3 2 3 2 5 2 5 2 5 2 5 2 5 2 5 2 5 2 5 2	10 m
Choice 4	4
Choice 5 SALLARY SALARY SA	-
Choice 6 A A A A A A A A A A A A A A A A A A	

(a) First choice blocks' feature maps without ELS

Choice 0	急急	2 2 2 2 C	E A A	St. 28 5	1 A B	1. S. S. D.	A A	<u>A</u> <u>A</u> <u>A</u>	222	808228
Choice 1	S.A.	12 A.A.	R & B	2. 2	and i		S B AN	R S.A	2. 2. 2	22222 B
Choice 2	R.a.	1 A. 2.	R & B	32 2 2	(A. A.		37 28-05-0	A. S. A	2. 2. 2	\$ 2 1 2 A B
Choice 3	35.28	220	a a s	22 2 8	1. E.B.		No. Al M	A Sta	2. 2. 2	211228
Choice 4	S. A.	229	2 2 2	S. 28 8	1. 2. 2. 3	22	8 2 De D	A S.S.	2 2 2	22222A
Choice 5	30.28.	22	2 2 2	S 2 2	C. C. A.	2 1 1 B	N De De C	A B A	2. 2. 2	22222
Choice 6	Sec. Mar	m 2 2	R. A. S	and the	A. A. A.	A A A	2 2 20	282	8. 2. 2	22232

(b) First choice blocks' feature maps with ELS

Choice 0	an 📴 da da cara en estas 📴 e		1 2 12 3	2 2 2 2	1223.22
Choice 1		S 2 5 9	125.26	A STAND	10000
Choice 2		19 19 19 19 19	1999	1 2 2 2 1	12002
Choice 3	NA PARANCE IN MARKED	18 Sec. 2. 2.			1921-12
Choice 4		A PERSON	121.20	A BEL	See 22
Choice 5			121	NO.	A Call
Choice 6	这些现在的现在,我们就不是是你的问题。	12 B 12 8 8			2222

(c) High-level choice blocks' feature maps without ELS

Choice 0	PE15213	201920	17.194	SCALE.	5 N. 19 10		
Choice 1	A 124 121	S 5 154 ()	(1944))	<u> 1988</u>	S. S. S. S. S. S.	1-19- 2.	
Choice 2	(入生)的公共的	9 N 19 N	5 M 24	20121	52.525	1-1/2 P	고 아파 이 아파 이
Choice 3	·深望地望远。			20424	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	1 30 m 1 1	
Choice 4	A 14 19 19	S (18)	的加强。	A 44	22.20	1-12 2 2	
Choice 5	A 19 19 19 19 19			S 1 31	52.525	1-4/2 2 2	
Choice 6	25 19 12 B	25 1 2 3	17519 A	10121	S. N. S. S.	9 5-4 P 2 3	21 21 21 21 21 21 21 22

(d) High-level choice blocks' feature maps with ELS

Figure 2. Learned low-level and high-level features for the supernet with and without ELS.

Models	FLOPS (M)	$> madds_{max}$	Top-1 (%)	Top-5 (%)	$< acc_{min}$
All Identity	23	No	24.1	45.0	Yes
All K7E6	557	Yes	76.8	93.3	No

Table 2. Full train results of models with minimal and maximal FLOPS.

SCARLET-C uses three identity operations and utilizes small expansion ratio extensively to cut down the FLOPs, large expansion ratio at the tail stage whose resolution is 7×7 . It prefers large kernels before the downsampling layers. Besides, it makes an extensive use of squeeze and excitation to boost accuracy.

References

- Gabriel Bender, Pieter-Jan Kindermans, Barret Zoph, Vijay Vasudevan, and Quoc Le. Understanding and Simplifying One-Shot Architecture Search. In *ICML*, pages 549–558, 2018.
- [2] Han Cai, Ligeng Zhu, and Song Han. ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware. In *ICLR*, 2019.

- [3] Xiangxiang Chu, Bo Zhang, Ruijun Xu, and Jixiang Li. Fair-NAS: Rethinking Evaluation Fairness of Weight Sharing Neural Architecture Search. arXiv preprint. arXiv:1907.01845, 2019.
- [4] 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.
- [5] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-Excitation Networks. In CVPR, pages 7132–7141, 2018.
- [6] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: Differentiable Architecture Search. In *ICLR*, 2019.
- [7] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobileNetV2: Inverted Residuals and Linear Bottlenecks. In CVPR, pages 4510– 4520, 2018.
- [8] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, and Quoc V Le. Mnasnet: Platform-Aware Neural Architecture Search for Mobile. In *CVPR*, 2019.