



# DSNAS: Direct Neural Architecture Search without Parameter Retraining

报告人：时间：2020年7月9日

# Outline

---

- Motivation
  - Search space
  - Method
  - Experiments
-

# motivation

---

## Problem:

Most of the NAS methods proposed recently are implicitly two-stage methods. These two stages are searching and evaluation.

The correlation between the performance at the end of searching and after training in evaluation is fairly low.

## Contribution:

propose a plug-and-play NAS framework, DSNAS, as an efficient solution to this problem in large scale.

---

# Outline

---

- Motivation
  - **Search space**
  - Method
  - Experiments
-

# Search Space

Input	Block	Channels	Repeat	Stride
$224^2 \times 3$	$3 \times 3$ Conv	16	1	2
$112^2 \times 16$	CB	64	4	2
$56^2 \times 64$	CB	160	4	2
$28^2 \times 160$	CB	320	8	2
$14^2 \times 320$	CB	640	4	2
$7^2 \times 640$	$1 \times 1$ Conv	1024	1	1
$7^2 \times 1024$	GAP	-	1	-
1024	FC	1000	1	-

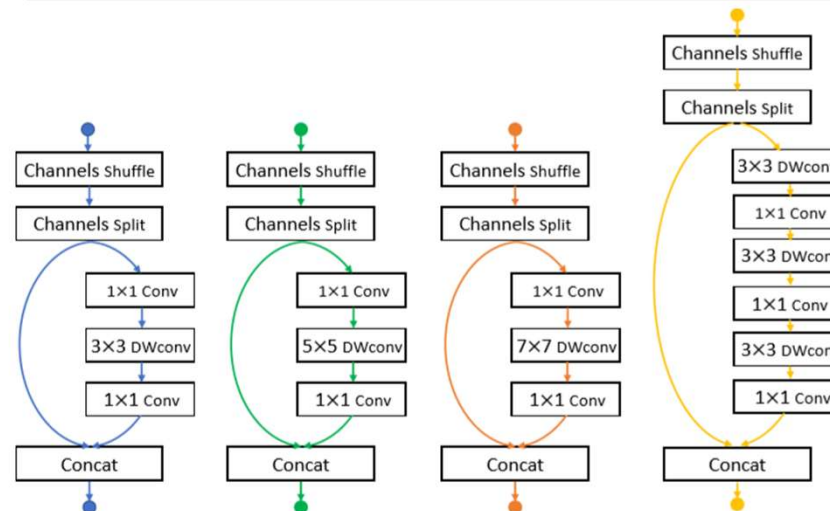


Figure 4. Choice blocks with stride=1. Choice blocks in search space. From left to right: Choice\_3, Choice\_5, Choice\_7, Choice\_x.

# Outline

---

- Motivation
  - Search space
  - **Method**
  - Experiments
-

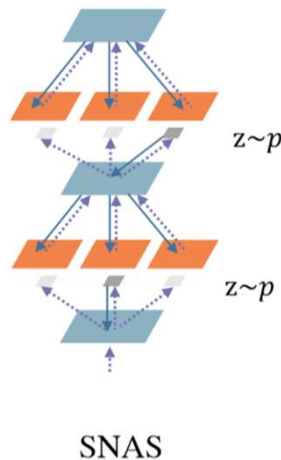
# Method

---

$$\tilde{O}_{i,j}(\cdot) = \mathbf{Z}_{i,j}^T \mathbf{O}_{i,j}(\cdot). \quad (1) \quad \text{Forward, } \mathbf{Z} \text{ is random variables}$$

$$\mathbb{E}_{\mathbf{Z} \sim p_{\alpha}(\mathbf{Z})} [L_{\theta}(\mathbf{Z})], \quad (2) \quad \text{Monte Carlo estimate of the expectation of task objective } L_{\theta}(\mathbf{Z})$$

$$\frac{\partial \mathcal{L}}{\partial \alpha_{i,j}^k} = \frac{\partial \mathcal{L}}{\partial x_j} \mathbf{O}_{i,j}^T(x_i) (\delta(k' - k) - \tilde{\mathbf{Z}}_{i,j}) \mathbf{Z}_{i,j}^k \frac{1}{\lambda \alpha_{i,j}^k}, \quad (3) \quad \text{Gradient, } \lambda \text{ is temperature}$$



# Method

---

$$\begin{aligned} & \mathbb{E}_{\tilde{\mathbf{Z}} \sim p(\tilde{\mathbf{Z}})} \left[ \frac{\partial \mathcal{L}}{\partial \alpha_{i,j}^k} \right] \\ &= \mathbb{E}_{\tilde{\mathbf{Z}} \sim p(\tilde{\mathbf{Z}})} \left[ \nabla_{\alpha_{i,j}^k} \log p(\tilde{\mathbf{Z}}) \left[ \frac{\partial \mathcal{L}}{\partial x_j} \mathbf{O}_{i,j}^T(x_i) \tilde{\mathbf{Z}}_{i,j} \right]_c \right], \end{aligned} \quad (4) \quad \text{Policy gradient}$$

$$\begin{aligned} & \lim_{\lambda \rightarrow 0} \mathbb{E}_{\tilde{\mathbf{Z}} \sim p(\tilde{\mathbf{Z}})} \left[ \frac{\partial \mathcal{L}}{\partial \alpha_{i,j}^k} \right] \\ &= \lim_{\lambda \rightarrow 0} \mathbb{E}_{\tilde{\mathbf{Z}} \sim p(\tilde{\mathbf{Z}})} \left[ \nabla_{\alpha_{i,j}^k} \log p(\tilde{\mathbf{Z}}_{i,j}) \left[ \frac{\partial \mathcal{L}}{\partial x_j} \tilde{\mathbf{O}}_{i,j}(x_i) \right]_c \right] \\ &= \mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z})} \left[ \nabla_{\alpha_{i,j}^k} \log p(\mathbf{Z}_{i,j}) \left[ \frac{\partial L}{\partial x_j} \mathbf{O}_{i,j}^T(x_i) \mathbf{Z}_{i,j} \right]_c \right] \\ &= \mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z})} \left[ \nabla_{\alpha_{i,j}^k} \log p(\mathbf{Z}_{i,j}) \left[ \frac{\partial L}{\partial x_j} \sum_k \mathbf{O}_{i,j}^k(x_i) \mathbf{Z}_{i,j}^k \right]_c \right], \end{aligned} \quad (5)$$

---



# Method

---

$$\mathbb{E}_{\mathbf{Z} \sim p(\mathbf{Z})} [\nabla_{\alpha_{i,j}^k} \log p(\mathbf{Z}_{i,j}) [\frac{\partial L}{\partial x_j} \sum_k O_{i,j}^k(x_i) Z_{i,j}^k]_c],$$

(5)

$$\begin{aligned} C(\mathbf{Z}_{i,j}) &= \sum_k \frac{\partial L}{\partial x_j} O_{i,j}^k(x_i) Z_{i,j}^k \\ &= \frac{\partial L}{\partial x_j^i} O_{i,j}^s(x_i) Z_{i,j}^s = \frac{\partial L}{\partial x_j^i} x_j^i \\ &= \frac{\partial L}{\partial x_j^i} \frac{\partial x_j^i}{\partial Z_{i,j}^s} = \frac{\partial L}{\partial Z_{i,j}^s}, \end{aligned}$$

(6)  $x_j^i = O_{i,j}^s(x_i) Z_{i,j}^s$  is the output of the operation  $O_{i,j}^s$  Chosen at edge(i,j)

---

# Method

---

## Algorithm 1 Discrete SNAS

---

**Require:** parent network, operation parameters  $\theta$  and categorical arch distribution  $p_{\alpha}(\mathbf{Z})$

Initialize  $\theta, \alpha$

**while** not converged **do**

    Sample one-hot random variables  $\mathbf{Z}$  from  $p_{\alpha}(\mathbf{Z})$

    Construct child network with  $\theta$  according to  $\mathbf{Z}$ , multiply a  $1^{dummy}$  after each feature map  $X$

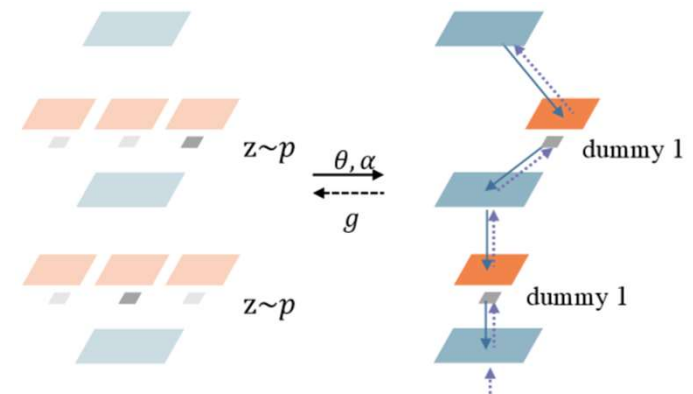
    Get a batch from data and forward to get  $L$

    Backward  $L$  to both  $\theta$  and  $1^{dummy}$ , backward  $\log p_{\alpha}(\mathbf{Z})$  to  $\alpha$

    Update  $\theta$  with  $\frac{\partial L}{\partial \theta}$ , update  $\alpha$  with  $\frac{\partial \log p_{\alpha}(\mathbf{Z})}{\partial \alpha} \frac{\partial L}{\partial 1^{dummy}}$

**end while**

---



DSNAS

# Method

---

Method	Forward time	Backward time	Memory
Subnetwork	$O(P)$	$O(Q)$	$O(M)$
SNAS	$O(nP)$	$O(nQ)$	$O(nM)$
ProxylessNAS*	$O(nP)$	$O(nQ)$	$O(nM)$
ProxylessNAS	$O(2P)$	$O(2Q)$	$O(2M)$
DSNAS	$O(P)$	$O(Q)$	$O(M)$

Table 1. Computation complexity comparison between SNAS, ProxylessNAS and DSNAS. ProxylessNAS\* indicates its theoretical complexity.

---

# Method

---

## **Progressive early stop:**

The non-zero entropy keeps the sampling going on until the end, regardless of the fact that sampling at that uncertainty level does not bring significant gains.

$$\min\{\alpha_{i,j}^k - \alpha_{i,j}^m, \forall m \neq k \mid \alpha_{i,j}^k = \max\{\alpha_{i,j}\}\} \geq h. \quad (8)$$

## **Comparison with one-shot NAS:**

one-shot NAS only do architecture optimization in one-shot, all of them are two-stage methods.

DSNAS optimizes architecture alongside with parameters, which is expected to save some resources in the pretraining stage.

superiority:

One-shot methods can also do finetuning, but still parameters of the optimal networks are updated less frequently than DSNAS

---

# Outline

---

- Motivation
  - Search space
  - Method
  - Experiments
-

# Experiment

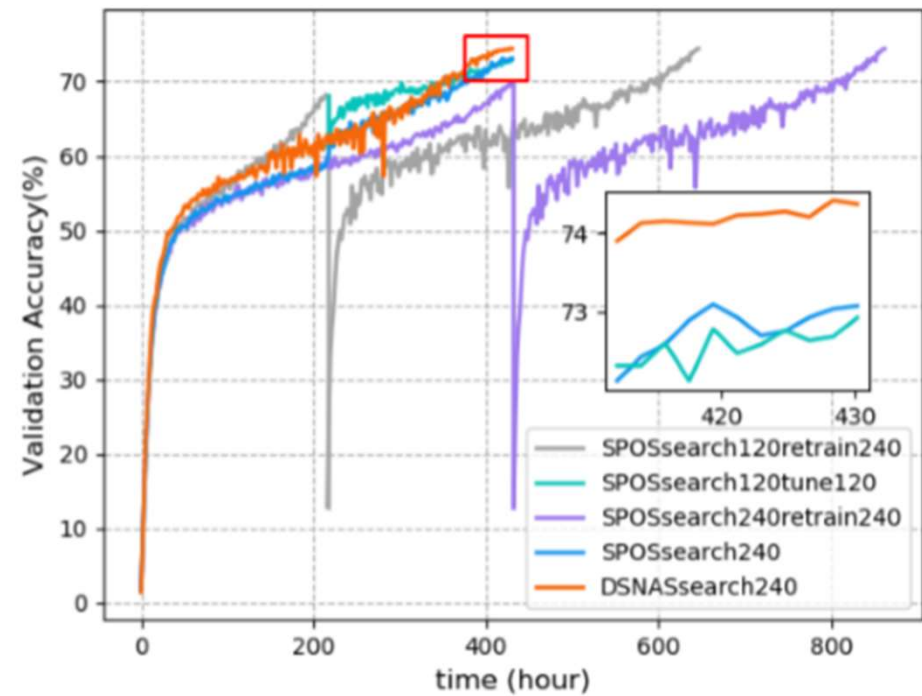
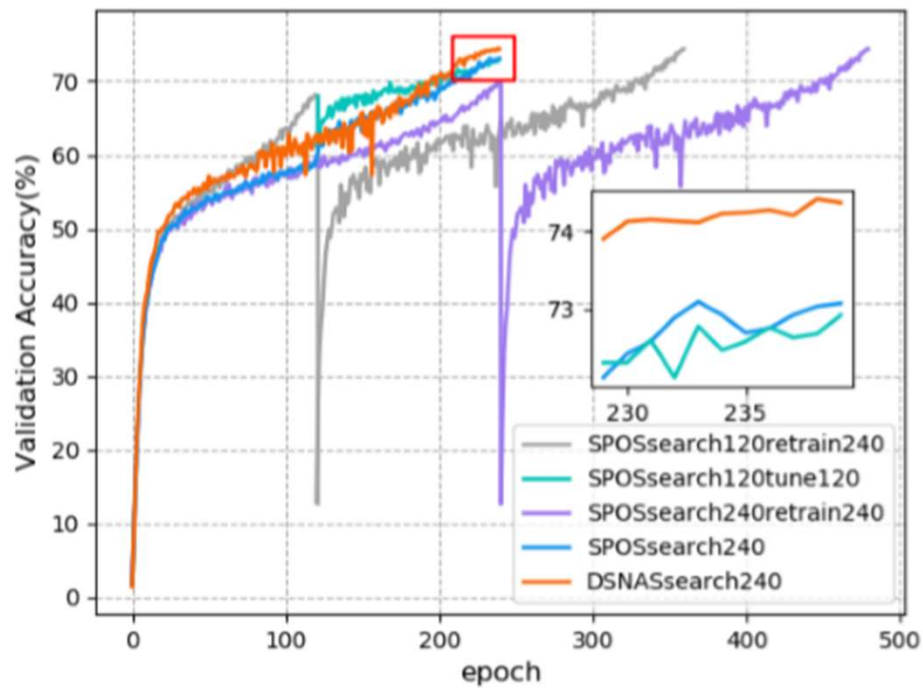
---

$$\tau = \frac{2(N_{\text{concordant}} - N_{\text{disconcordant}})}{N(N - 1)}, \quad (9)$$

Model	$\tau_{\text{inter}}$	$\tau_{\text{intra}}$
Single Path One-Shot[7]	0.33	-0.07
ProxylessNas [3]	-	-0.33

# Experiment

---



# Experiment

Model	FLOPS	Search		Retrain		No proxy	Time (GPU hour)	
		Top-1 acc(%)	Top-5 acc(%)	Top-1 acc(%)	Top-5 acc(%)		Search	Retrain
MobileNet V1 (0.75x)[8]	325M	Manual		68.4	-		Manual	
MobileNet V2 (1.0x)[19]	300M	Manual		72.0	91.00		Manual	
ShuffleNet V2 (1.5x)[15]	299M	Manual		72.6	90.60		Manual	
NASNET-A(4@1056)[29]	564M	-	-	74.0	91.60	False	48000	-
PNASNET[13]	588M	-	-	74.2	91.90	False	5400	-
MnasNet-A1[23]	312M	-	-	75.2	92.50	False	40000 <sup>§</sup>	-
DARTS[14]	574M	-	-	73.3	91.30	False	24	288
SNAS[27]	522M	-	-	72.7	90.80	False	36	288
Proxyless-R (mobile)[3]	320M	62.6*	84.7*	74.6	92.20	True	300 <sup>‡</sup>	≥384
Single Path One-Shot[7]	319M	68.7 <sup>†</sup>	-	74.3	-	True	250	384
Single Path One-Shot*	323M	68.2±0.38	88.28	74.3±0.36	91.79	True	250	384
Random Search	≤330M	≤68.2	≤88.31	≤73.9	≤91.8	True	250	384
DSNAS	324M	<b>74.4±0.22</b>	91.54	74.3±0.27	91.90	True	420	



# Experiment

---

Method	Train			Test
	Forward	Backward	Opt	
SNAS	0.26s	0.46s	0.14s	0.18s
DSNAS	0.05s	0.07s	0.13s	0.04s

Table 4. Computation time of SNAS and DSNAS

Method	Train			Test
	Forward	Backward	Opt	
ProxylessNAS	3.3s	2.3s	3.6s	1.2s
DSNAS	1.9s	1.3s	2.6s	0.9s

Table 5. Computation time of ProxylessNAS and DSNAS

---

---

**The End!**

---