

# DetNAS: Neural Architecture Search on Object Detection

Yukang Chen<sup>1†</sup>, Tong Yang<sup>2†</sup>, Xiangyu Zhang<sup>2</sup>, Gaofeng Meng<sup>1</sup>, Chunhong Pan<sup>1</sup>, Jian Sun<sup>2</sup>

<sup>1</sup>*National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences*

<sup>2</sup>*Megvii Inc*

{yukang.chen, gfmeng, chpan}@nlpr.ia.ac.cn {yangtong, zhangxiangyu, sunjian}@megvii.com

# Overview

- The performance of object detectors highly rely on features extracted by backbones. Most works on object detection directly use networks designed for classification as backbone feature extractors, e.g., ResNet.
- Image classification focus on "what" the main object of an image is, while object detection aims at finding "where" and "what" each object instance in an image.
- The architectures optimized on image classification can not guarantee the performance on object detection.

# Neural architecture search

- NAS has achieved great progress on image classification and semantic segmentation.
- Some NAS work directly applies the architecture searched on CIFAR-10 classification on object detection, which is sub-optimal.

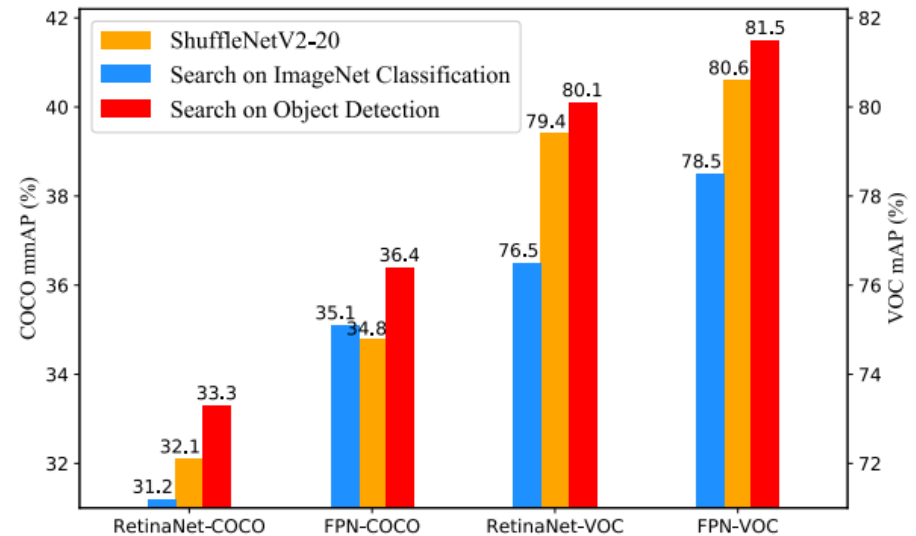


Figure 1. Backbone comparisons in various detectors and datasets. Architectures searched on the target task achieves consistent improvements over the hand-craft baseline and the ImageNet classification searched counterparts. More details are given in Table 2.

# DetNAS

- DetNAS is to learn a backbone network for object detection task, and conducts neural architecture search directly on target tasks.
- To improve computation efficiency, we formulate this problem into searching the optimal path in a large graph or supernet.
- Briefly speaking, DetNAS consists of three steps:
  - training a supernet that includes all sub-networks in search space;
  - searching for the sub-network with the highest performance on the validation set with EA;
  - retraining the result network and evaluating it on the test set.

# Overall pipeline

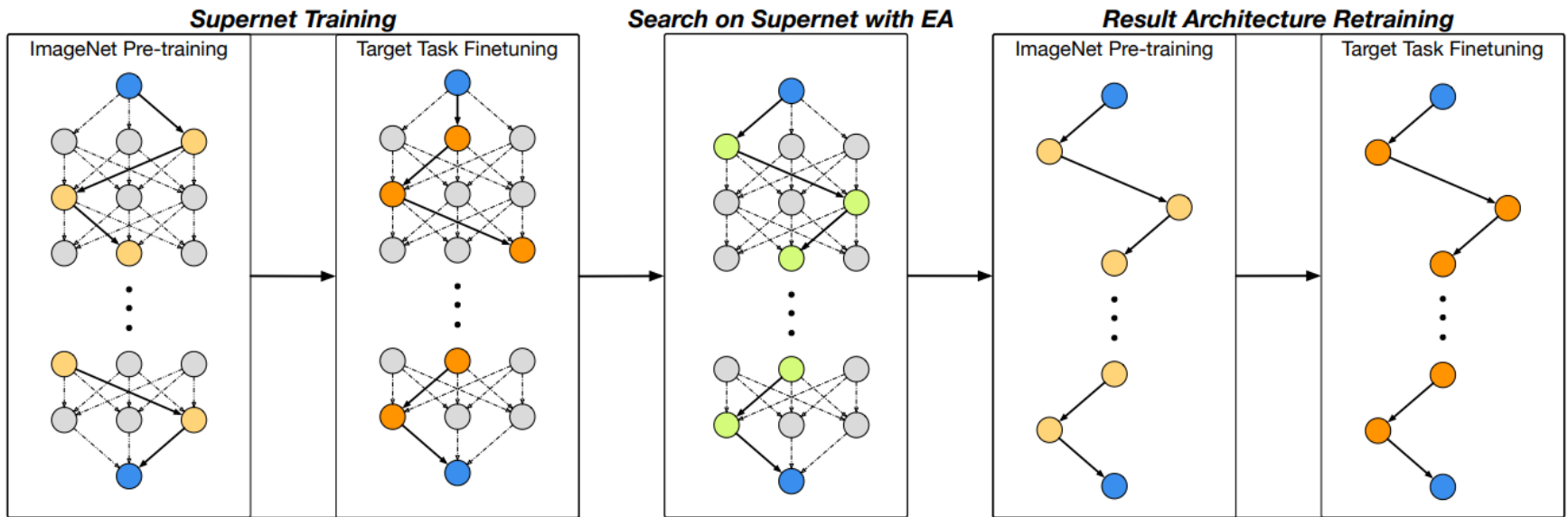


Figure 2. The pipeline of DetNAS. It involves three steps: supernet training, search on supernet with EA and result architecture retraining. The training step includes both ImageNet pre-training and target task finetuning. If we search in the from scratch scheme, no ImageNet pre-training needs during supernet training and the iterations on target task finetuning would be doubled. More details are listed in Table 1.

# Search Space Supernet

- Training object detectors is known to be computationally expensive.
- The search space is formulated into a large supernet graph to relief the issue of computation cost.

---

**Algorithm 1:** SuperNet Training

---

**Input** : Search space  $\mathbb{S}$ , Detector  $\mathcal{D}$ ,  
ImageNet pretrain iterations  $I_P$ ,  
Detection finetuning iterations  $I_F$ ,  
ImageNet pretrain training set  $\mathbb{T}_{Pre}$ ,  
Detection finetuning training set  $\mathbb{T}_{Det}$ .

```
1  $S \leftarrow \text{initialize}(\mathbb{S})$ 
2  $i \leftarrow 0$ 
3 while  $i < I_P$  do
4    $\theta \leftarrow \text{random-path}(\mathbb{S})$ 
5    $S \leftarrow \text{training}(S(\theta), \mathbb{T}_{Pre})$ 
6    $i \leftarrow i + 1$ 
7 end
8 while  $i < I_P + I_F$  do
9    $\theta \leftarrow \text{random-path}(\mathbb{S})$ 
10   $S \leftarrow \text{training}(\mathcal{D}(S(\theta)), \mathbb{T}_{Det})$ 
11   $i \leftarrow i + 1$ 
12 end
```

**Output:** The trained supernet model  $S$

---

# Search Space Supernet

- Key points:
  - The training process is path-wise  
In each iteration, only one single path is sampled for forward and back propagation. No gradient or weight update acts on other paths or nodes in the supernet graph.
  - Synchronized Batch Normalization (SyncBN) is necessary for finetuning on detection datasets instead of common Batch Normalization (BN).

---

**Algorithm 1:** SuperNet Training

---

**Input** : Search space  $\mathbb{S}$ , Detector  $\mathcal{D}$ ,  
ImageNet pretrain iterations  $I_P$ ,  
Detection finetuning iterations  $I_F$ ,  
ImageNet pretrain training set  $\mathbb{T}_{Pre}$ ,  
Detection finetuning training set  $\mathbb{T}_{Det}$ .

```
1  $S \leftarrow \text{initialize}(\mathbb{S})$ 
2  $i \leftarrow 0$ 
3 while  $i < I_P$  do
4    $\theta \leftarrow \text{random-path}(\mathbb{S})$ 
5    $S \leftarrow \text{training}(S(\theta), \mathbb{T}_{Pre})$ 
6    $i \leftarrow i + 1$ 
7 end
8 while  $i < I_P + I_F$  do
9    $\theta \leftarrow \text{random-path}(\mathbb{S})$ 
10   $S \leftarrow \text{training}(\mathcal{D}(S(\theta)), \mathbb{T}_{Det})$ 
11   $i \leftarrow i + 1$ 
12 end
```

**Output:** The trained supernet model  $S$

---

# Search Space Supernet

Table 2. ShuffleNetv2 search space backbone

Stage	Operation	$c_1$	$n_1$	$c_2$	$n_2$	$s$
0	Conv $3\times3$ -BN-ReLU	16	1	48	1	2
1	Shuffle block (search)	64	4	96	8	2
2	Shuffle block (search)	160	4	240	8	2
3	Shuffle block (search)	320	8	480	16	2
4	Shuffle block (search)	640	4	960	8	2

\* Each stage has  $n$  blocks. All blocks in one stage has output channels  $c$ .  
The first block of each stage has a stride  $s$  and all others use stride 1.

- The ShuffleNetv2 block is a efficient and lightweight convolution architectures. It involves channel split and shuffle operation.
- Search space 1: 4 stages, 20 ShuffleNetv2 blocks, 4 choices/block,  $4^{20}\sim10^{12}$  possible architectures



# Search Strategy

**1-3:** A population of networks  $P$  is initialized randomly. Each individual  $P$  in the population  $P$  consists of its architecture  $P.\theta$  and its fitness  $P.f$ .

**5-11:** After initialization, they start to evaluate the individual architecture  $P.\theta$  to obtain its fitness  $P.f$  on the validation set  $V_{\text{det}}$ . **(6)** Before that, they collect the batch normalization statistics for each network on a small subset of training set  $T_{\text{Dets}}$ .

**12-16:** Among these evaluated networks, we select the top  $|P|$  as parents to generate child networks. **(14-15)** This step involves two common transformations, mutation and crossover.

---

**Algorithm 2:** Search on Supernet with EA

---

**Input** : Trained supernet model  $S$ , Detector  $\mathcal{D}$ , Population size  $|P|$ , Parent size  $|P|$ , Detection validation set  $V_{\text{Det}}$ , Small subset of detection training set  $T_{\text{Dets}}$ , Evolution iterations  $I_E$ , FLOPs constraint  $\eta$ .

```
1  $P^{(0)} \leftarrow \text{random-initialize}(|P|, \eta)$ 
2  $j \leftarrow 0$ 
3  $f_{\text{best}}, \theta_{\text{best}} \leftarrow (0, \emptyset)$ 
4 while  $j < I_E$  do
5   for  $P \in P^{(j)}$  do
6      $S(P.\theta) \leftarrow \text{collectBN}(\mathcal{D}(S(P.\theta)), T_{\text{Dets}})$ 
7      $P.f \leftarrow \text{evaluate}(\mathcal{D}(S(P.\theta)), V_{\text{Det}})$ 
8     if  $P.f > f_{\text{best}}$  then
9        $f_{\text{best}}, \theta_{\text{best}} \leftarrow (P.f, P.\theta)$ 
10    end
11  end
12   $\mathcal{P} \leftarrow \text{select-topk}(P^{(j)}, |P|)$ 
13   $j \leftarrow j + 1$ 
14   $P_{\text{mut}} \leftarrow \text{mutate}(\mathcal{P}, \frac{|P|}{2}, \eta)$ 
15   $P_{\text{crs}} \leftarrow \text{crossover}(\mathcal{P}, \frac{|P|}{2}, \eta)$ 
16   $P^{(j)} \leftarrow \{P_{\text{mut}}, P_{\text{crs}}\}$ 
17 end
```

**Output:** The best architecture  $\theta_{\text{best}}$

---

# Search Strategy

- Key points:
  - The search process is also pathwise. To evaluate each path, validation data is only passed through the specific path.
  - Before the path is evaluated on the validation set, training data needs to be passed on the path for several iterations. This involves only forward propagation and no gradient backward. It aims to compute batch statistics.

---

**Algorithm 2:** Search on Supernet with EA

---

**Input** : Trained supernet model  $S$ , Detector  $\mathcal{D}$ ,  
Population size  $|\mathbf{P}|$ , Parent size  $|\mathcal{P}|$ ,  
Detection validation set  $\mathbb{V}_{Det}$ , Small subset  
of detection training set  $\mathbb{T}_{Det_S}$ , Evolution  
iterations  $I_E$ , FLOPs constraint  $\eta$ .

```
1  $\mathbf{P}^{(0)} \leftarrow \text{random-initialize}(|\mathbf{P}|, \eta)$ 
2  $j \leftarrow 0$ 
3  $f_{best}, \theta_{best} \leftarrow (0, \emptyset)$ 
4 while  $j < I_E$  do
5   for  $P \in \mathbf{P}^{(j)}$  do
6      $S(P.\theta) \leftarrow \text{collectBN}(\mathcal{D}(S(P.\theta)), \mathbb{T}_{Det_S})$ 
7      $P.f \leftarrow \text{evaluate}(\mathcal{D}(S(P.\theta)), \mathbb{V}_{Det})$ 
8     if  $P.f > f_{best}$  then
9        $f_{best}, \theta_{best} \leftarrow (P.f, P.\theta)$ 
10    end
11  end
12   $\mathcal{P} \leftarrow \text{select-topk}(\mathbf{P}^{(j)}, |\mathcal{P}|)$ 
13   $j \leftarrow j + 1$ 
14   $\mathbf{P}_{mut} \leftarrow \text{mutate}(\mathcal{P}, \frac{|\mathbf{P}|}{2}, \eta)$ 
15   $\mathbf{P}_{crs} \leftarrow \text{crossover}(\mathcal{P}, \frac{|\mathbf{P}|}{2}, \eta)$ 
16   $\mathbf{P}^{(j)} \leftarrow \{\mathbf{P}_{mut}, \mathbf{P}_{crs}\}$ 
17 end
```

**Output:** The best architecture  $\theta_{best}$

---

# Batch Normalization

- Typically, the parameters of BN, during fine-tuning, are fixed as the pre-training batch statistics.
- In DetNAS, the supernet is trained in a path-wise manner. The batch statistics of one node are obviously not same on different paths.
- When finetuning the supernet on the detection dataset, they replace the conventional BN with SyncBN to compute batch statistics across multiple GPUs. This enables the supernet finetuning by increasing the effective batch size for BN.

# Results

Table 3. Results of different backbones

Backbone	ImageNet Classification		Object Detection with FPN on COCO					
	FLOPs	Accuracy	mAP	$AP_{50}$	$AP_{75}$	$AP_s$	$AP_m$	$AP_l$
ResNet-50	3.8G	76.15	36.7	58.4	39.6	21.1	39.8	48.1
ResNet-101	7.6G	<b>77.37</b>	39.4	61.2	43.4	22.6	42.9	51.4
DetNet-59	4.8G	76.50	37.9	60.1	41.2	22.7	41.2	48.3
ShuffleNetV2-40	<b>1.3G</b>	77.18	39.2	60.8	42.4	23.6	42.3	52.2
DetNASNet	<b>1.3G</b>	77.20	<b>40.0</b>	61.5	43.6	23.3	42.5	53.8

Table 4. Comparisons in ShuffleNetV2 search space.

	FPN		RetinaNet	
	VOC (mAP, %)	COCO (mmAP, %)	VOC (mAP, %)	COCO (mmAP, %)
ShuffleNetV2-20	80.6 <sub>(73.1)</sub>	34.8 <sub>(73.1)</sub>	79.4 <sub>(73.1)</sub>	32.1 <sub>(73.1)</sub>
ClsNASNet	78.5 <sub>(74.3)</sub>	35.1 <sub>(74.3)</sub>	76.5 <sub>(74.3)</sub>	31.2 <sub>(74.3)</sub>
Random	80.9±0.2 <sub>(73.9±0.2)</sub>	35.9±0.3 <sub>(73.9±0.2)</sub>	79.0±0.7 <sub>(73.9±0.2)</sub>	32.5±0.4 <sub>(73.9±0.2)</sub>
DetNAS-S	81.1 <sub>(73.8)</sub>	<b>35.9</b> <sub>(74.3)</sub>	79.9 <sub>(74.2)</sub>	32.8 <sub>(74.2)</sub>
DetNAS-P	<b>81.5</b> <sub>(74.0)</sub>	<b>36.4</b> <sub>(74.0)</sub>	<b>80.1</b> <sub>(74.1)</sub>	<b>33.3</b> <sub>(73.9)</sub>

\* The subscript numbers are the ImageNet classification accuracy of the pretrained models. All models in this figure are trained with hyper-parameters described in Section 4.1. Their FLOPs are all similar and under 300M.

# Summary

- The authors propose DetNAS to search architectures for backbones of detectors.
- The authors prove the effectiveness of DetNAS with multidimensional experiments. It works in different detectors (two-stage FPN and one-stage RetinaNet), various datasets (COCO and VOC) and different scheme (ImageNet pre-training and training from scratch).
- The searched networks have consistently better performance than the architectures searched from ImageNet on the target task. DetNASNet, the main result architecture, achieves 40.0 mAP on COCO and outperforms ResNet-101 with 5x less FLOPs complexity.