



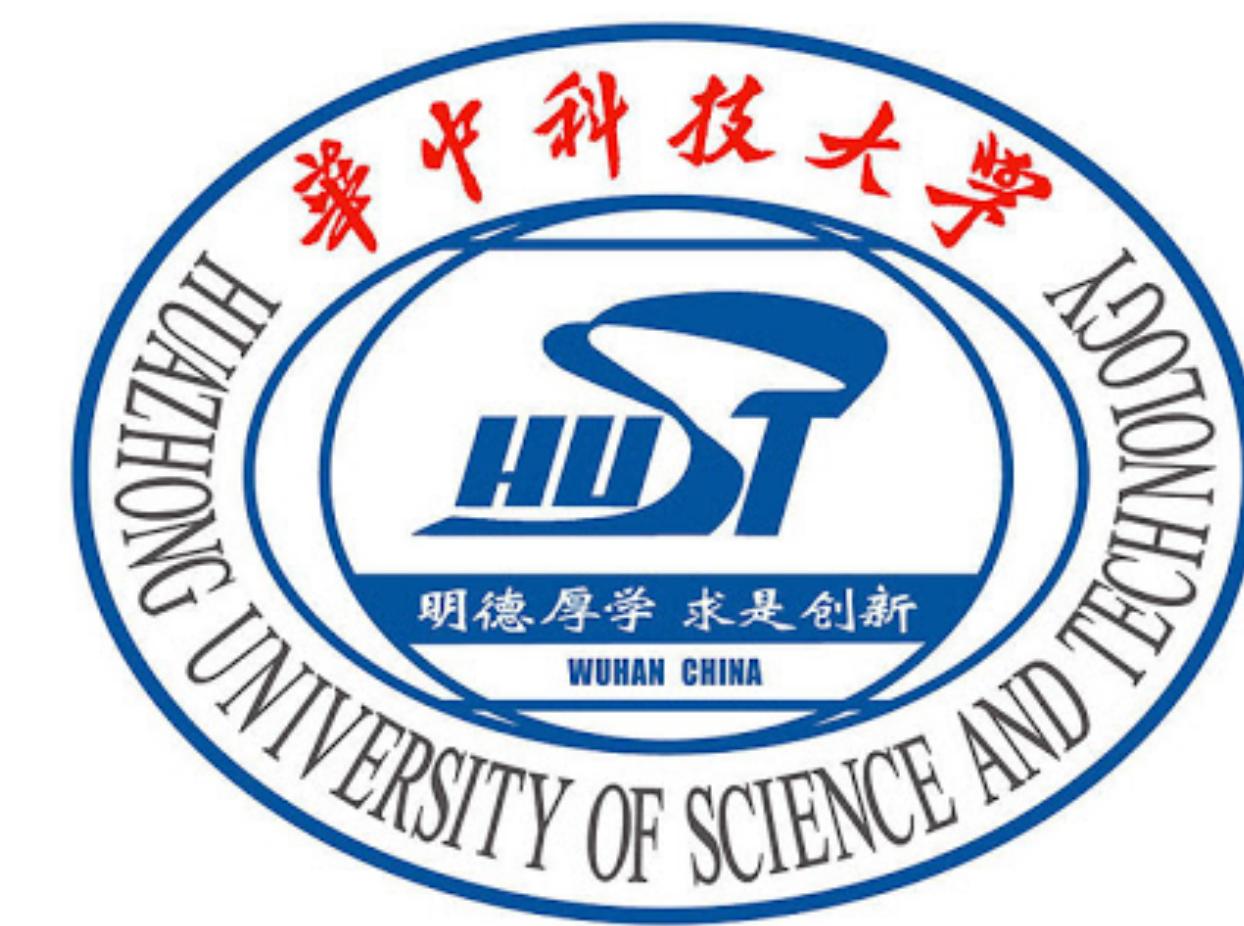
# GreedyNAS: Towards Fast One-Shot NAS with Greedy Supernet

Shan You<sup>1,2\*</sup>, Tao Huang<sup>1,3\*</sup>, Mingmin Yang<sup>1\*</sup>, Fei Wang<sup>1</sup>, Chen Qian<sup>1</sup>, Changshui Zhang<sup>2</sup>

<sup>1</sup>SenseTime <sup>2</sup>Department of Automation, Tsinghua University <sup>3</sup>Dian Group, School of CST,

Huazhong University of Science and Technology \*Equal contributions

{youshan,huangtao,yangmingmin,wangfei,qianchen}@sensetime.com zcs@mail.tsinghua.edu.cn



## Motivation

**One-shot NAS:** Based on the weight-sharing paradigm, One-shot NAS methods model NAS as a one-shot training process of an over-parameterized supernet, where various architectures can be directly derived.

**Supernet:** matters as a fundamental performance estimator of different architectures (paths).

**Target Assumption:** The supernet should estimate the (relative) performance accurately for *all* paths, and thus all paths are treated equally and trained simultaneously.

### Issues:

- It is harsh for a single supernet to evaluate accurately on such a huge-scale search space (e.g.,  $7^{21}$ ).
- There exist many architectures of inferior quality in terms of accuracy performance.
- Since the weights of all paths are highly shared, if a weak path is sampled and gets trained, it would disturb the weights of those potentially-good paths.
- Training on those weak paths actually involves unnecessary update of weights, and slows down the training efficiency more or less.

## Intuition

- block the training of these weak paths.
- Consider a complete partition of search space  $\mathcal{A}$  of two subsets  $\mathcal{A}_{good}$  and  $\mathcal{A}_{weak}$ :

$$\mathcal{A} = \mathcal{A}_{good} \cup \mathcal{A}_{weak}, \quad \mathcal{A}_{good} \cap \mathcal{A}_{weak} = \emptyset,$$

where for an Oracle supernet  $\mathcal{N}_o$ ,

$$ACC(\mathbf{a}, \mathcal{N}_o, \mathcal{D}_{val}) \geq ACC(\mathbf{b}, \mathcal{N}_o, \mathcal{D}_{val})$$

holds for all  $\mathbf{a} \in \mathcal{A}_{good}$ ,  $\mathbf{b} \in \mathcal{A}_{weak}$  on validation dataset  $\mathcal{D}_{val}$ .

- Idea: just sample from the potentially-good paths  $\mathcal{A}_{good}$  instead of all paths  $\mathcal{A}$ ,

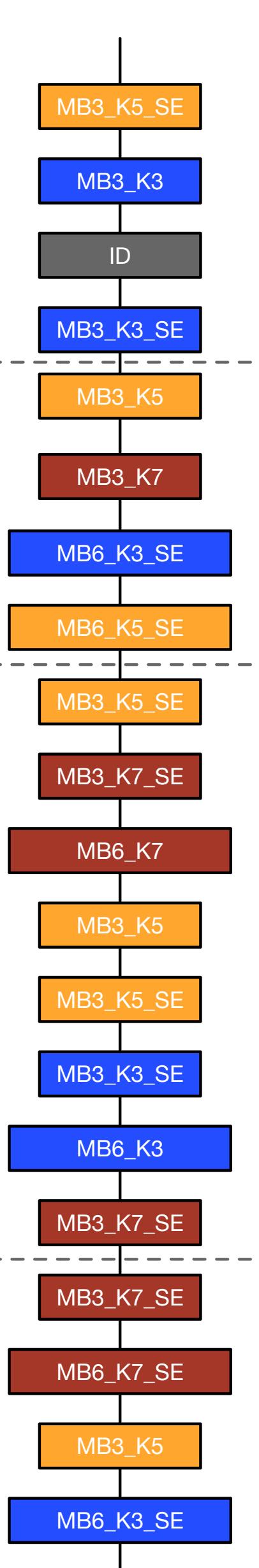
$$p(\mathbf{a}; \mathcal{N}_o, \mathcal{D}_{val}) = \frac{1}{|\mathcal{A}_{good}|} \mathbb{I}(\mathbf{a} \in \mathcal{A}_{good}).$$

### Problems:

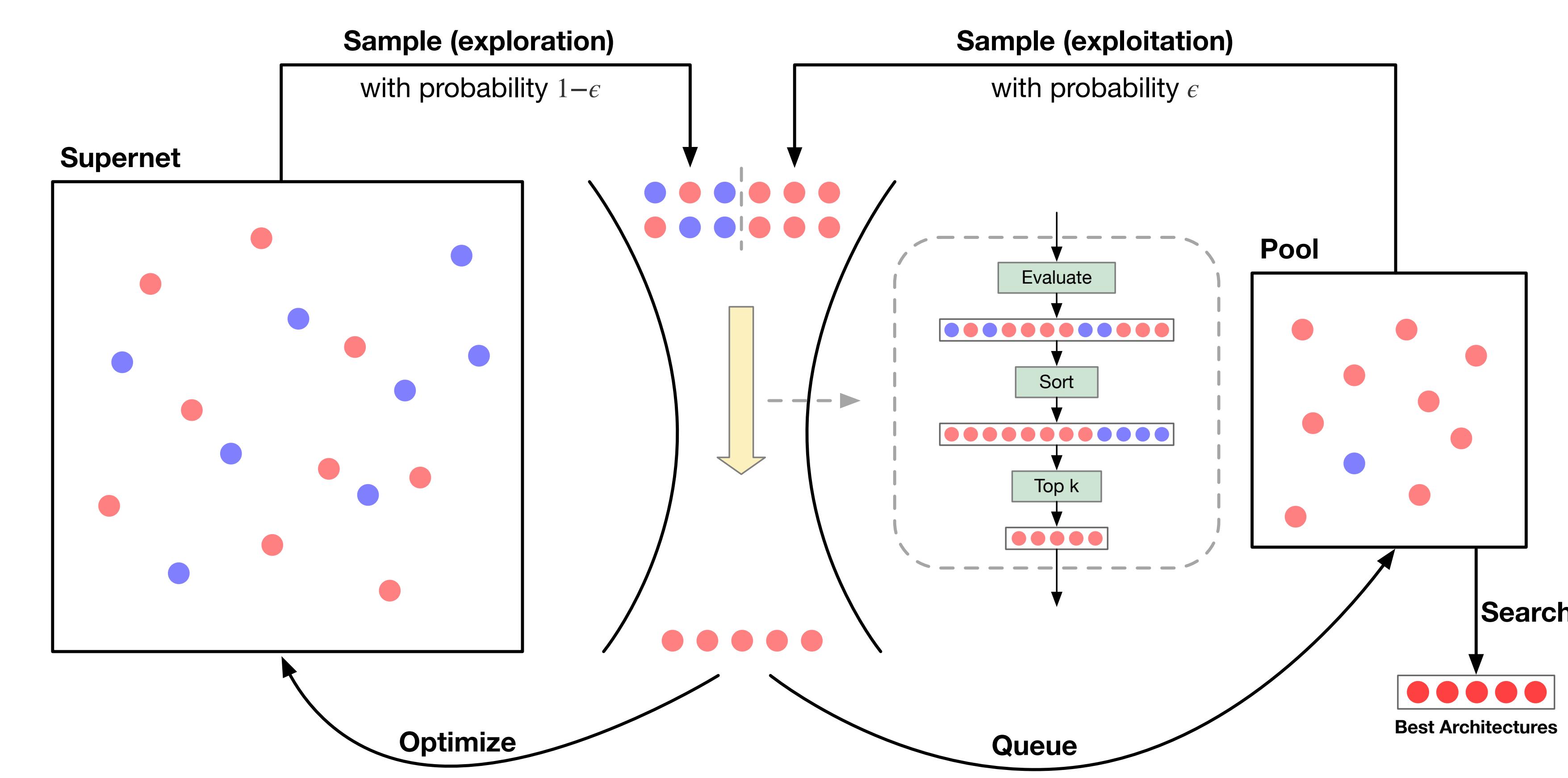
- Q: Oracle supernet  $\mathcal{N}_o$  is unknown.  
A: greedily use current supernet  $\mathcal{N}_t$  as a proxy
- Q: How can we accurately identify whether a path is from  $\mathcal{A}_{good}$  or  $\mathcal{A}_{weak}$ ?  
A: multi-path sampling with rejection.

## Experimental Settings

- Search space: MobileNetV2 inverted bottleneck with CNN kernel  $\{3, 5, 7\}$  and expansion ratio  $\{3, 6\}$ . Size  $7^{21}$  with identity. Larger size  $13^{21}$  with SE.
- ImageNet dataset: 50K validation, 50K testing
- Supernet: sample 10 paths and filter 5, 1K images for path filtering, pool size 1K, SGD optimizer
- Searching: evolutionary NSGA-II
- Retraining: following Proxyless-NAS without SE and Mnasnet with SE.



## Framework of GreedyNAS



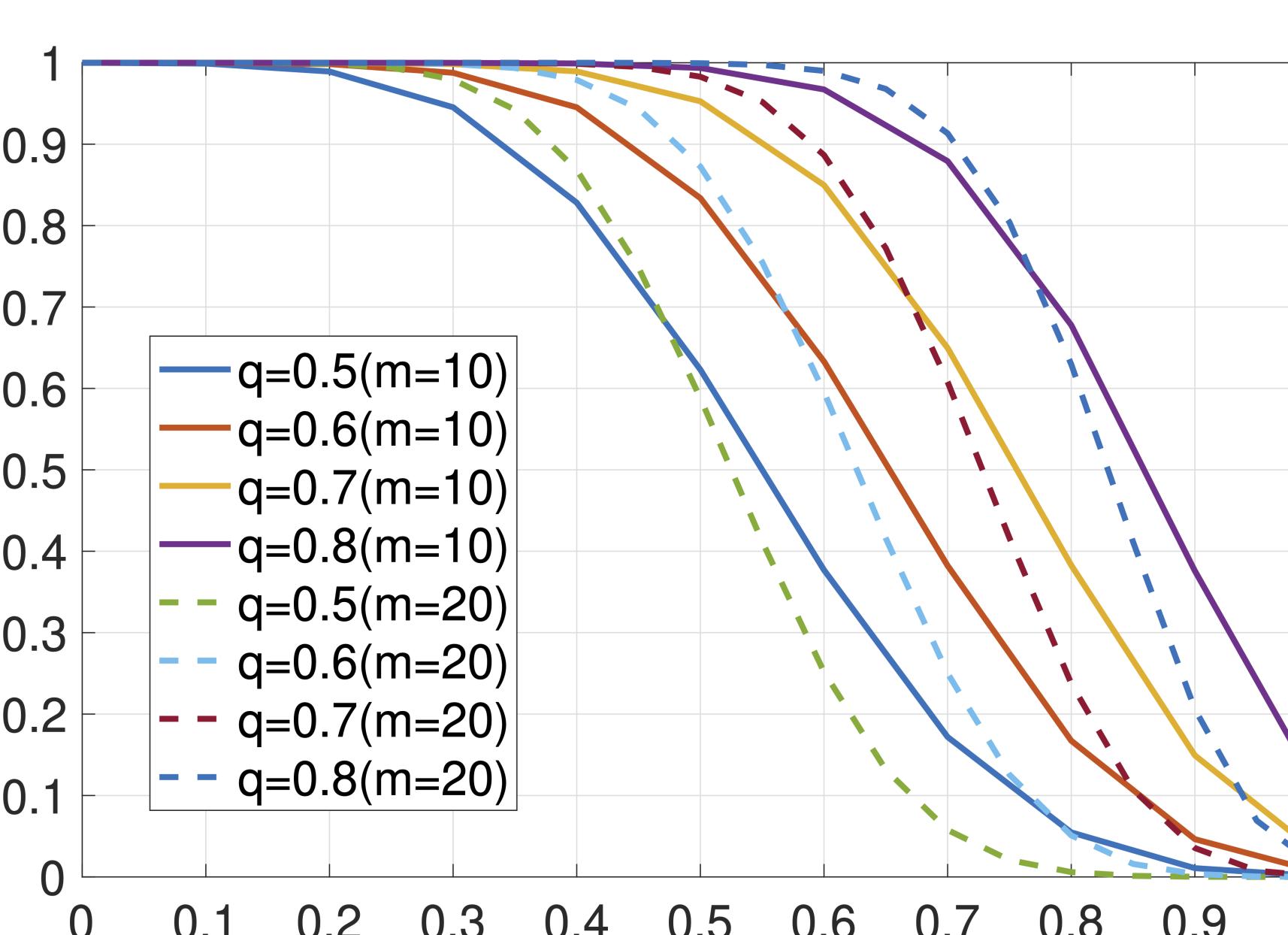
The supernet greedily shrinks its training space from all paths (red and blue dots) into potentially-good paths (red dots), and further into candidate pool.

## Multi-path Sampling with Rejection

**Theorem 1** If  $m$  paths are sampled uniformly i.i.d. from  $\mathcal{A}$ , then it holds that at least  $k$  ( $k \leq m$ ) paths are from  $\mathcal{A}_{good}$  with probability

$$\sum_{j=k}^m C_m^j q^j (1-q)^{m-j}, \quad (1)$$

where  $q = |\mathcal{A}_{good}| / |\mathcal{A}|$ .



Eq.(1) can be large; with  $q = 0.6$ , it has 83.38% confidence to say at least 5 out of 10 paths are from  $\mathcal{A}_{good}$ .

Solution: just rank the sampled  $m$  paths using a small portion of validation data  $\mathcal{D}_{val}$ , keep the Top- $k$  paths and reject the remaining paths.

**Input:** number of sampled multiple paths  $m$ , number of kept paths  $k$ , candidate pool  $\mathcal{P}$

- if without candidate pool  $\mathcal{P}$  then
- sample  $m$  paths  $\{\mathbf{a}_i\}_{i=1}^m$  i.i.d. w.r.t.  $\mathbf{a}_i \sim U(\mathcal{A})$
- else
- sample  $m$  paths  $\{\mathbf{a}_i\}_{i=1}^m$  i.i.d. w.r.t.  $\mathbf{a}_i \sim (1-\epsilon) \cdot U(\mathcal{A}) + \epsilon \cdot U(\mathcal{P})$
- end if
- randomly sample a batch  $\hat{\mathcal{D}}_{val}$  in  $\mathcal{D}_{val}$
- evaluate the loss  $\ell_i$  of each path  $\mathbf{a}_i$  on  $\hat{\mathcal{D}}_{val}$
- rank the paths by  $\ell_i$ , and get Top- $k$  indexes  $\{t_i\}_{i=1}^k$
- return  $k$  paths  $\{\mathbf{a}_{t_i}\}_{i=1}^k$  and filter the rest

## Exploration and Exploitation Training with Candidate Path Pool

We introduce a candidate path pool to store the discovered good paths, and sample from it,

$$\mathbf{a} \sim (1-\epsilon) \cdot U(\mathcal{A}) + \epsilon \cdot U(\mathcal{P}), \quad (2)$$

### Four advantages:

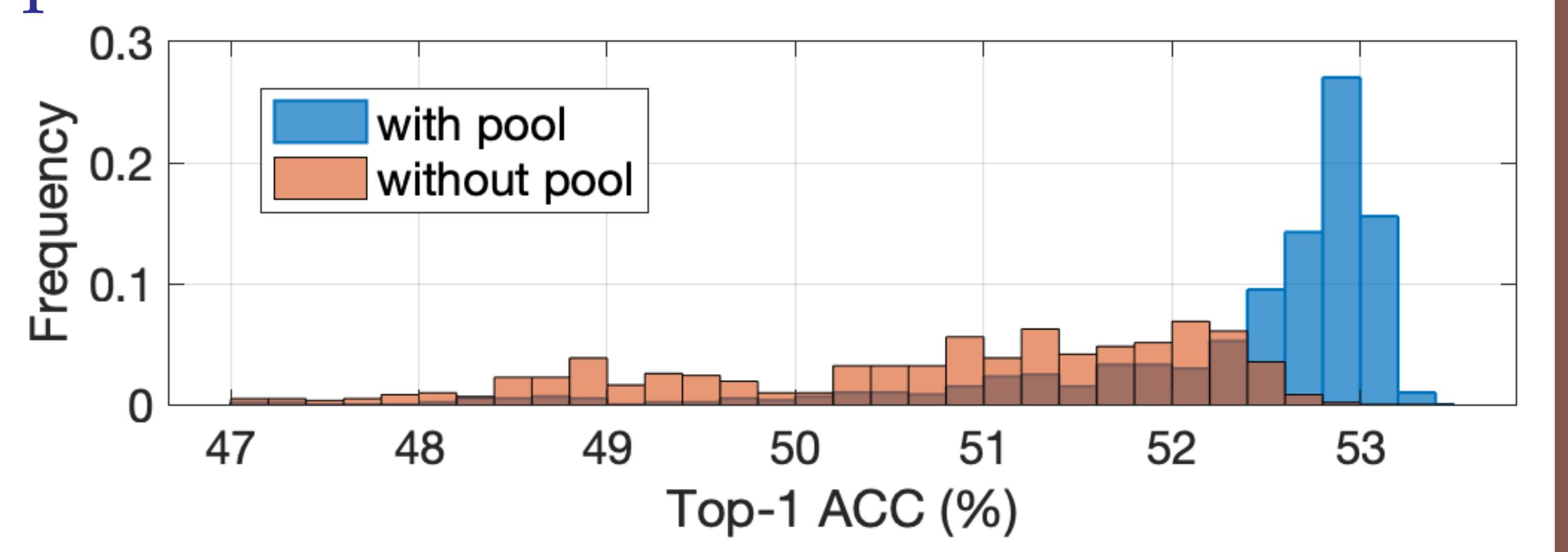
- boosting the training efficiency
- increasing the probability of sampling good paths  $q = \epsilon + (1-\epsilon)|\mathcal{A}_{good}| / |\mathcal{A}|$ , e.g. from 83.38% to 99.36% for 5/10 with  $\epsilon = 0.5$

- stopping principle via candidate pool

Stop by observing the steadiness of pool

$$\pi := \frac{|\mathcal{P}_t \cap \mathcal{P}|}{|\mathcal{P}|} \leq \alpha$$

- searching by initialization with candidate pool



## Searching Results with Same Search Space on ImageNet

Methods	performance			supernet training efficiency		
	Top-1 (%)	FLOPs	latency	#optimization	#evaluation	corrected #optimization
Proxyless-R (mobile)	74.60	320M	79 ms	-	-	-
Random Search	74.07	321M	69 ms	1.23M × 120	-	147.6M
Uniform Sampling	74.50	326M	72 ms	1.23M × 120	-	147.6M
FairNAS-C	74.69	321M	75 ms	1.23M × 150	-	184.5M
Random Search-E	73.88	320M	91 ms	1.23M × 73	-	89.8M
Uniform Sampling-E	74.17	320M	94 ms	1.23M × 73	-	89.8M
<b>GreedyNAS</b>	<b>74.85</b>	320M	89 ms	1.23M × 46	2.40M × 46	<b>89.7M</b>
<b>GreedyNAS</b>	<b>74.93</b>	324M	78 ms	1.23M × 46	2.40M × 46	<b>89.7M</b>

## Comparison with State-of-the-art NAS Methods on ImageNet

Methods	Top-1 (%)	FLOPs (M)	latency (ms)	Params (M)	training (Gdays)	search (Gdays)
SCARLET-C	75.6	280	67	6.0	10	12
MnasNet-A1	75.2	312	55	3.9	288 <sup>‡</sup>	-
<b>GreedyNAS-C</b>	<b>76.2</b>	284	70	4.7	7	< 1
FairNAS-C	74.7	321	75	4.4	10	2
SCARLET-B	76.3	329	104	6.5	10	12
<b>GreedyNAS-B</b>	<b>76.8</b>	324	110	5.2	7	< 1
SCARLET-A	76.9	365	118	6.7	10	12
EfficientNet-B0	76.3	390	82	5.3	-	-
DARTS	73.3	574	-	4.7	4 <sup>†</sup>	-
<b>GreedyNAS-A</b>	<b>77.1</b>	366	77	6.5	7	< 1

Rank correlation coefficient of 1000 paths measured by the loss (ACC) of 1K vs 50K validation images w.r.t. different types of supernets.

Spearman rho	Kendall tau
random uniform(ACC)	greedy
0.155	0.968(0.869)
<b>0.997</b>	0.113
0.851(0.699)	<b>0.961</b>

