

# **COMP5212: Machine Learning**

**Lecture 20**

**Minhao Cheng**

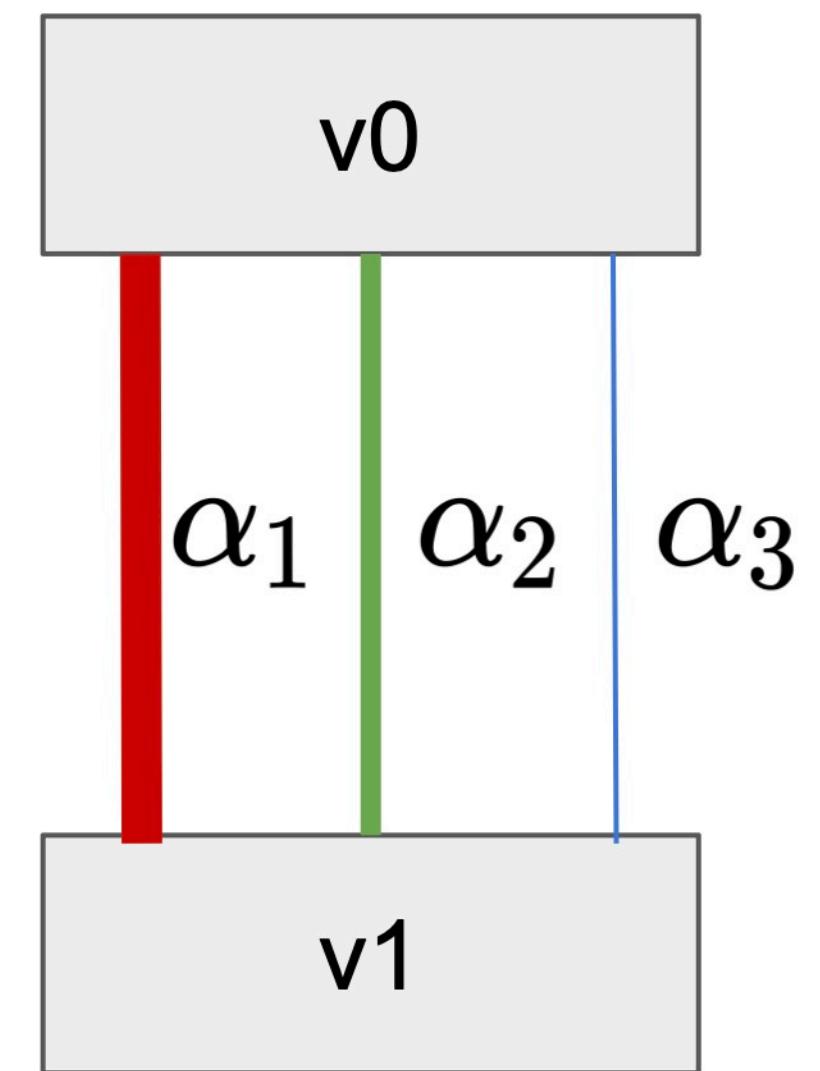
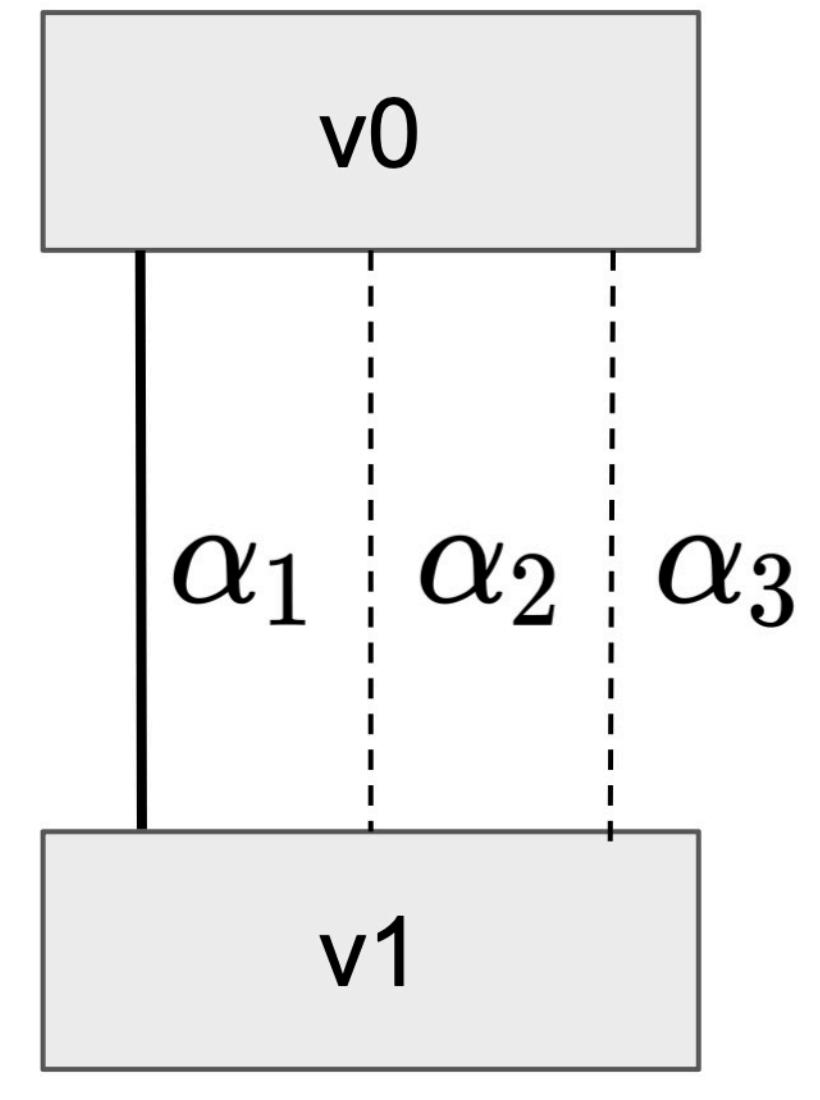
# Differentiable NAS

## Continuous Relaxation

- Final architecture:  $[\alpha_1, \alpha_2, \alpha_3]$  is a one-hot vector
- Relax to continuous values in the search phase=> Bi-level optimization for finding  $\alpha$

$$\min_{\alpha} L_{\text{val}}(w^*(\alpha), \alpha)$$

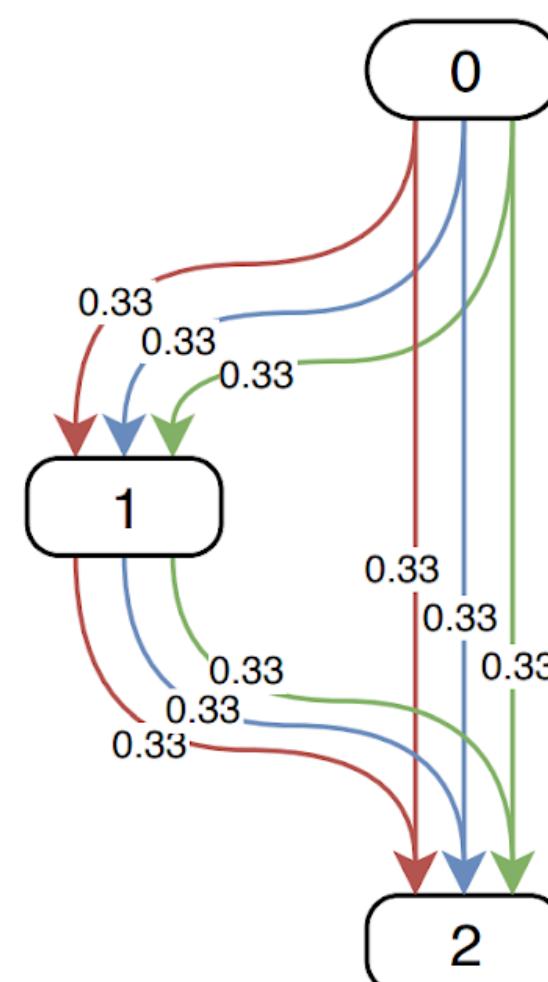
$$\text{s.t. } w^*(\alpha) = \arg \min_w L_{\text{train}}(w, \alpha)$$



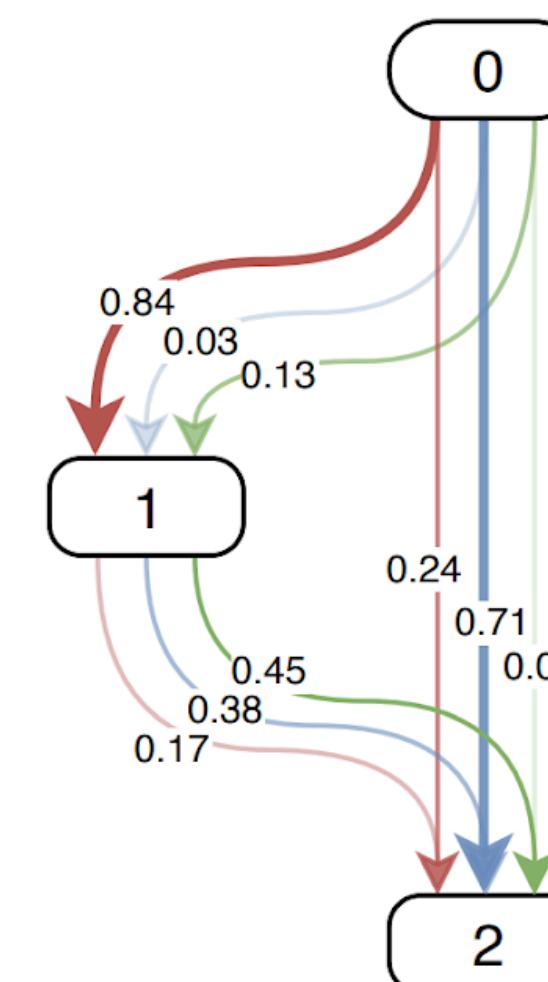
# Differentiable NAS

## Differentiable Neural Architecture Search (DARTS)

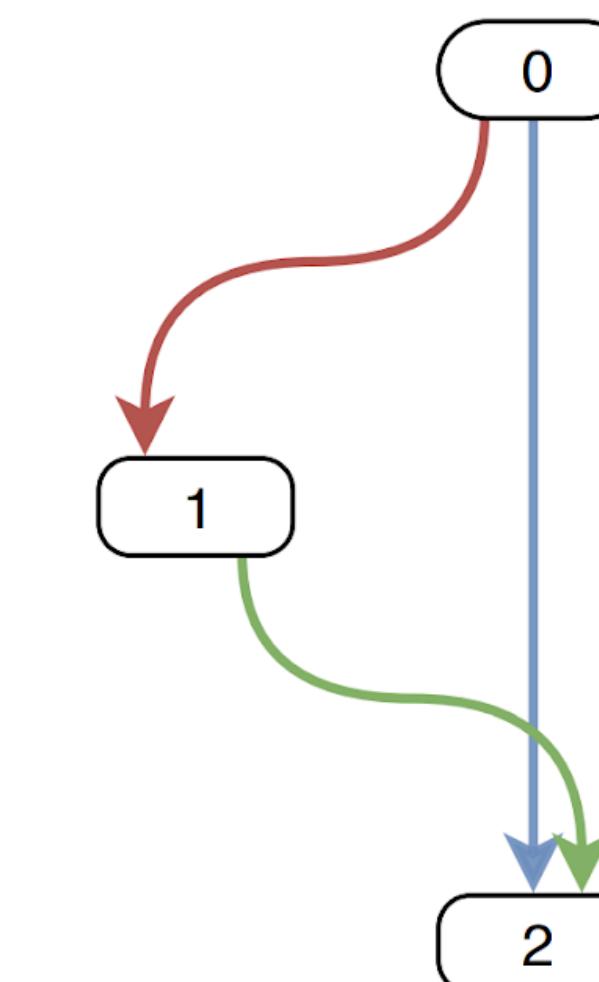
- Solve the bi-level optimization problem to obtain  $(\alpha^*, w^*)$  (supernet)
- Use magnitude of  $\alpha^*$  to choose the final architecture



(d) Search start



(e) Search end



(f) Final cell

# Differentiable NAS

## How to solve bi-level optimization?

$$\min_{\alpha} L_{\text{val}}(w^*(\alpha), \alpha)$$

$$\text{s.t. } w^*(\alpha) = \arg \min_w L_{\text{train}}(w, \alpha)$$

- Iteratively update  $w$  and  $\alpha$
- Update  $w$ :
  - Time consuming to compute  $w^*$  exactly => approximate by one SGD step
    - $w' \leftarrow w - \eta \nabla_w L_{\text{train}}(w, \alpha)$
- Update  $\alpha$ :
  - First order DARTS: assume  $w$  is constant w.r.t.  $\alpha$ 
    - $\alpha \leftarrow \alpha - c \nabla_\alpha L_{\text{val}}(w', \alpha)$

# Differentiable NAS

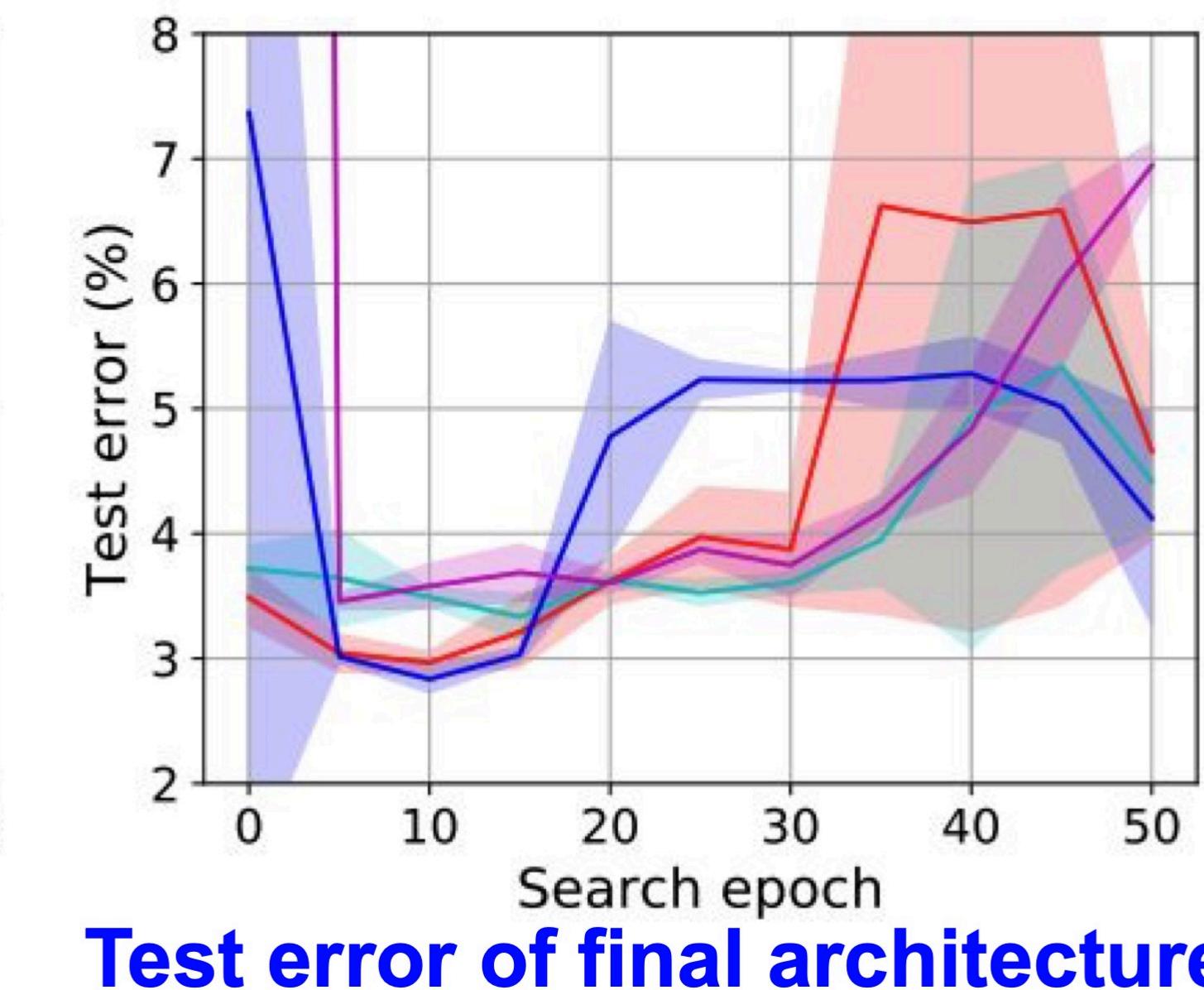
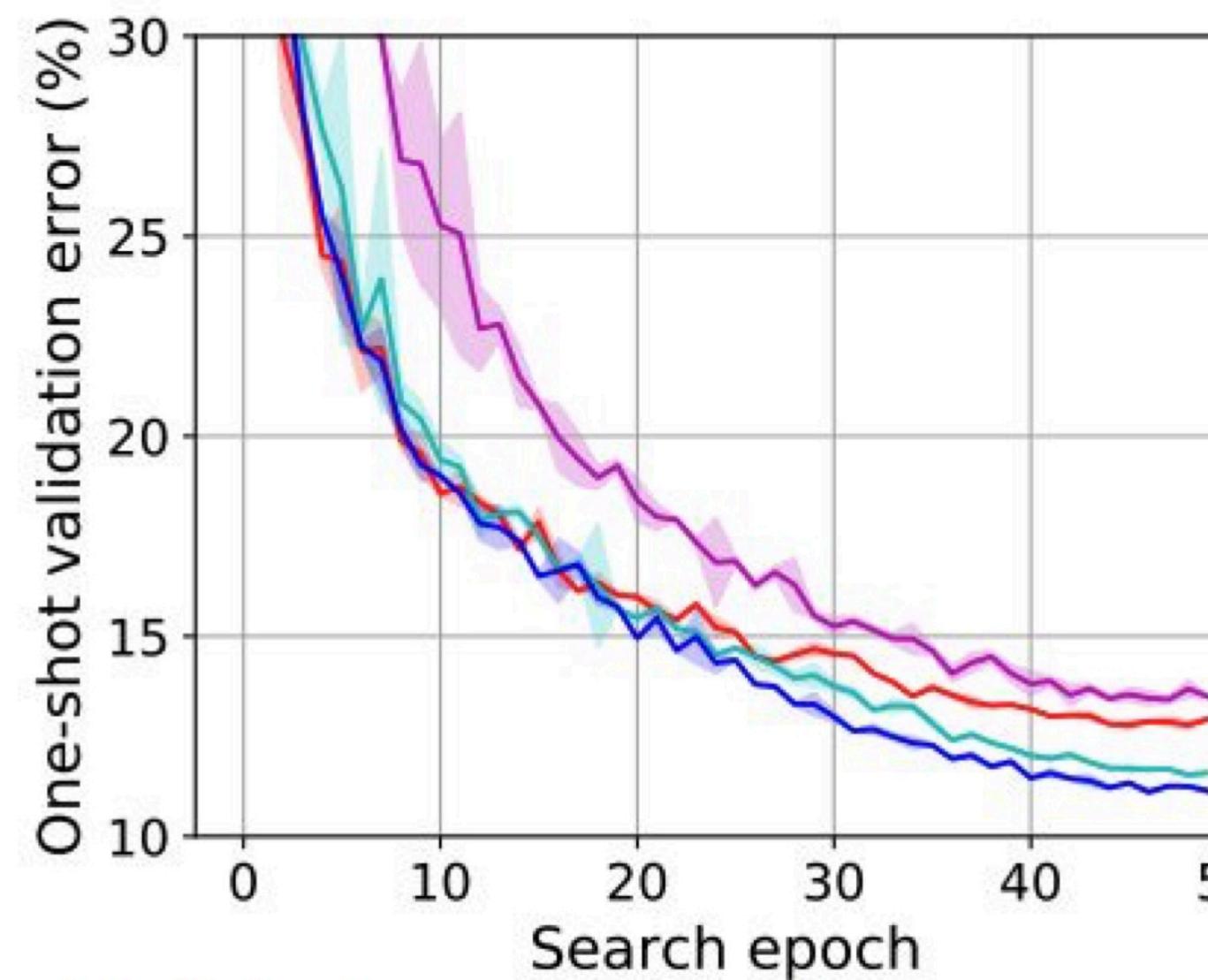
## Complexity of DARTS

- Time complexity: training the supernet **only once**
  - Supernet is a network with  $K$  operations with each edge  
=> **only  $K$  times slower than standard training**
  - Usually good enough
- Memory complexity (GPU memory):
  - Backprop on all the operations on each edge  
=>  **$K$  times memory consumption**
  - Prohibits for many problems

# Differentiable NAS

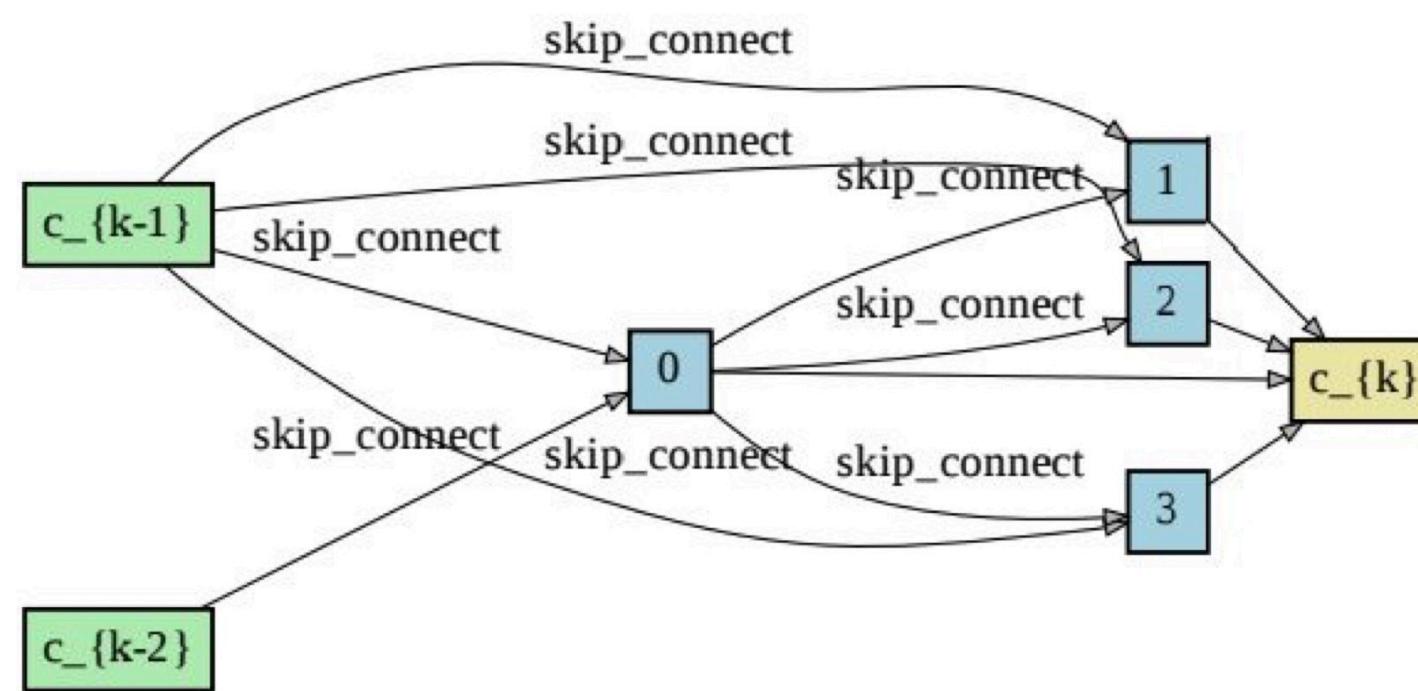
## DARTS fails in many simple cases

- Space 1: 2 operations per edge (selected from the original DARTS supernet)
- Space 2: 2 operations per edge {Conv3x3, skip\_connect}
- Space 3: 3 operations per edge {Conv3x3, skip\_connect, Zero}
- Space 4: 2 operations per edge {Conv3x3, Gaussian\_noise}

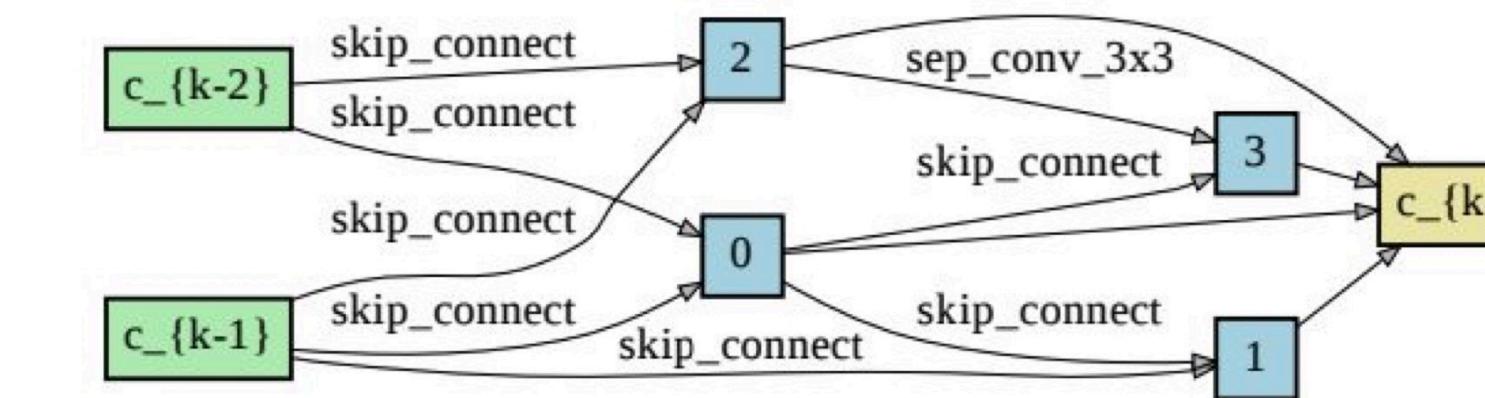


# Differentiable NAS

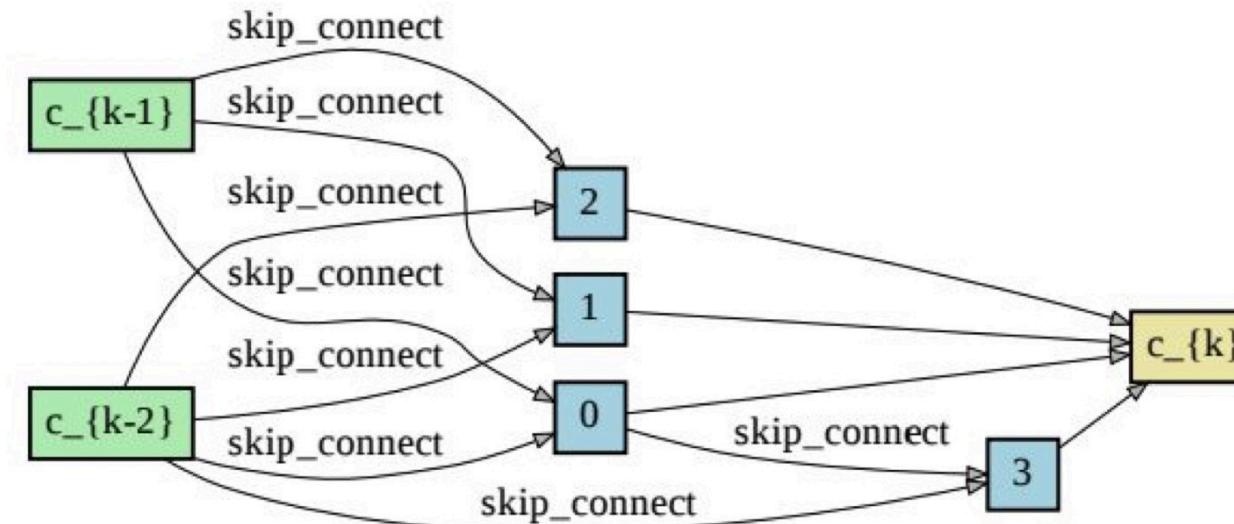
## DARTS leads to degenerated solutions



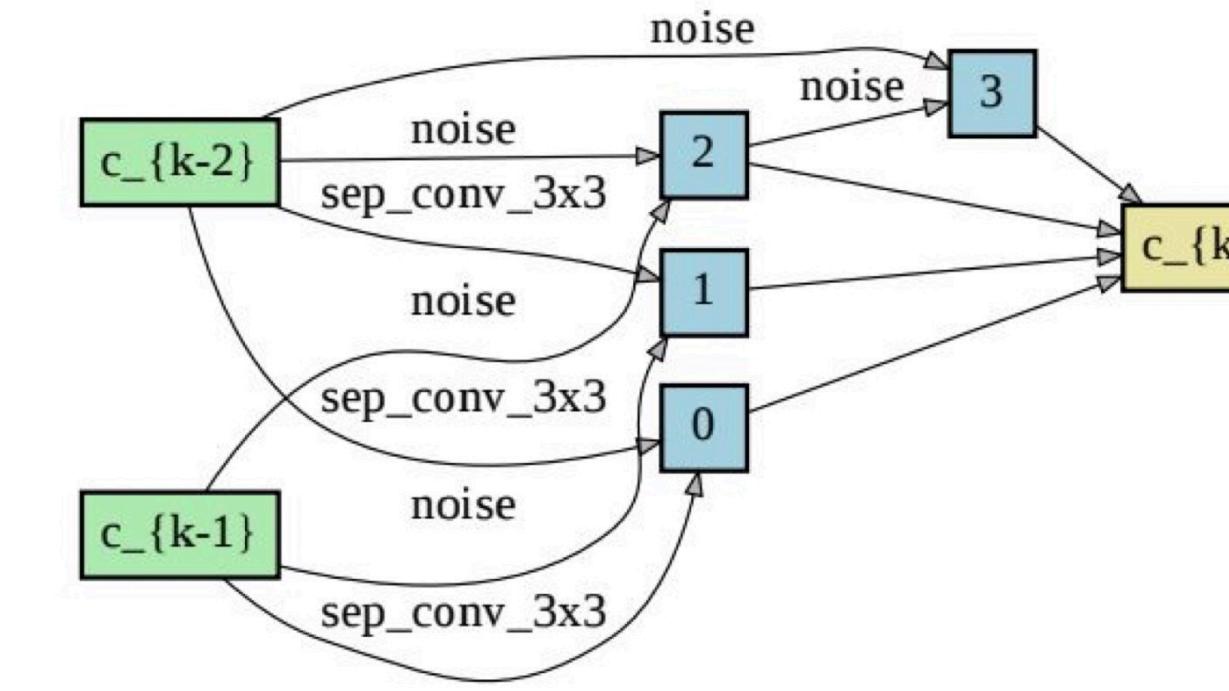
S1



S2



S3



S4

# Differentiable NAS

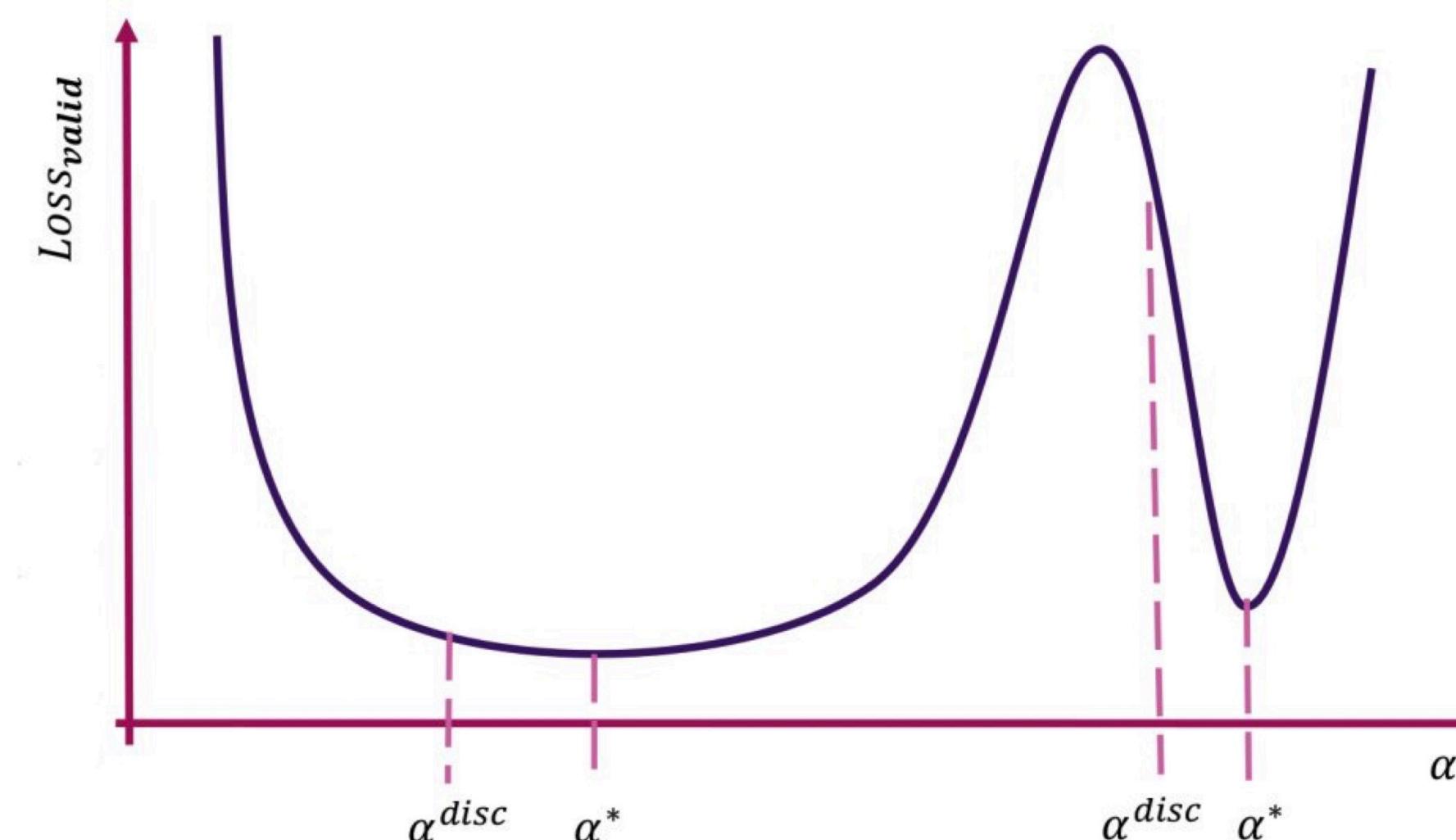
## Reason 1: sharpness of the solution

- A good **continuous solution** doesn't imply a good **discrete solution**
- Gap between continuous and discrete solutions can be estimated by sharpness
  - Assume  $\alpha^*$  is the continuous solution and  $\bar{\alpha}$  is the discrete solution
  - Based on Taylor expansion:
$$L_{\text{val}}(w^*, \bar{\alpha}) \approx L_{\text{val}}(w^*, \alpha^*) + \frac{1}{2}(\bar{\alpha} - \alpha^*)^T H(\bar{\alpha} - \alpha^*)$$
where  
 $H = \nabla_{\alpha}^2 L_{\text{val}}(w^*, \alpha^*)$  is the Hessian
  - Standard DARTS lead to “Sharp solutions” (large Hessian)

# Differentiable NAS

## Reason 1: sharpness of the solution

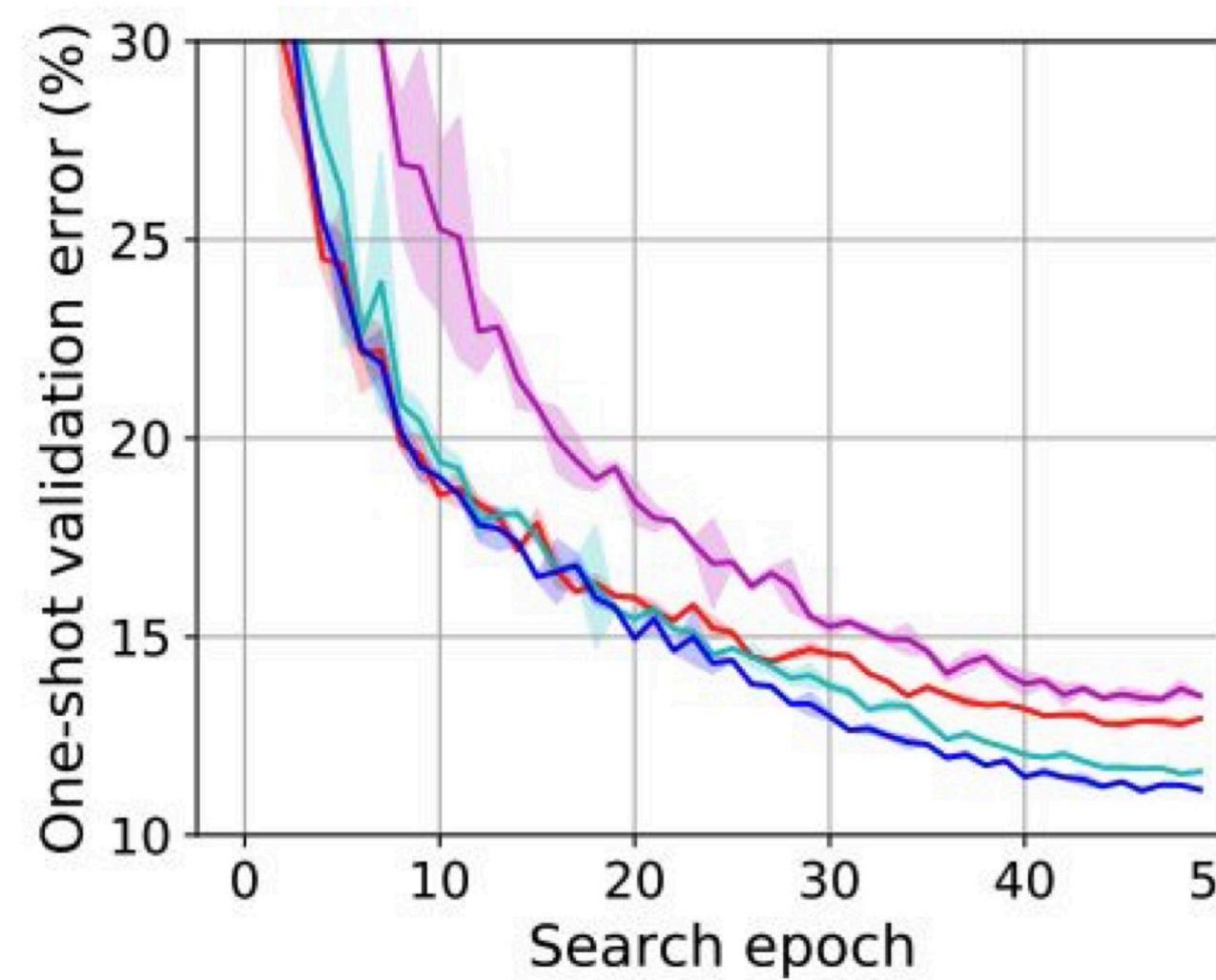
- A good **continuous solution** doesn't imply a good **discrete solution**
- Gap between continuous and discrete solutions can be estimated by sharpness
- Standard DARTS lead to “Sharp solutions” (large Hessian)



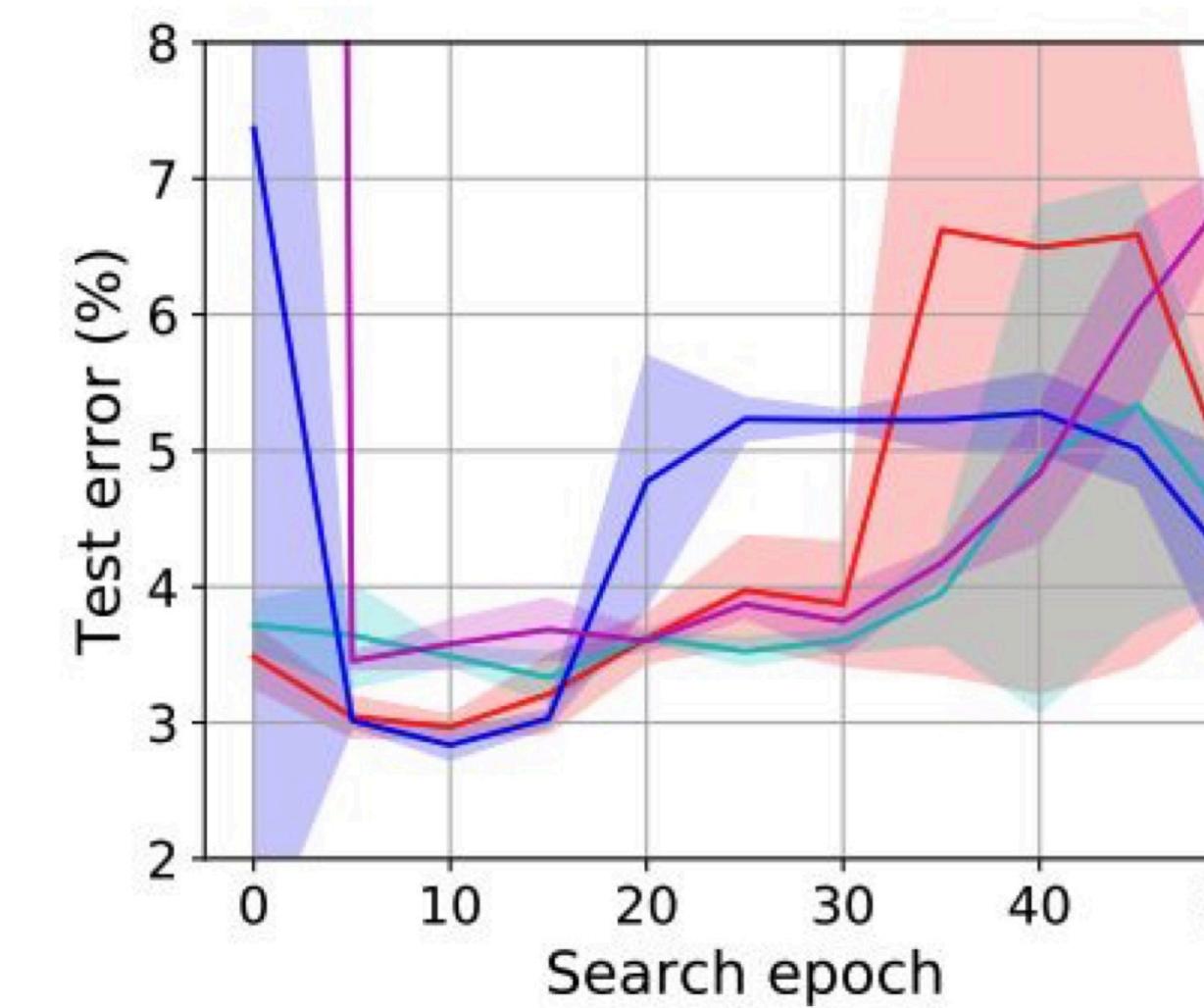
# Differentiable NAS

## Reason 1: sharpness of the solution

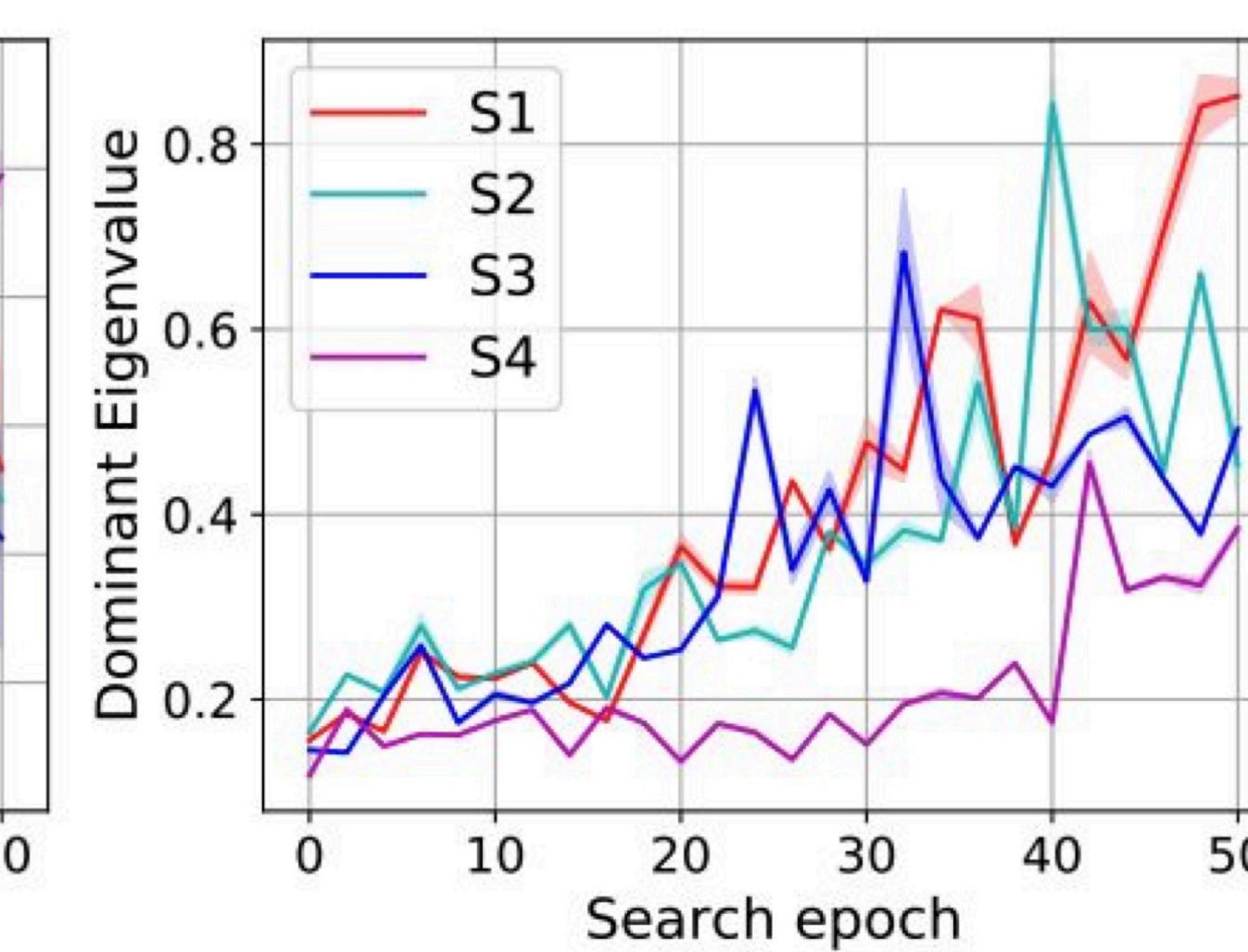
- DARTS training leads to sharp local minimums



Validation error of supernet



Test error of final architecture

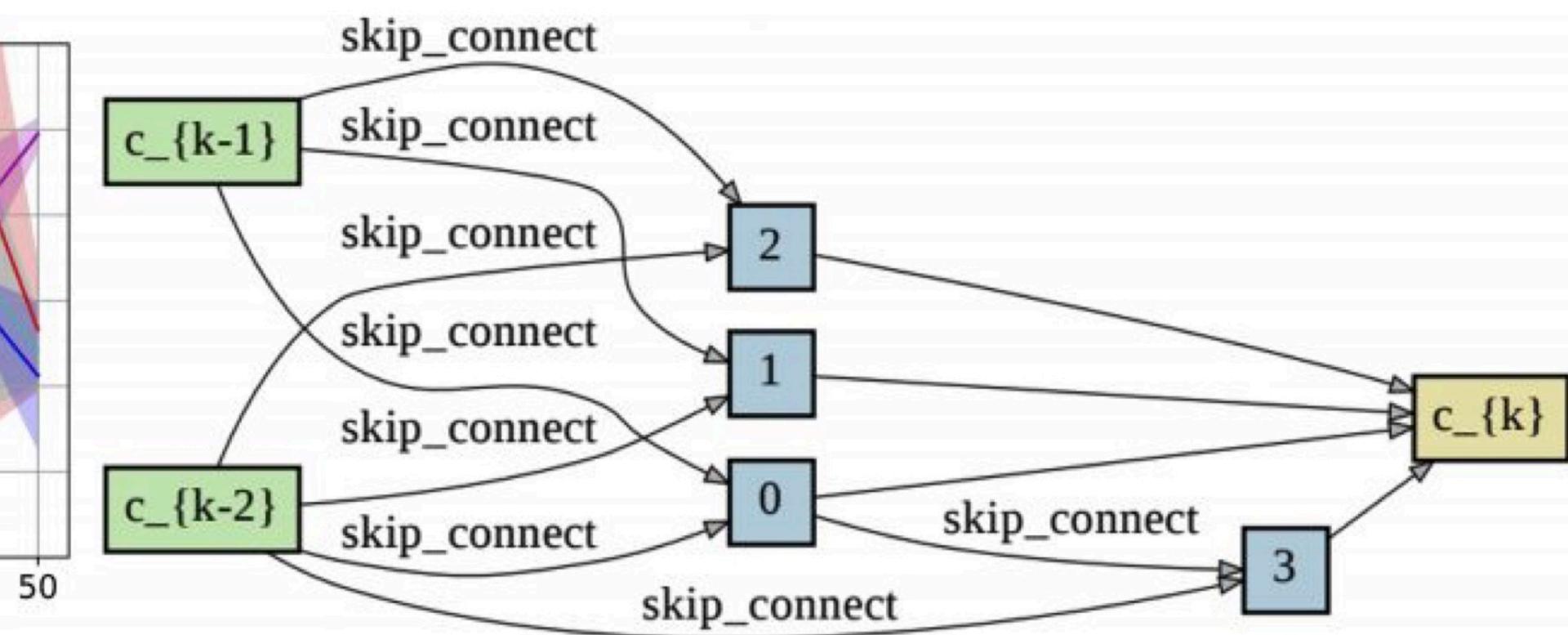
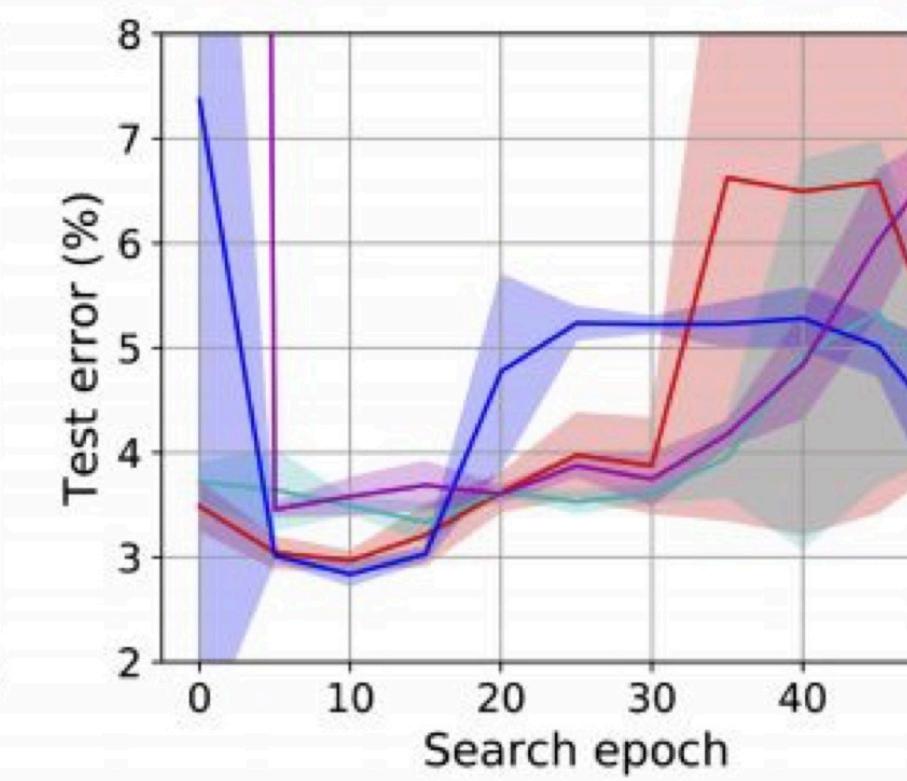
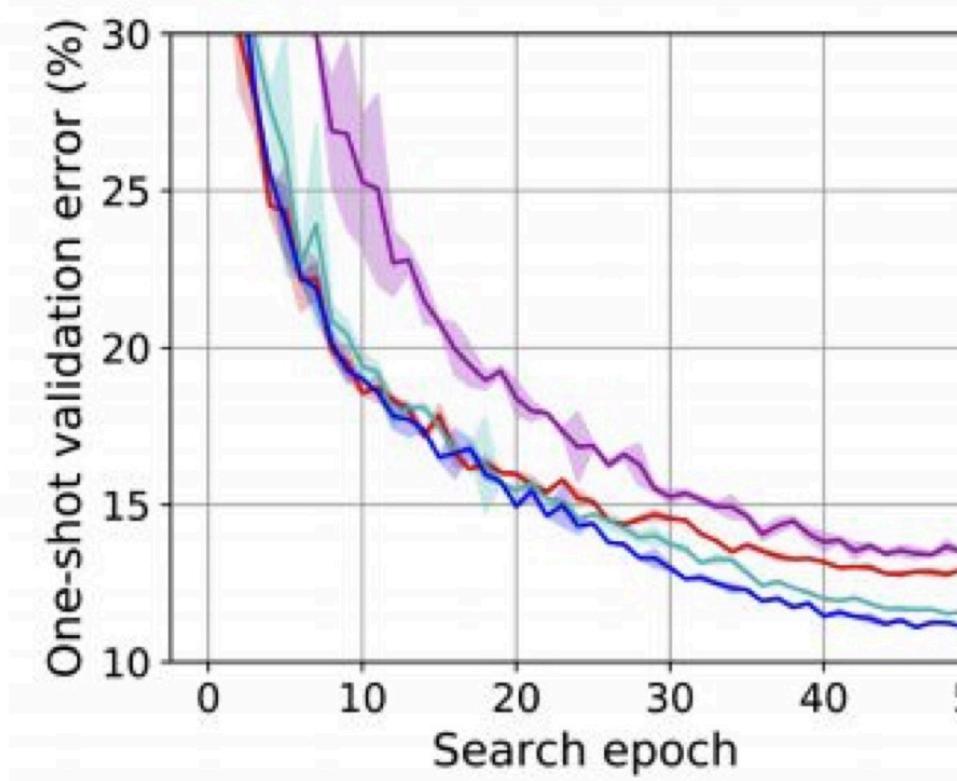


Dominant eigenvalue of Hessian

# Differentiable NAS

## Reason 2: Skip connection domination

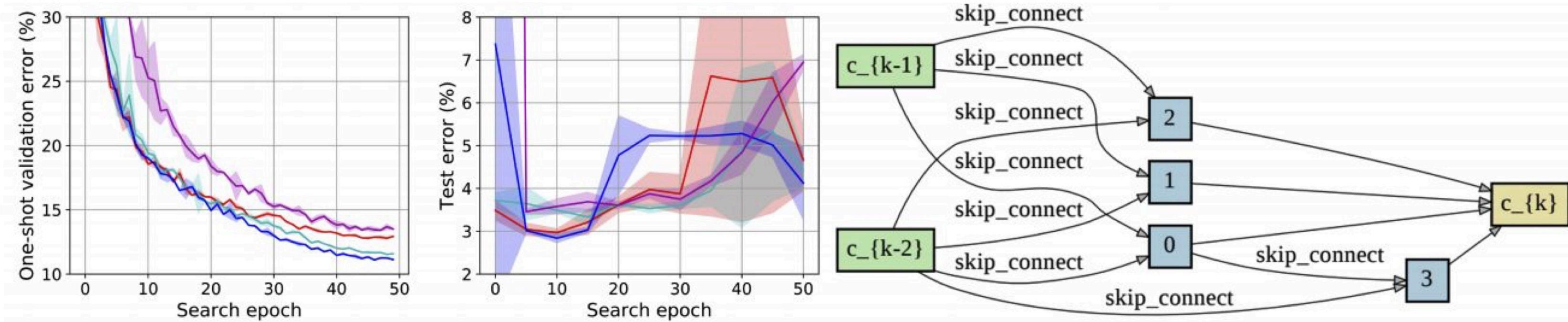
- Supernet accuracy ↑
- Weight for skip connection ↑
- Weight for convolution ↓



# Differentiable NAS

## Reason 2: Skip connection domination

- Supernet accuracy  $\uparrow$
- Weight for skip connection  $\uparrow$
- Weight for convolution  $\downarrow$



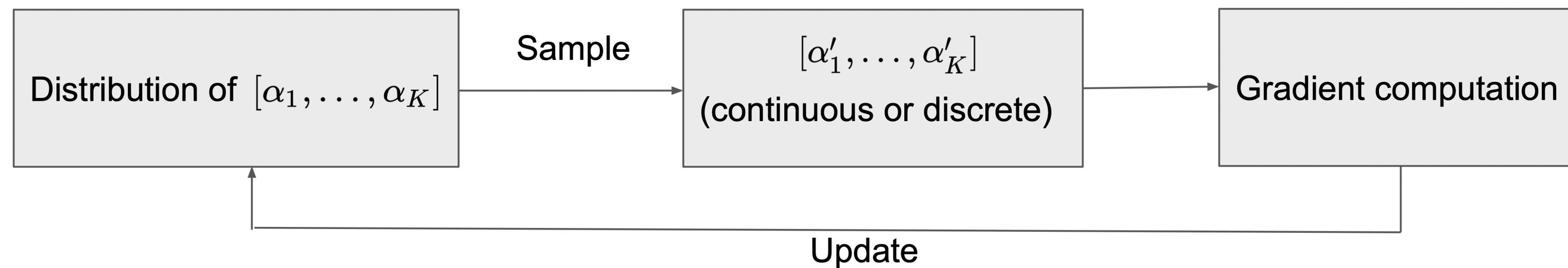
- Formally, we proved that for the optimal supernet, as number of layers goes to infinity,  $\alpha_{\text{skip}} \uparrow 1$  and  $\alpha_{\text{conv}} \downarrow 0$

# Improvements over DARTS

- Supernet Training
  - Usually aim to make superset more “discretizable”
  - Balance exploration and exploitation
- Scalability
  - How to use more blocks in searching?
  - Reduce memory overhead to directly search on larger problems
- Architecture Selection
  - Does architecture weight  $\alpha$  really indicate their performance

# Supernet training: Distribution learning

- Rethink DARTS as a distribution learning problem
  - For each edge,  $[\alpha_1, \dots, \alpha_k]$  defines a distribution over operations
  - We eventually “sample” an architecture from this distribution
  - How to learn  $[\alpha_1, \dots, \alpha_k]$  based on gradient-based optimization?
- Benefits:
  - Performance will be preserved better after discretization
  - Reduced training time in some cases



# Supernet training: Distribution learning

## Gumbel softmax

- Sampling from a distribution  $i \sim \alpha_i / \sum_{i'} \alpha_{i'}$  (can't backprop from  $i$  to  $\alpha$ )
- Gumbel-max: this is equivalent to

$$i = \arg \max_{i'} \{G_{i'} + \log(\alpha_i)\}$$

where each  $G_{i'} \sim \text{Gumbel}(0, 1)$

- Gumbel-softmax: using softmax with temperature annealed to be close to zero

$$z_i = \frac{\exp(G_i + \log(\alpha_i)) / \gamma}{\sum_{i'} \exp(G_{i'} + \log(\alpha_{i'})) / \gamma}$$

- This enables back-propagation to  $[\alpha_1, \dots, \alpha_K]$  (reparameterization trick)
- SNAS: use Gumbel softmax with annealed temperature in DARTS

# Supernet training: Distribution learning

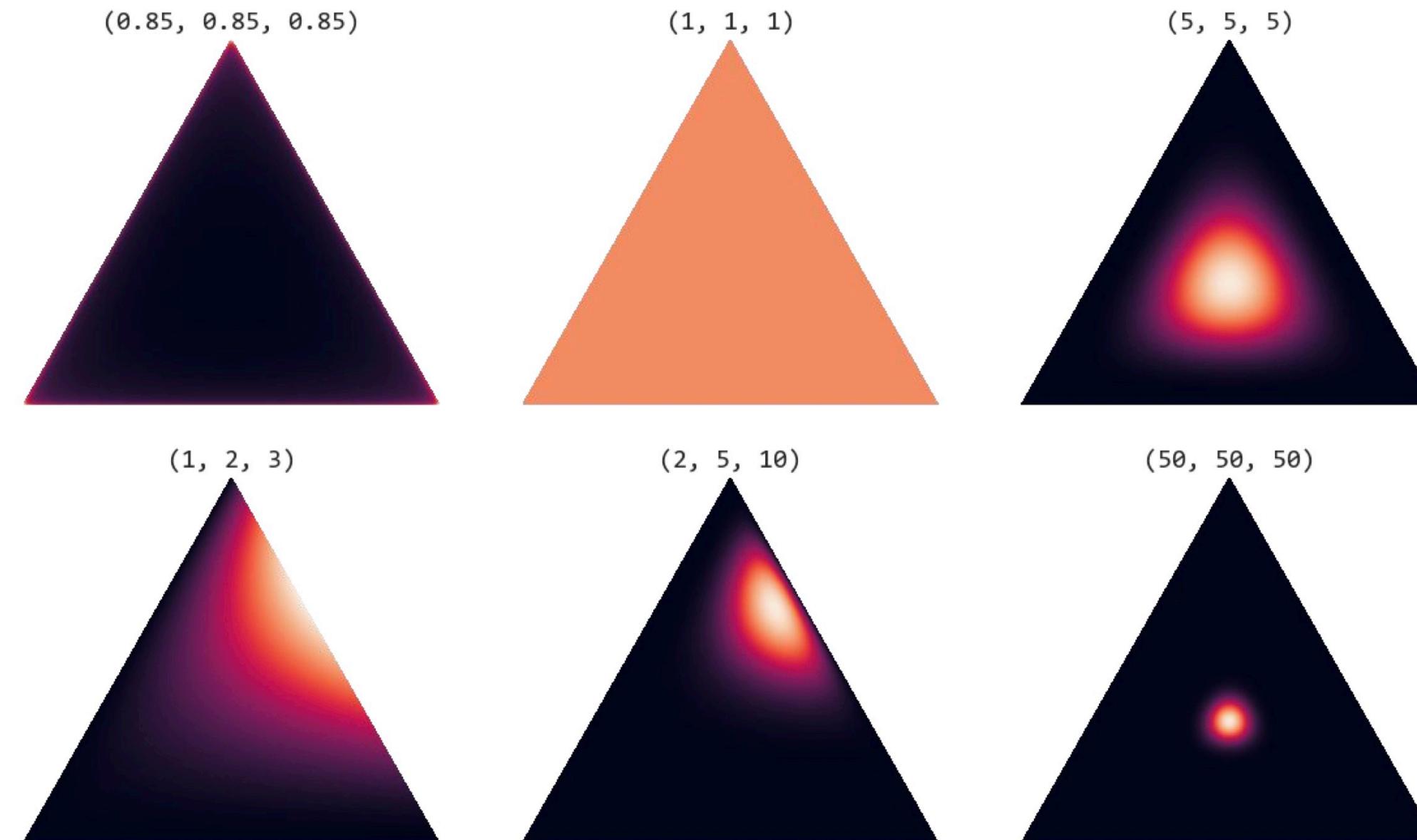
## DrNAS

- Assume architecture parameters  $[\alpha_1, \dots, \alpha_K]$  are sampled from Dirichlet Distribution:

$$[\alpha_1, \dots, \alpha_K] \sim \text{Dir}([\beta_1, \dots, \beta_K])$$

- Dirichlet distribution samples from the standard K-1 simplex

- $\beta \ll 1$  leads to **sparse** samples with high variance
- $\beta \gg 1$  leads to **dense** samples with low variance (for sufficient exploration)



# Supernet training: Distribution learning

## DrNAS

- DrNAS objective:

- Point estimation → distribution learning

$$\min_{\beta} E_{q(\alpha|\beta)}[L_{val}(w^*(\alpha), \alpha)] + \lambda d(\beta, \hat{\beta}), \quad s.t.$$

$$w^* = \arg \min_w L_{train}(w, \alpha), \quad q(\alpha|\beta) \sim Dir(\beta)$$

- Gradient computation:

$$\frac{d\alpha_i}{d\beta_j} = -\frac{\frac{\partial F_{Beta}}{\partial \beta_j}(\alpha_j|\beta_j, \beta_{tot} - \beta_j)}{f_{Beta}(\alpha_j|\beta_j, \beta_{tot} - \beta_j)} \times \left( \frac{\delta_{ij} - \alpha_i}{1 - \alpha_j} \right)$$

- Architecture selection: magnitude of  $\beta$

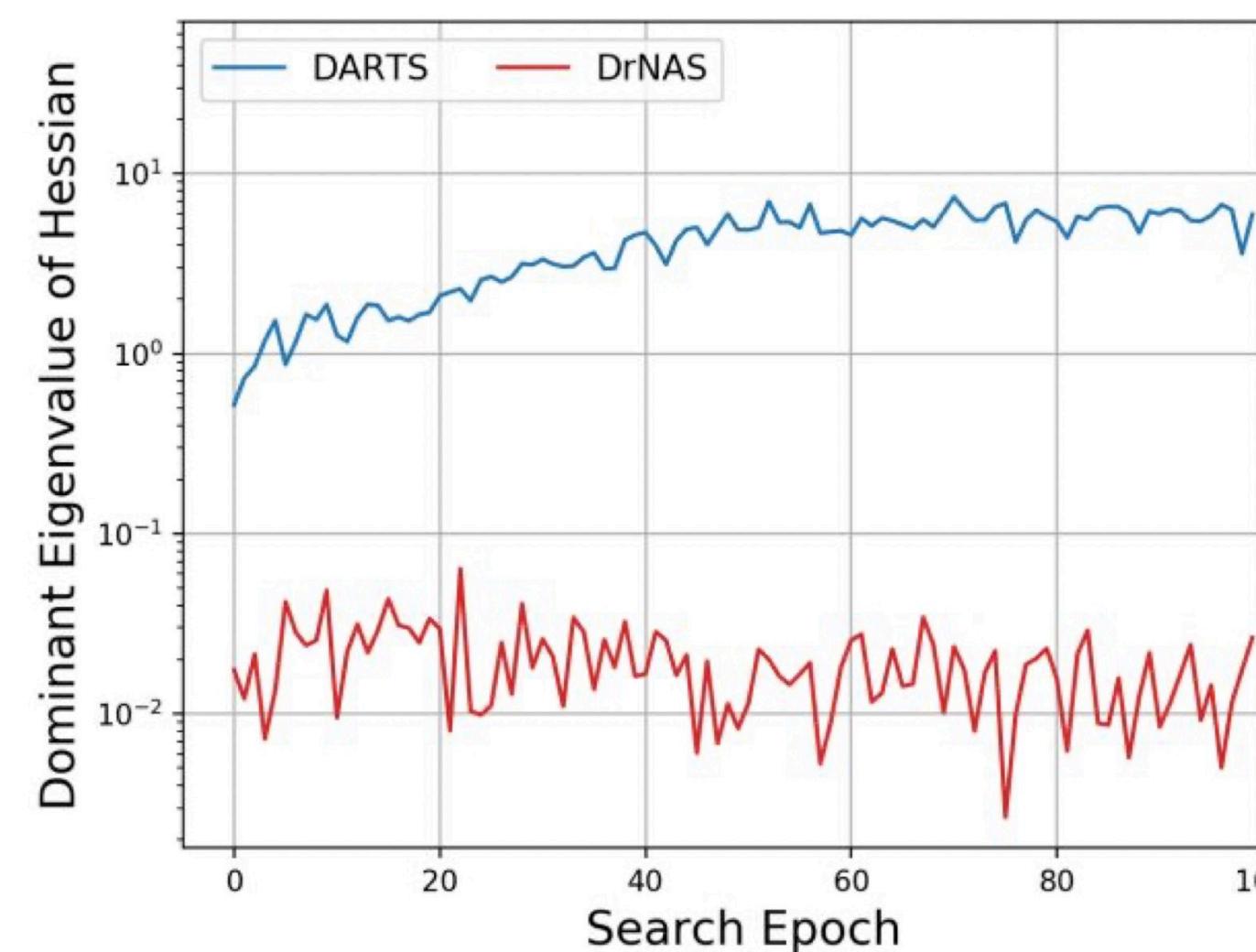
# Supernet training: Distribution learning

## DrNAS

- We show that minimizing the expected Lval controls the trace norm of Hessian:

$$E_{q(\alpha|\beta)}(L_{val}(w, \alpha)) \gtrsim \tilde{L}_{val}(w^*, \mu) + C \cdot \text{tr}(\nabla_\mu^2 \tilde{L}_{val}(w^*, \mu))$$

with  $\tilde{L}_{val}(w^*, \mu) = L_{val}(w^*, \text{Softmax}(\mu))$



# Supernet training: Distribution learning

## DrNAS

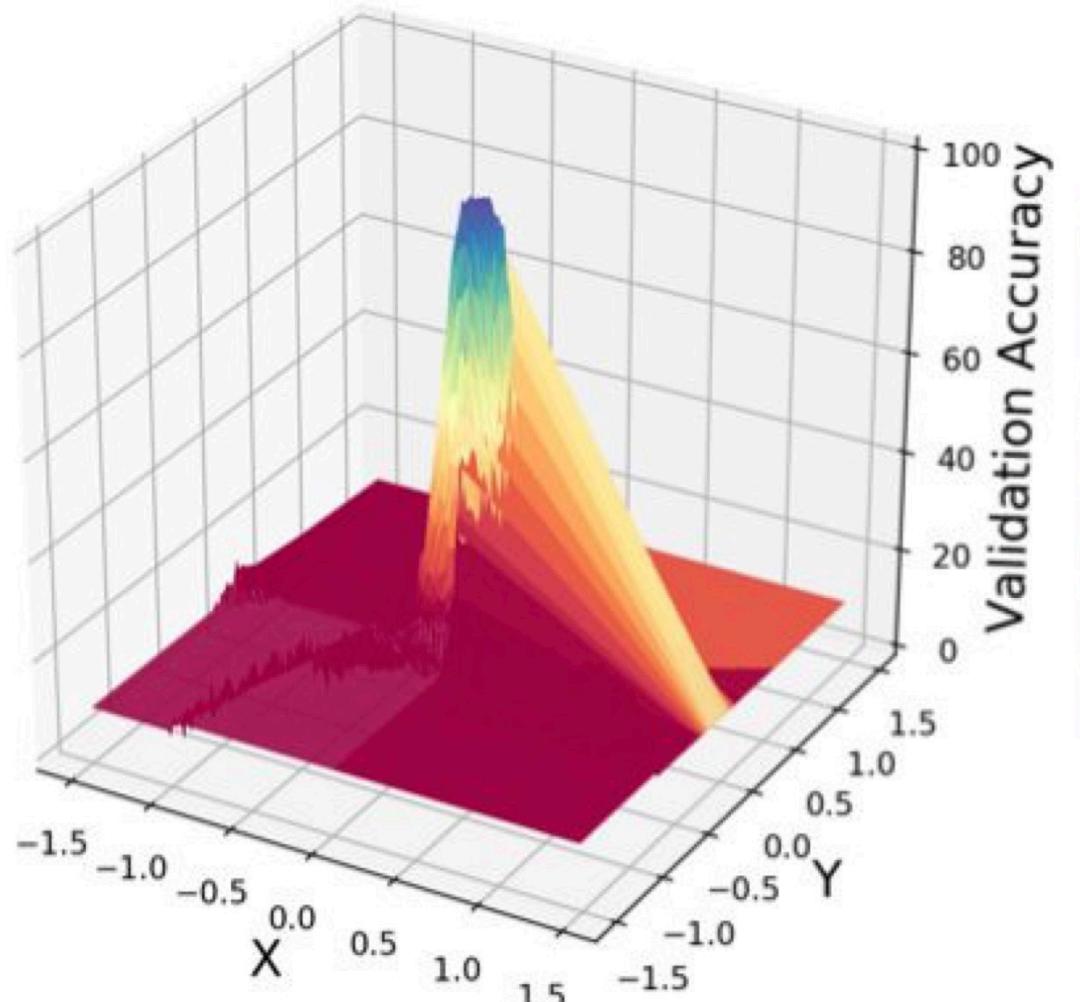
- On NAS-Bench-201
  - Achieve **oracle** when searching on CIFAR-100  
DrNAS (73.51) vs SNAS (69.34) vs DARTS (38.97)

Method	CIFAR-10		CIFAR-100		ImageNet-16-120	
	validation	test	validation	test	validation	test
ResNet	90.83	93.97	70.42	70.86	44.53	43.63
Random (baseline)	$90.93 \pm 0.36$	$93.70 \pm 0.36$	$70.60 \pm 1.37$	$70.65 \pm 1.38$	$42.92 \pm 2.00$	$42.96 \pm 2.15$
RSPS	$84.16 \pm 1.69$	$87.66 \pm 1.69$	$45.78 \pm 6.33$	$46.60 \pm 6.57$	$31.09 \pm 5.65$	$30.78 \pm 6.12$
Reinforce	$91.09 \pm 0.37$	$93.85 \pm 0.37$	$70.05 \pm 1.67$	$70.17 \pm 1.61$	$43.04 \pm 2.18$	$43.16 \pm 2.28$
ENAS	$39.77 \pm 0.00$	$54.30 \pm 0.00$	$10.23 \pm 0.12$	$10.62 \pm 0.27$	$16.43 \pm 0.00$	$16.32 \pm 0.00$
DARTS (1st)	$39.77 \pm 0.00$	$54.30 \pm 0.00$	$38.57 \pm 0.00$	$38.97 \pm 0.00$	$18.87 \pm 0.00$	$18.41 \pm 0.00$
DARTS (2nd)	$39.77 \pm 0.00$	$54.30 \pm 0.00$	$38.57 \pm 0.00$	$38.97 \pm 0.00$	$18.87 \pm 0.00$	$18.41 \pm 0.00$
GDAS	$90.01 \pm 0.46$	$93.23 \pm 0.23$	$24.05 \pm 8.12$	$24.20 \pm 8.08$	$40.66 \pm 0.00$	$41.02 \pm 0.00$
SNAS	$90.10 \pm 1.04$	$92.77 \pm 0.83$	$69.69 \pm 2.39$	$69.34 \pm 1.98$	$42.84 \pm 1.79$	$43.16 \pm 2.64$
DSNAS	$89.66 \pm 0.29$	$93.08 \pm 0.13$	$30.87 \pm 16.40$	$31.01 \pm 16.38$	$40.61 \pm 0.09$	$41.07 \pm 0.09$
PC-DARTS	$89.96 \pm 0.15$	$93.41 \pm 0.30$	$67.12 \pm 0.39$	$67.48 \pm 0.89$	$40.83 \pm 0.08$	$41.31 \pm 0.22$
<b>DrNAS</b>	$91.55 \pm 0.00$	$94.36 \pm 0.00$	$73.49 \pm 0.00$	$73.51 \pm 0.00$	$46.37 \pm 0.00$	$46.34 \pm 0.00$
<b>optimal</b>	91.61	94.37	73.49	73.51	46.77	47.31

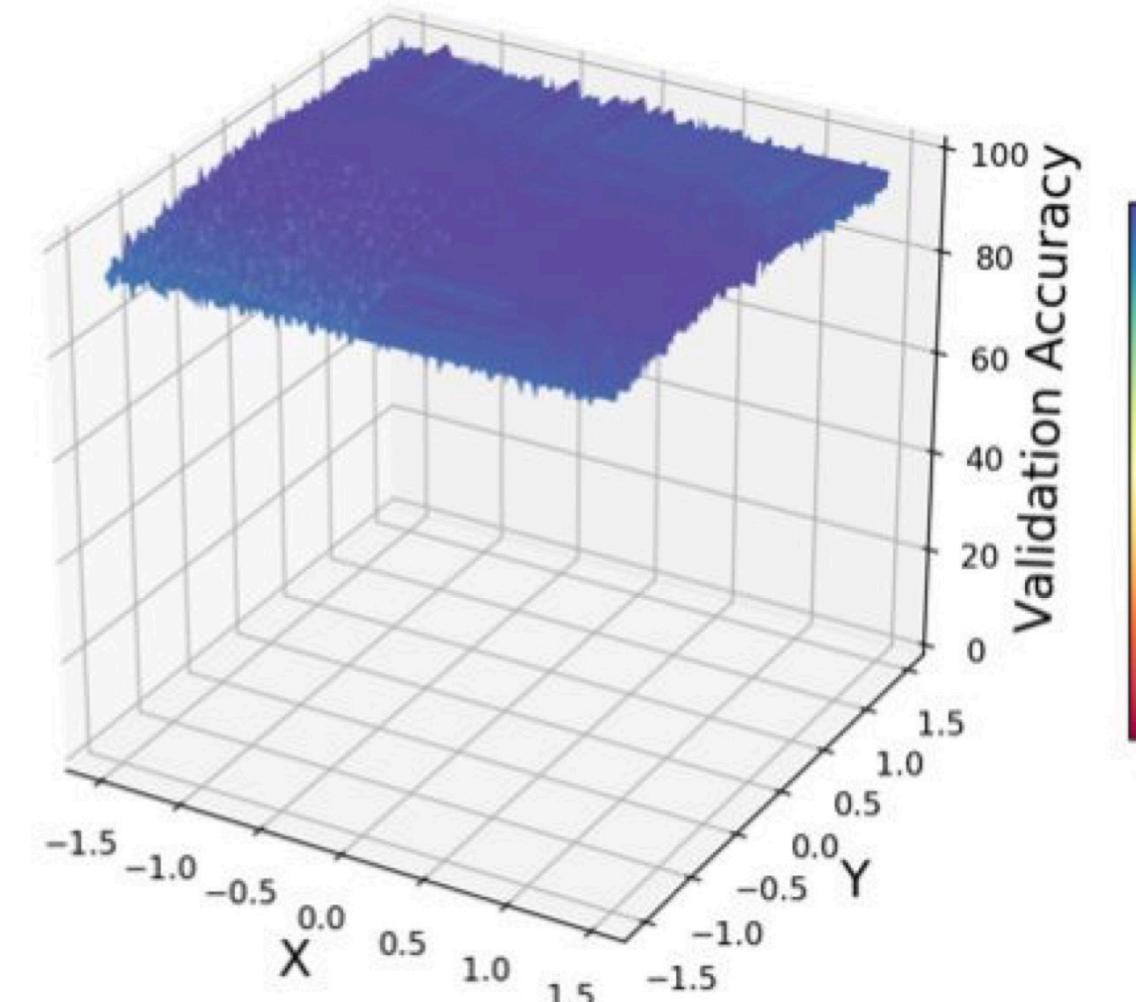
# Supernet training

## Perturbation-based regularization

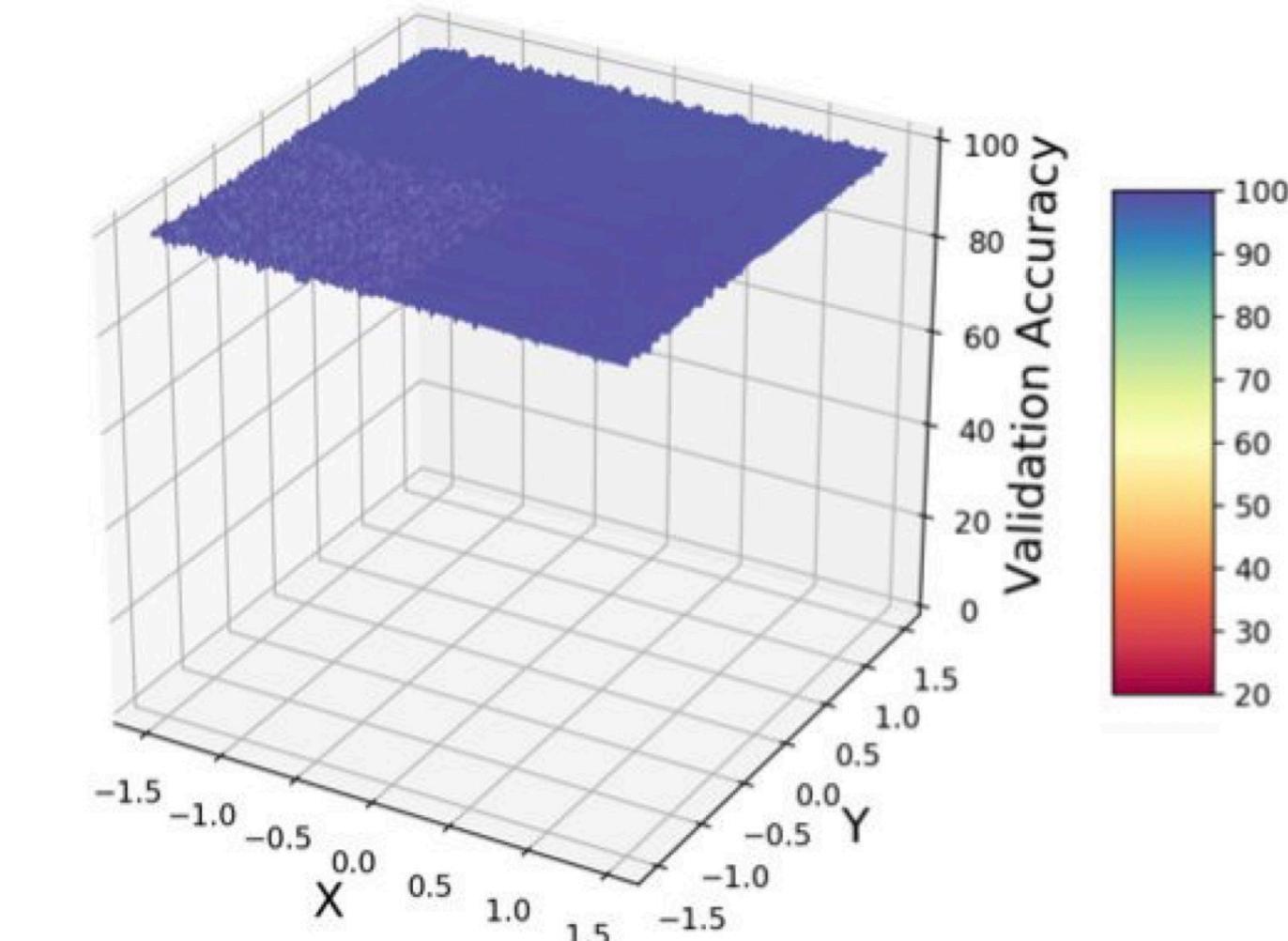
- A smoother landscape will make supernet robust to discretization



(a) DARTS



(b) SDARTS-RS



(c) SDARTS-ADV

# Supernet training

## Perturbation-based regularization

- Make supernet robust to  $\alpha$  perturbation
  - Since we need to perturb it to a discrete architecture in the final stage
  - Mathematically, we hope the superset robust to random or adversarial (worst-case) perturbation of  $\alpha$

# Supernet training

## Perturbation-based regularization

- Make supernet robust to  $\alpha$  perturbation
  - Since we need to perturb it to a discrete architecture in the final stage
  - Mathematically, we hope the superset robust to random or adversarial (worst-case) perturbation of  $\alpha$

$$\min_{\alpha} L_{\text{val}}(\bar{w}(\alpha), \alpha), \text{ s.t.}$$

SDARTS-RS:  $\bar{w}(\alpha) = \arg \min_w E_{\delta \sim U[-\epsilon, \epsilon]} L_{\text{train}}(w, A + \delta)$

SDARTS-Adv:  $\bar{w}(\alpha) = \arg \min_w \max_{\|\delta\| \leq \epsilon} L_{\text{train}}(w, A + \delta)$

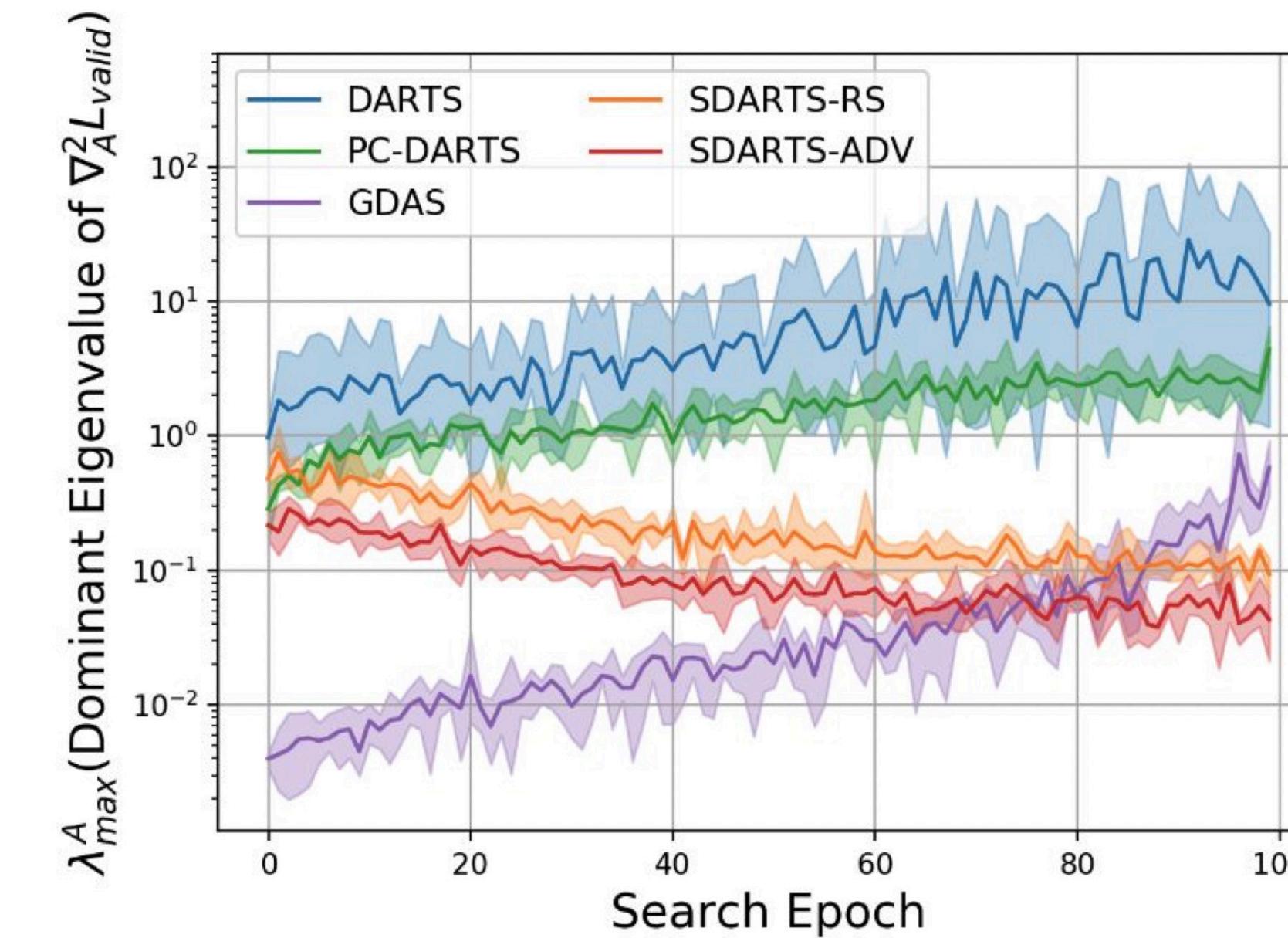
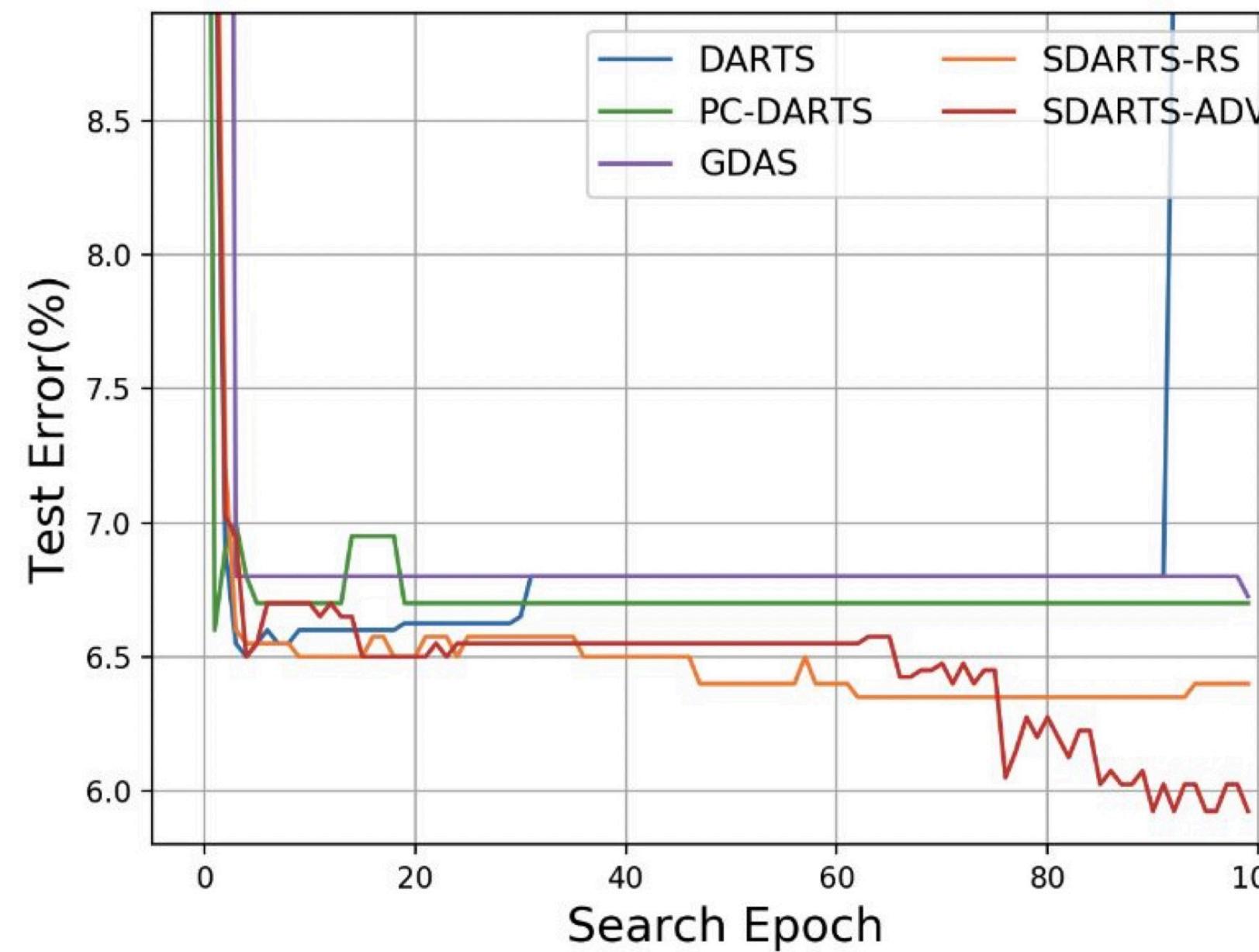
SDARTS: Each step

- Perturb  $\alpha$ 
  - ◆ Random:  $\alpha' \leftarrow \alpha + N(0, \sigma^2)$
  - ◆ Adversarial:  
 $\alpha' \leftarrow \alpha + \nabla_{\alpha} L_{\text{train}}(\alpha, w)$
- Update  $w$  based on  $\alpha'$
- Update  $\alpha$  based on  $w$

# Supernet training

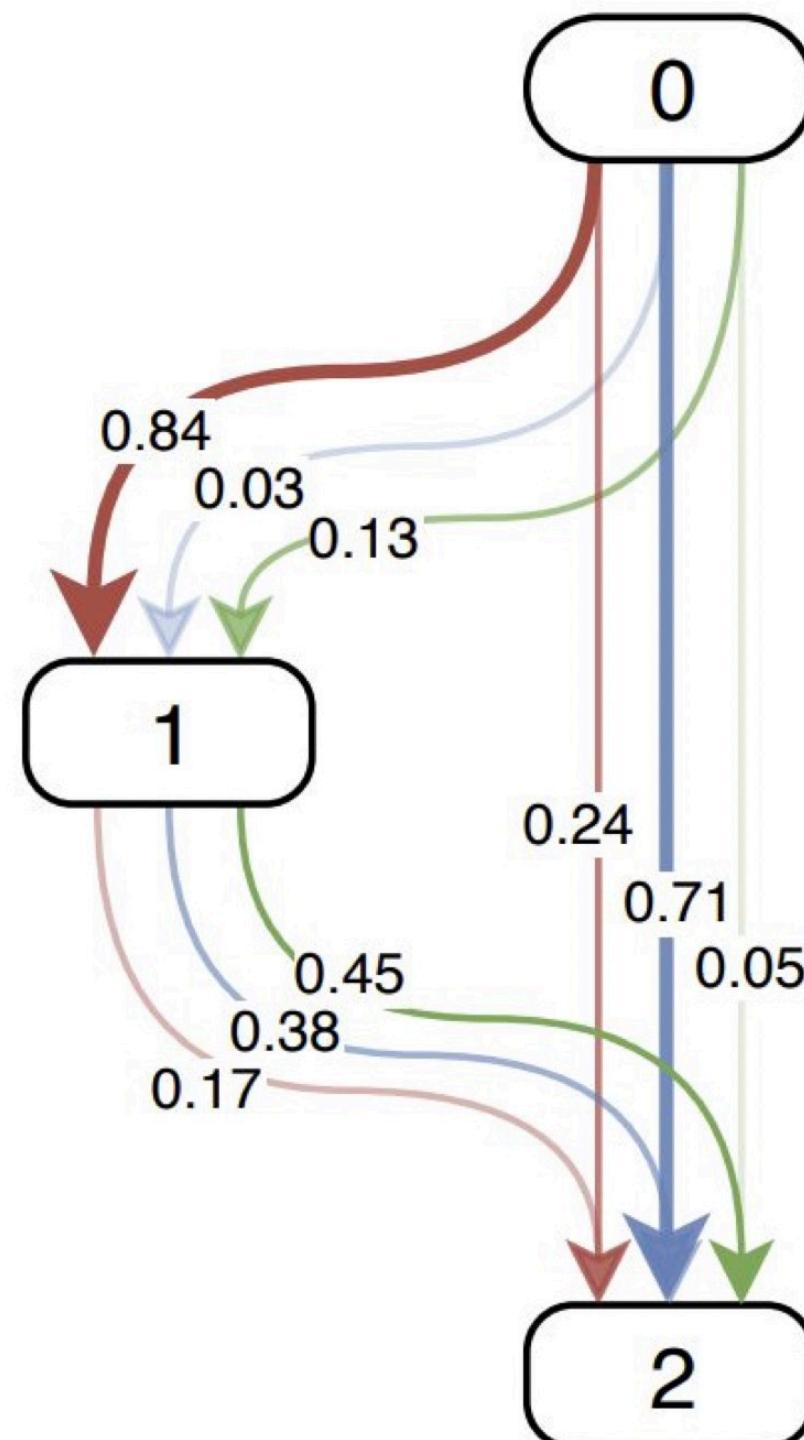
## SmoothDARTS

- On NAS-Bench-1Shot1
  - Continues to discover better architectures
  - Anneal Hessian to a low level

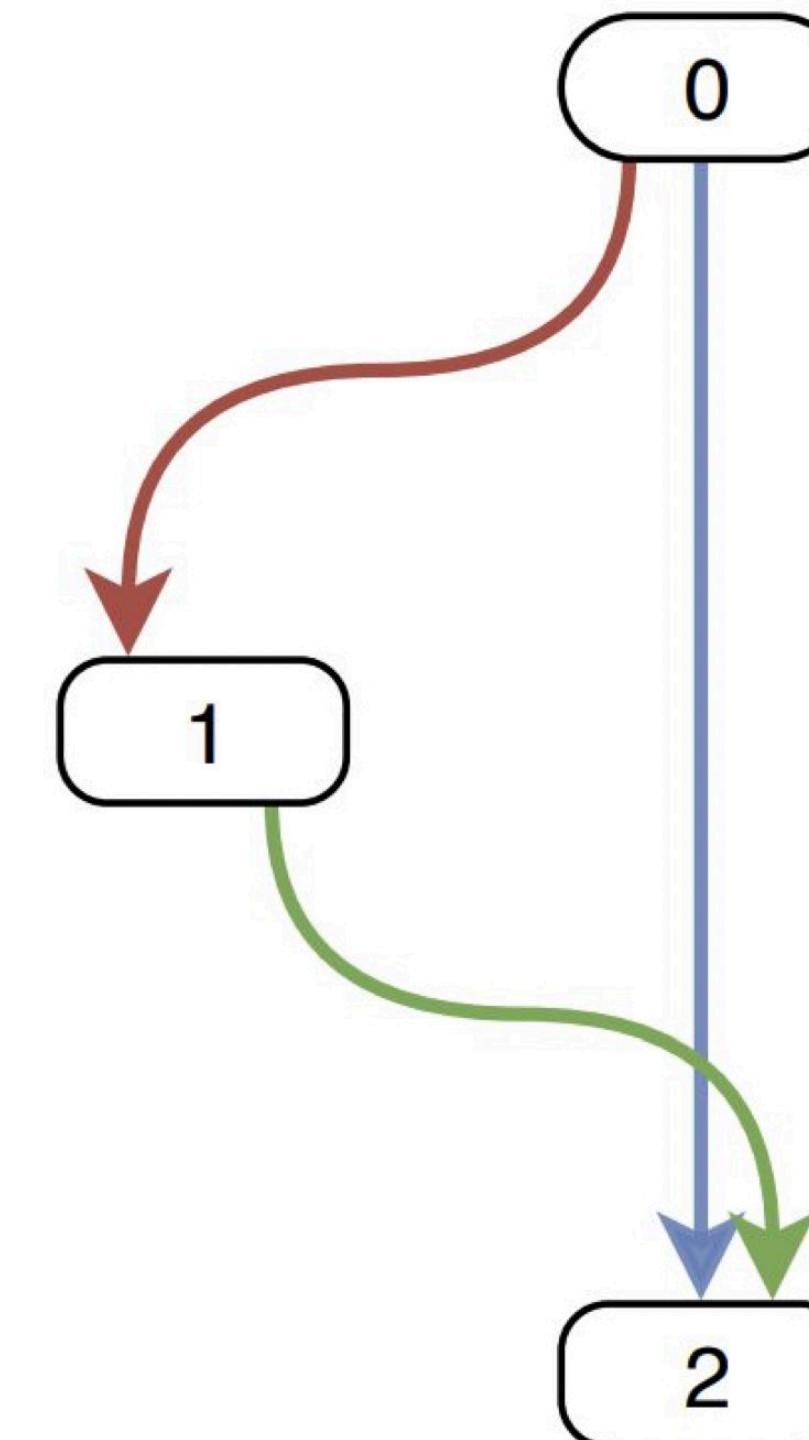


# Architecture Selection

## Architecture Selection in DARTS



(e) Search end



(f) Final cell

# Architecture Selection

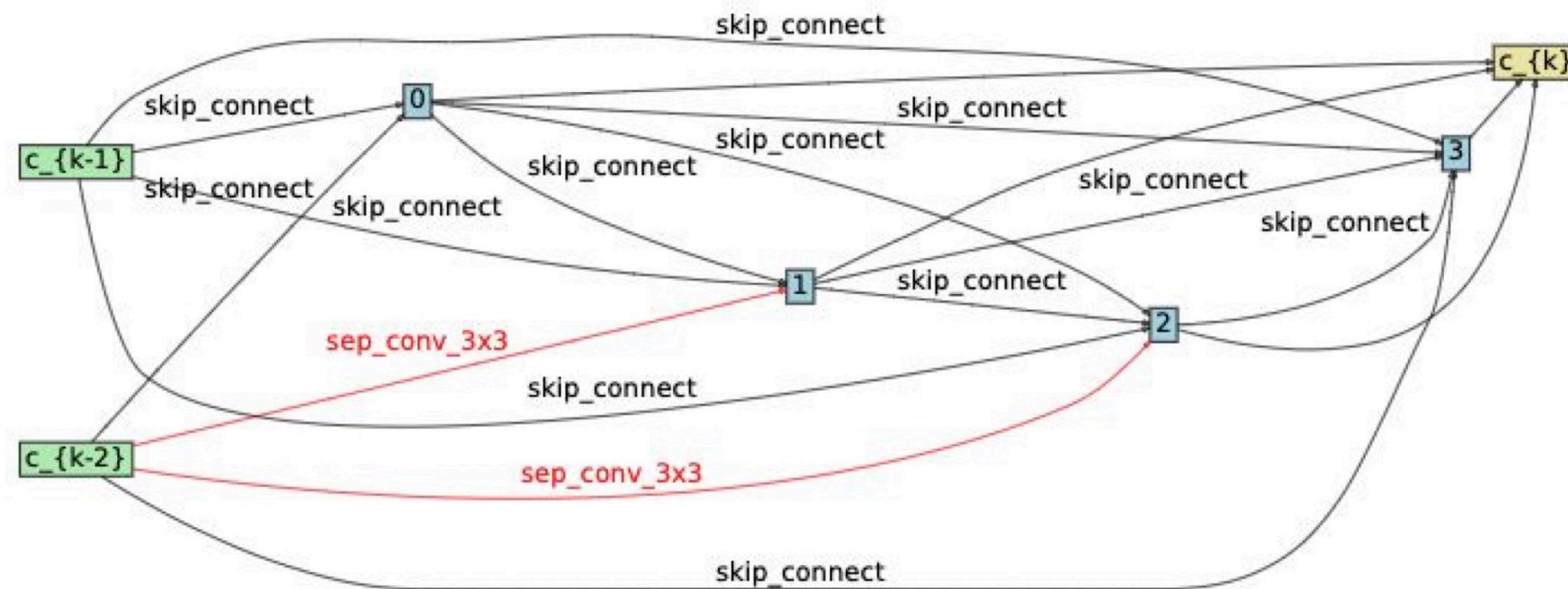
## Architecture Selection in DARTS

- Recall the skip-domination problem:
  - For the optimal supernet with infinite number of layers:  $\alpha_{\text{skip}} \uparrow 1$  and  $\alpha_{\text{conv}} \downarrow 0$
  - $a$  values may not really represent the “**importance**” of each operation
- Skip connection stands out if we select the best operation based on  $a$
- Does  $\alpha_{\text{skip}} > \alpha_{\text{conv}}$  mean skip connection is better than convolution?

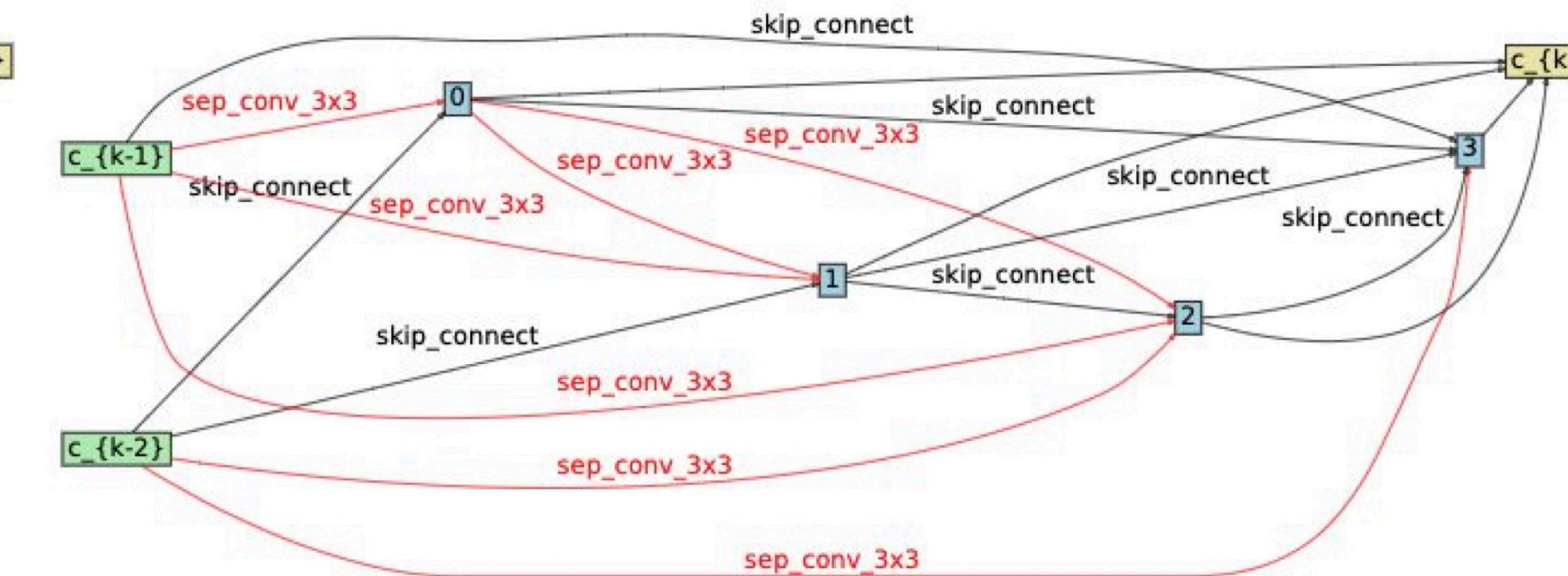
# Architecture Selection

## Does $\alpha$ represent operation strength?

- Probably **Not!**
- S2: (Skip\_connect, sep\_conv\_3x3)
  - Skip connections dominate according to  $\alpha$
  - But the accuracy of S2 supernet benefits from more convolutions



Magnitude-base selection



Progressive tuning selection

# Architecture Selection

Does  $\alpha$  represent operation strength?

- Same observations on large space: DARTS space
  - Magnitude of  $\alpha$  deviates from accuracy of the supernet
  - Some operations with small  $\alpha$  are in fact more important for supernet

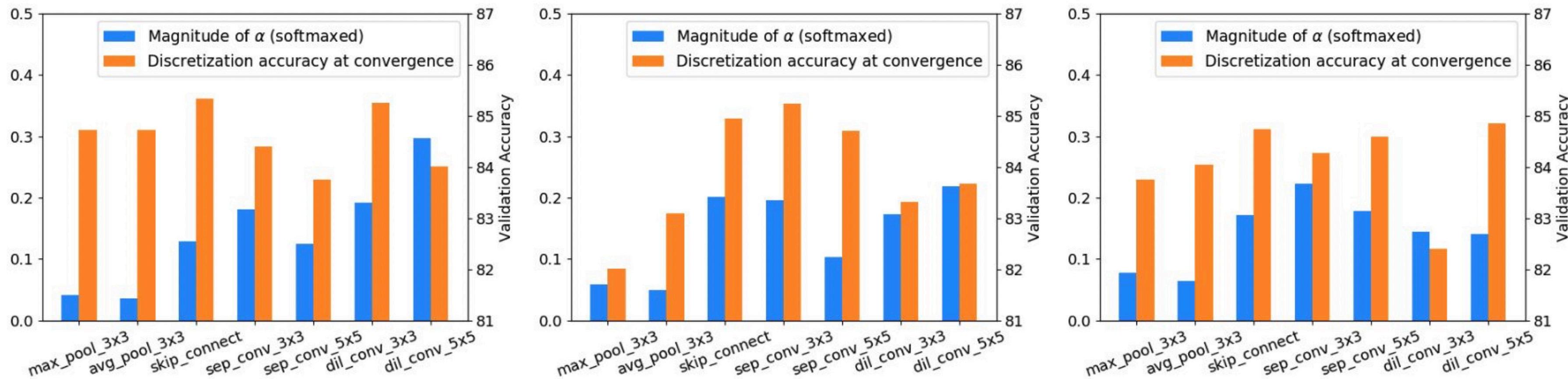


Figure: Magnitude of  $\alpha$  vs Accuracy after choosing one operation

# Architecture Selection

## A New architecture Selection Method

- Evaluate the importance of an operation  $\circ$  by:
  - Compute the drop of validation accuracy when  $\circ$  is removed (no need for further training)
- Use this to choose the best  $\circ$  for an edge
- Fine-tune the solution, and move to the next edge
- “Perturbation-based selection” (PT for short)

# Architecture Selection

## A New architecture Selection Method

- PT consistently improves over the original magnitude-based selection

<b>Dataset</b>	<b>Space</b>	<b>DARTS</b>	<b>DARTS+PT (Ours)</b>
C10	S1	3.84	3.50
	S2	4.85	2.79
	S3	3.34	2.49
	S4	7.20	2.64
C100	S1	29.46	24.48
	S2	26.05	23.16
	S3	28.90	22.03
	S4	22.85	20.80
SVHN	S1	4.58	2.62
	S2	3.53	2.53
	S3	3.41	2.42
	S4	3.05	2.42

# Architecture Selection

## A New architecture Selection Method

- Performance improves with more searching epochs

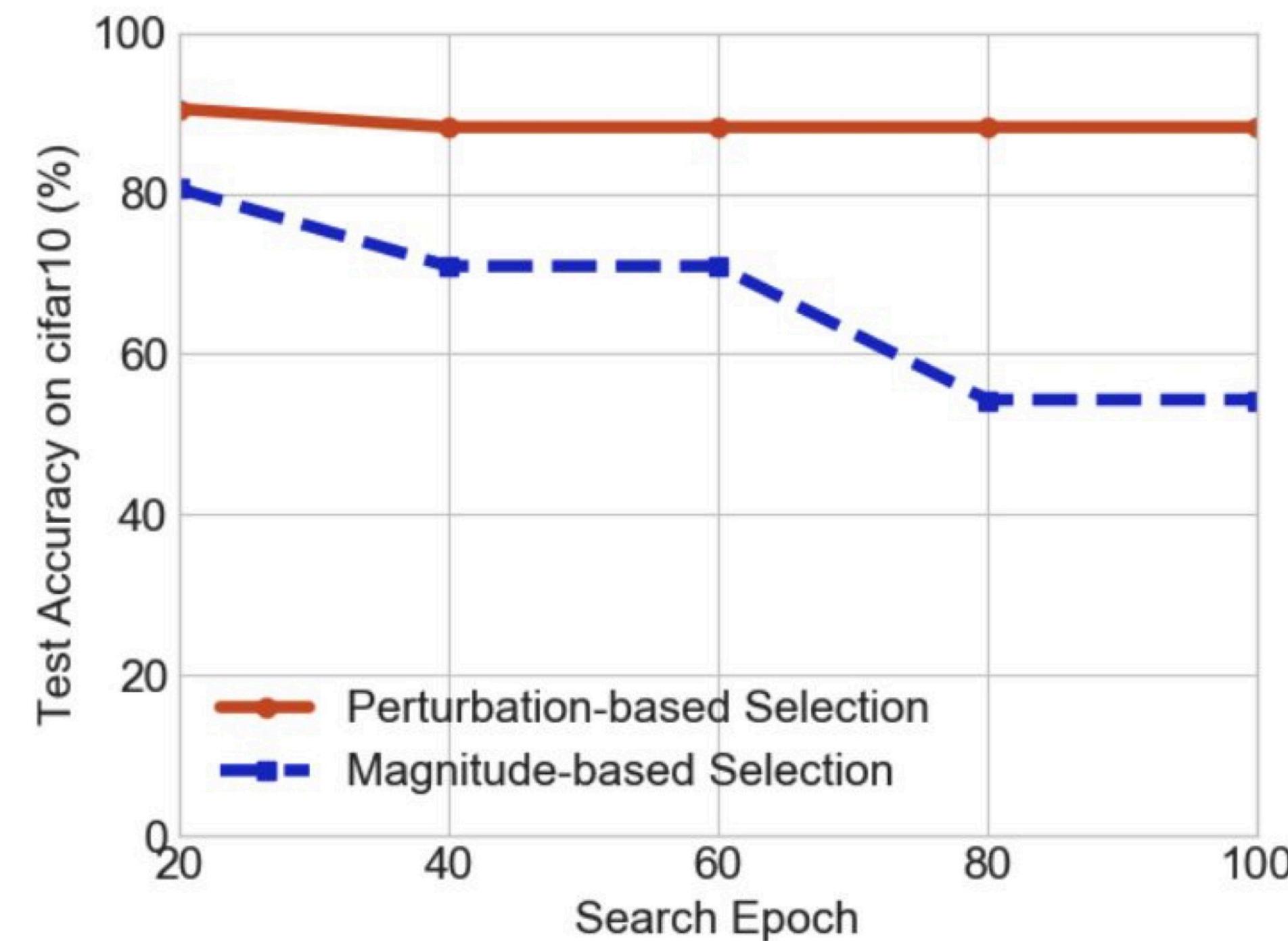


Figure: Test accuracy vs search epoch on NAS-Bench-201 space

# Architecture Selection

## A New architecture Selection Method

Architecture	Test Error (%)	Search Cost (GPU days)	Search Method
DARTS (1st) (Liu et al., 2019)	$3.00 \pm 0.14$	0.4	differentiable
DARTS (2nd) (Liu et al., 2019)	$2.76 \pm 0.09$	1.0	differentiable
SNAS (moderate) (Xie et al., 2019)	$2.85 \pm 0.02$	1.5	differentiable
DrNAS (Chen et al., 2020)	$2.54 \pm 0.03$	0.4	differentiable
NASP (Yao et al., 2019)	$2.83 \pm 0.09$	0.1	differentiable
SDARTS-ADV (Chen & Hsieh, 2020)	$2.61 \pm 0.02$	1.3	differentiable
ProxylessNAS (Cai et al., 2019) <sup>†</sup>	2.08	4.0	differentiable
PC-DARTS (Xu et al., 2020)	$2.57 \pm 0.07$	0.1	differentiable
DrNAS (with progressive learning)	$2.46 \pm 0.03$	0.6	differentiable
DARTS+PT (Wang et al., 2020)	$2.61 \pm 0.08$	0.8	differentiable
SDARTS-ADV+PT	$2.54 \pm 0.01$	0.8	differentiable

<sup>†</sup> Obtained on a different space with PyramidNet as the backbone.

# Architecture Selection

## A New architecture Selection Method

Table 3: Darts+PT on S1-S4 (test error (%)).

Dataset	Space	DARTS	Darts+PT (Ours)	Darts+PT (fix $\alpha$ )*
C10	S1	3.84	3.50	2.86
	S2	4.85	2.79	2.59
	S3	3.34	2.49	2.52
	S4	7.20	2.64	2.58
C100	S1	29.46	24.48	24.40
	S2	26.05	23.16	23.30
	S3	28.90	22.03	21.94
	S4	22.85	20.80	20.66
SVHN	S1	4.58	2.62	2.39
	S2	3.53	2.53	2.32
	S3	3.41	2.42	2.32
	S4	3.05	2.42	2.39

