

# Firefly Neural Architecture Descent: a General Approach for Growing Neural Networks

Lemeng Wu, Bo Liu, Peter Stone, Qiang Liu  
University of Texas at Austin  
(NeurIPS 2020)

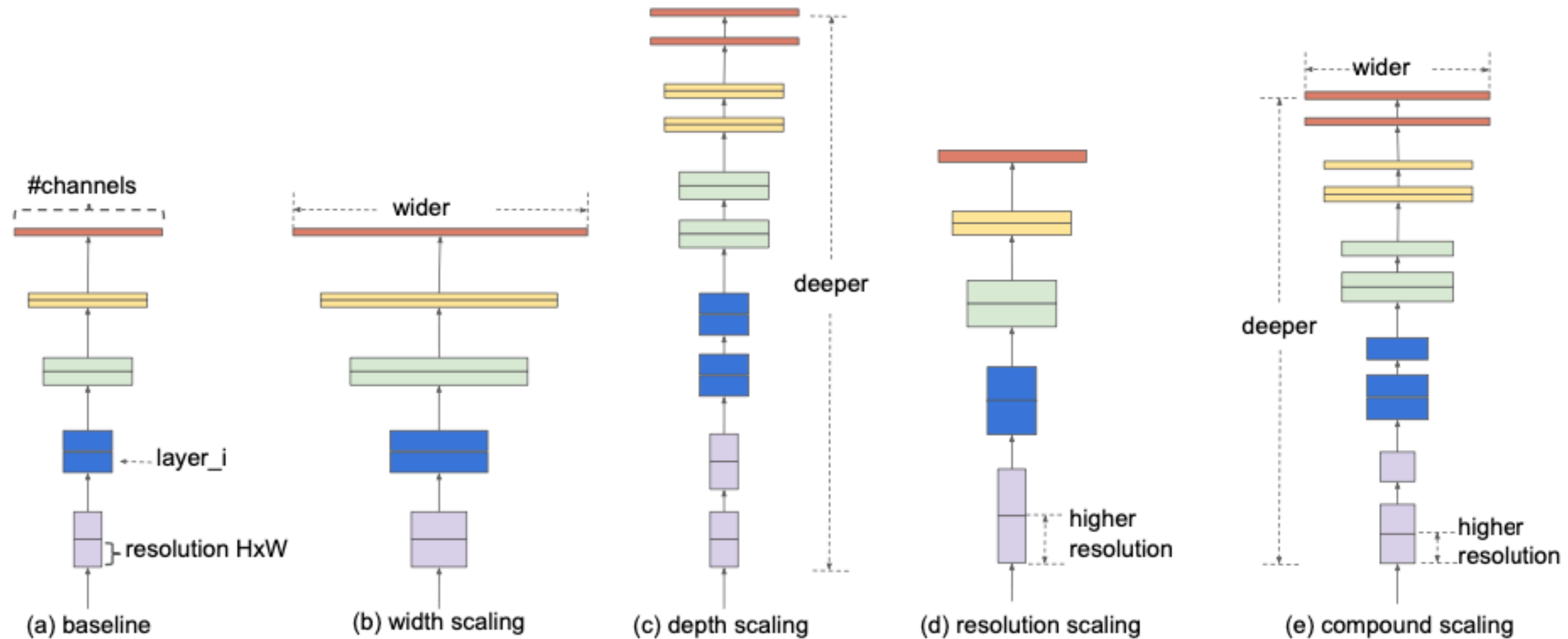
Reporter: Yinghuan Zhang, Mar 9th 2021

# Background

## Model Scaling:

ResNet can be scaled down (e.g., ResNet-18) or up (e.g., ResNet-200)

Scale a ConvNet for different resource constraints



# Background

## Model Scaling:

ResNet can be scaled down (e.g., ResNet-18) or up (e.g., ResNet-200)

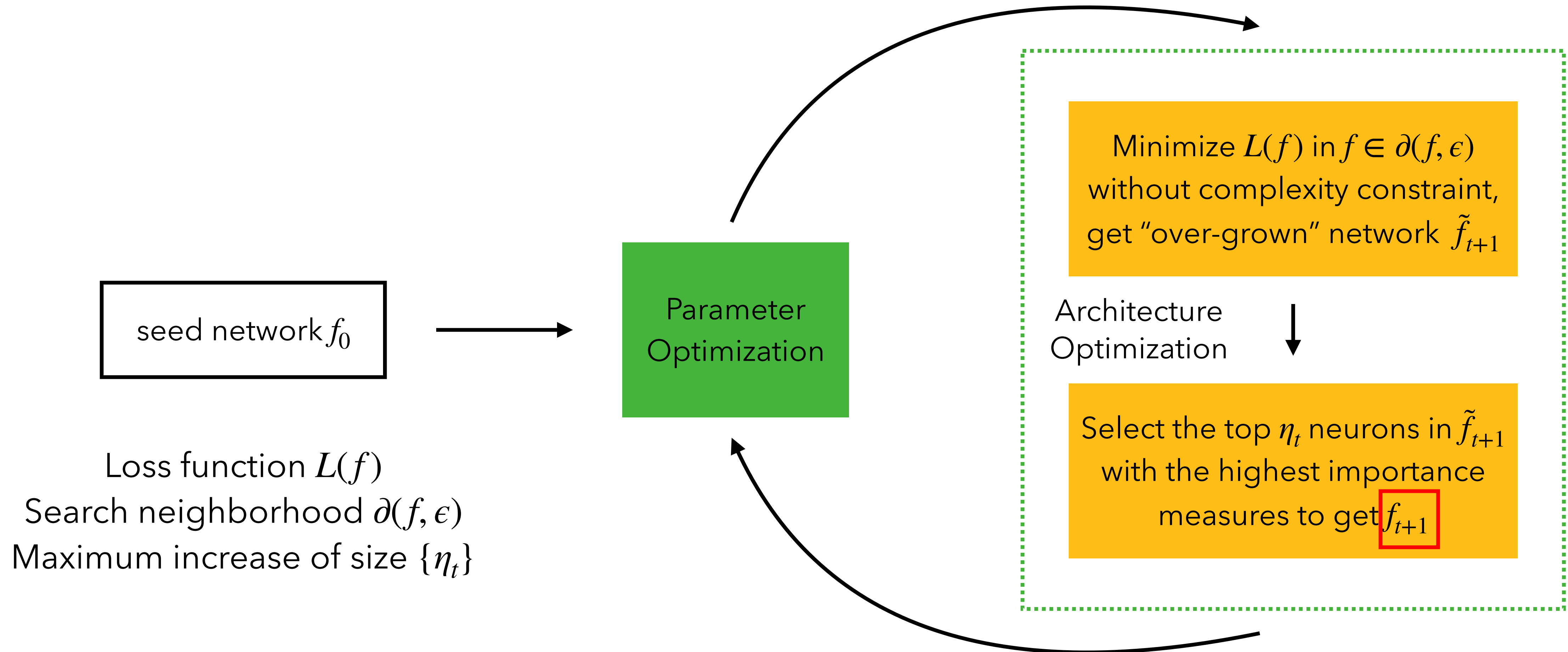
Scale a ConvNet for different resource constraints

## Challenges:

- Grow strategy
- Accuracy drop brought by newly initialized parameters in grown network

***Firefly Neural Architecture Descent*** — a general framework for progressively and dynamically growing neural networks to jointly optimize the networks' parameters and architectures

# Framework Overview



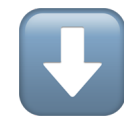
# Framework Overview

## Neighborhood:

$$f(x) = f_t(x) + O(\epsilon), \text{ for } \forall x \in \partial(f_t, \epsilon)$$

## Goal:

$$\arg \min_f \{L(f) \text{ s.t. } f \in \Omega, C(f) \leq \eta\}$$



$$f_{t+1} = \arg \min_f \{L(f) \text{ s.t. } f \in \partial(f_t, \epsilon), C(f) \leq C(f_t) + \eta_t\}$$

Minimize  $L(f)$  in  $f \in \partial(f, \epsilon)$   
without complexity constraint,  
get "over-grown" network  $\tilde{f}_{t+1}$

Architecture  
Optimization



Select the top  $\eta_t$  neurons in  $\tilde{f}_{t+1}$   
with the highest importance  
measures to get  $f_{t+1}$

## Optimization Goal:

$$f_{t+1} = \arg \min_f \{L(f) \text{ s.t. } f \in \partial(f_t, \epsilon), C(f) \leq C(f_t) + \eta_t\}$$

## Goal in Architecture Optimization:

- Instantiate neighborhood  $\partial(f_t, \epsilon)$
- Solve the optimization in  $f_{t+1}$

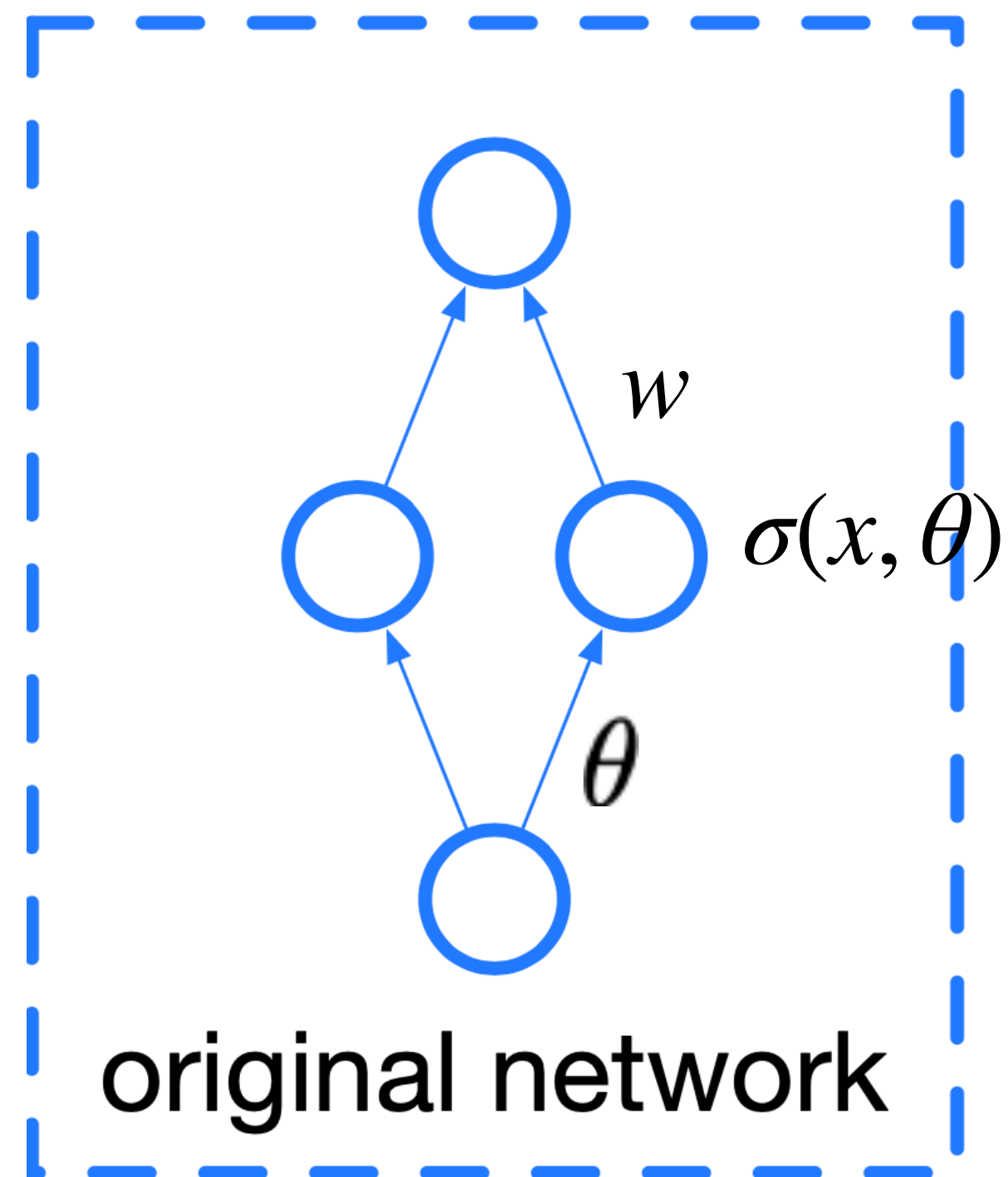
Minimize  $L(f)$  in  $f \in \partial(f, \epsilon)$   
without complexity constraint,  
get "over-grown" network  $\tilde{f}_{t+1}$

Architecture  
Optimization



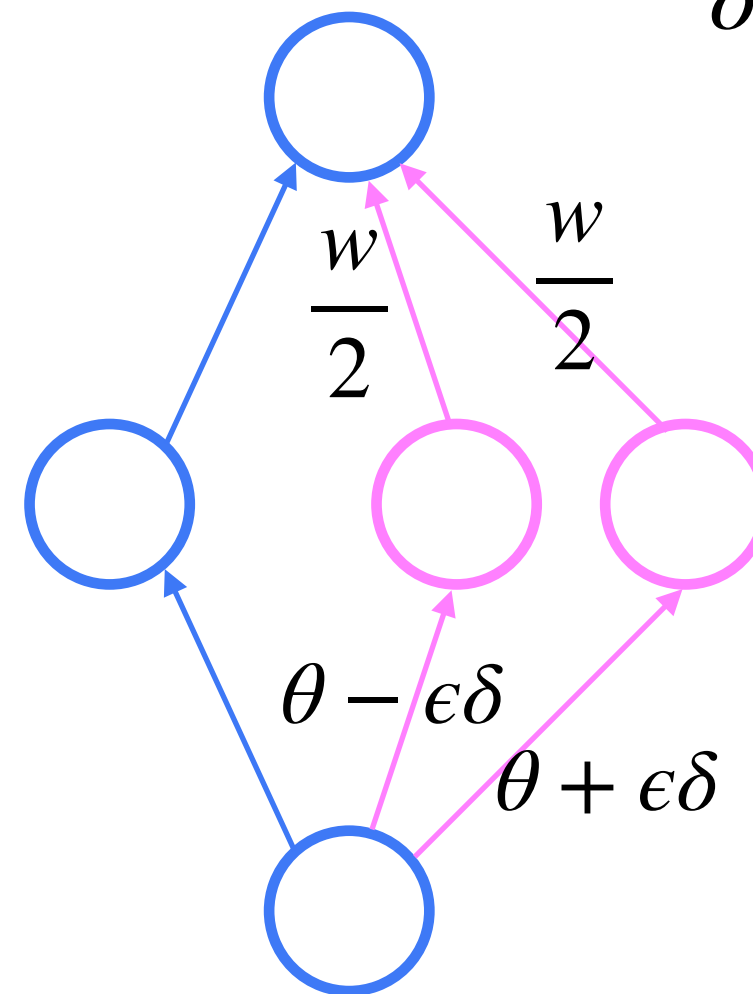
Select the top  $\eta_t$  neurons in  $\tilde{f}_{t+1}$   
with the highest importance  
measures to get  $f_{t+1}$

# Growing Network Width



## Splitting Existing Neurons

Growing New Neurons



$\epsilon$ : step size, indicating how much the network is changed

$\delta$ : update direction, indicating how the network is changed

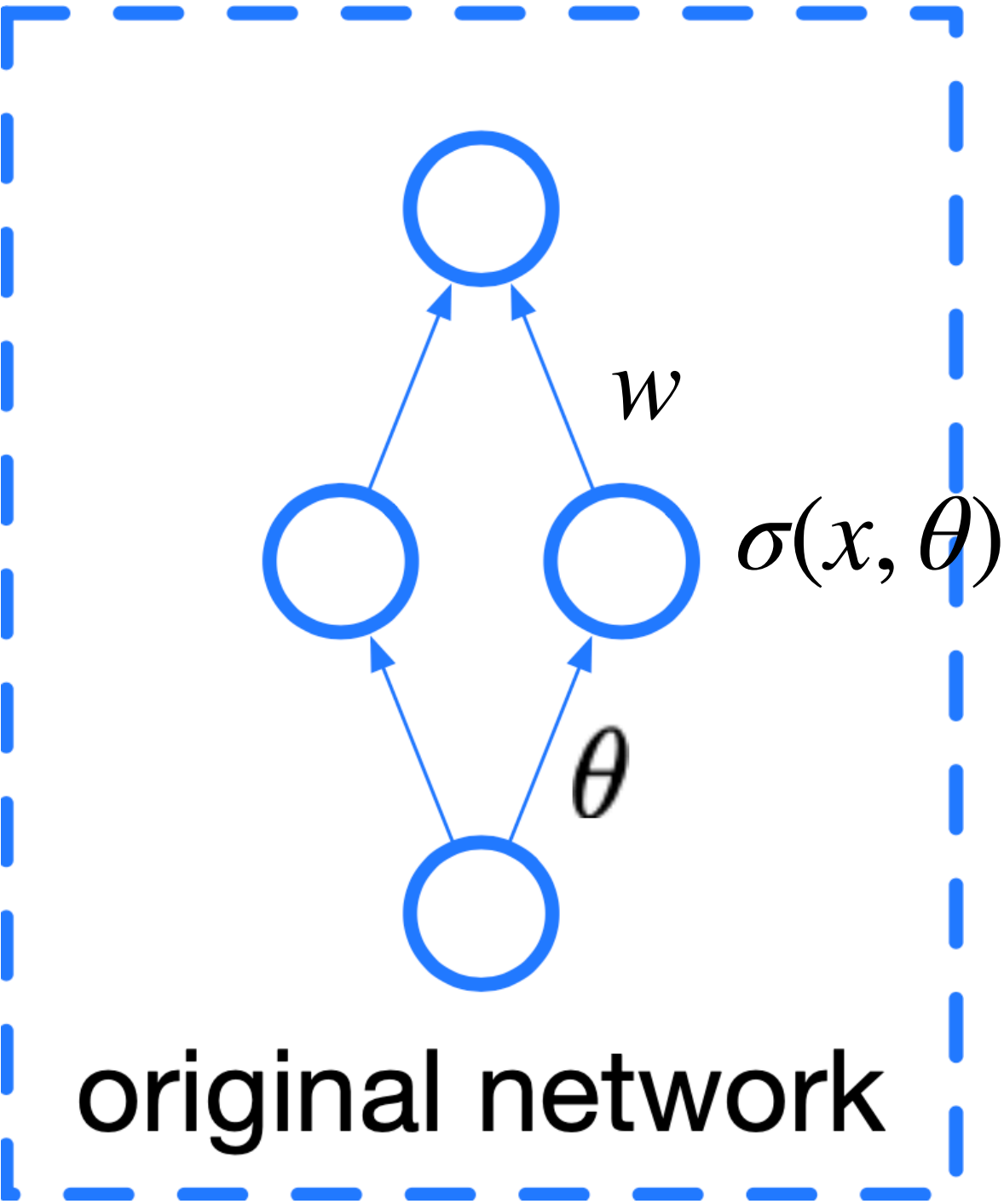
$$\sigma_1(x, \theta - \epsilon\delta) \approx \sigma(x, \theta)$$

$$\sigma_2(x, \theta + \epsilon\delta) \approx \sigma(x, \theta)$$

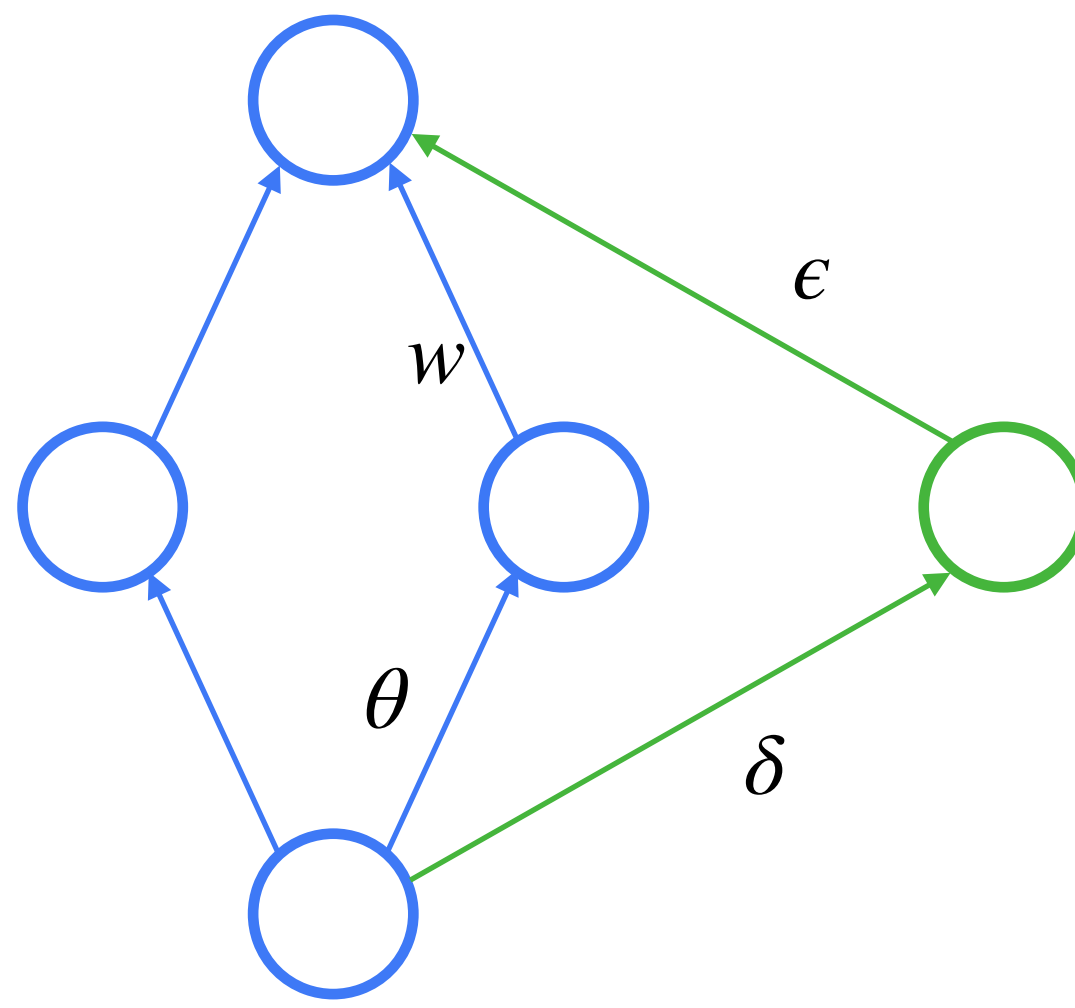
$$\frac{w}{2} \cdot \sigma_1 + \frac{w}{2} \cdot \sigma_2 \approx \frac{w}{2} \cdot \sigma + \frac{w}{2} \cdot \sigma = w \cdot \sigma(x, \theta)$$

$$f_{new} \approx f_{orig}$$

# Growing Network Width



Splitting Existing Neurons  
**Growing New Neurons**



$\epsilon$ : make sure the new network is close to original network

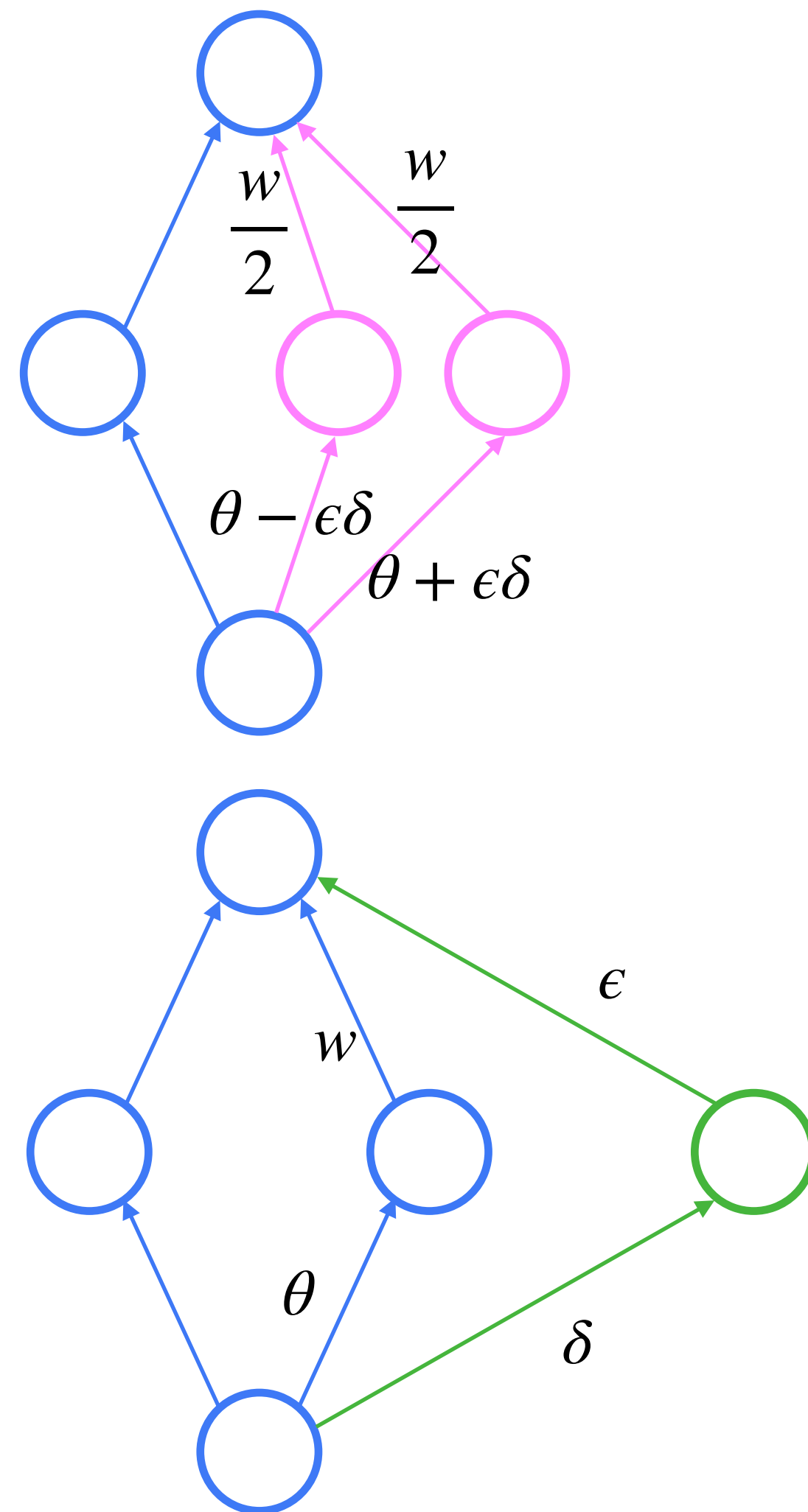
$\delta$ : a trainable parameter of the new neuron

$$\epsilon \cdot \sigma(x, \delta)$$

$$f_{new} \approx f_{orig}$$



# Growing Network Width



Overall, to grow  $f_t(x) = \sum_i \sigma(x; \theta_i)$  wider,  
the neighborhood set  $\partial(f_t, \epsilon)$  can include functions of the form

$$f_{\epsilon, \delta}(x) = \sum_{i=1}^m \frac{1}{2} (\sigma(x, \theta_i + \epsilon_i \delta_i) + \sigma(x, \theta_i - \epsilon_i \delta_i)) + \sum_{i=m+1}^{m+m'} \epsilon_i \sigma(x, \delta_i)$$

Up to  $m$  neurons to split

Up to  $m'$  neurons to add

Goal in Architecture Optimization:

- Instantiate neighborhood  $\partial(f_t, \epsilon)$  ✓
- Solve the optimization in  $f_{t+1}$

# Growing Network Width

Goal:

$$f_{t+1} = \arg \min_f \{L(f) \text{ s.t. } f \in \partial(f_t, \epsilon), C(f) \leq C(f_t) + \eta_t\}$$

$$f_{\epsilon, \delta}(x) = \sum_{i=1}^m \frac{1}{2} (\sigma(x, \theta_i + \epsilon_i \delta_i) + \sigma(x, \theta_i - \epsilon_i \delta_i)) + \sum_{i=m+1}^{m+m'} \epsilon_i \sigma(x, \delta_i)$$



Optimization Goal:

$$\min_{\epsilon, \delta} \left\{ L(f_{\epsilon, \delta}) \text{ s.t. } \|\epsilon\|_0 \leq \eta_t, \quad \|\epsilon\|_\infty \leq \epsilon, \quad \|\delta\|_{2, \infty} \leq 1 \right\}$$

# Growing Network Width

**Step One.** Optimize  $\delta$  and  $\epsilon$  without complexity constraints

**“Over-Grown” network**

$$[\tilde{\epsilon}, \tilde{\delta}] = \arg \min_{\epsilon, \delta} \left\{ L(f_{\epsilon, \delta}) \quad s.t. \quad \|\epsilon\|_{\infty} \leq \epsilon, \quad \|\delta\|_{2, \infty} \leq 1 \right\}.$$

$$\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_{m+m'}), \delta = (\delta_1, \delta_2, \dots, \delta_{m+m'})$$

$$\nabla_{\epsilon, \delta} L(f) = \left[ \frac{\partial L(f)}{\partial \epsilon_1}, \dots, \frac{\partial L(f)}{\partial \epsilon_{m+m'}}, \frac{\partial L(f)}{\partial \delta_1}, \dots, \frac{\partial L(f)}{\partial \delta_{m+m'}} \right]^T$$

Minimize  $L(f)$  in  $f \in \partial(f, \epsilon)$   
without complexity constraint,  
get “over-grown” network  $\tilde{f}_{t+1}$

Architecture  
Optimization



Select the top  $\eta_t$  neurons in  $\tilde{f}_{t+1}$   
with the highest importance  
measures to get  $f_{t+1}$

# Growing Network Width

**Step Two.** Re-Optimize  $\epsilon$  with Taylor approximation on the loss

$$L(f_{\epsilon, \tilde{\delta}}) = L(f) + \sum_{i=1}^{m+m'} \epsilon_i s_i + O(\epsilon^2), \quad \boxed{s_i} = \frac{1}{\tilde{\epsilon}_i} \int_0^{\tilde{\epsilon}_i} \nabla_{\zeta_i} L(f_{[\tilde{\epsilon}_{-i}, \zeta_i], \tilde{\delta}}) d\zeta_i$$

$[\tilde{\epsilon}_{-i}, \zeta_i]$  denotes replacing the  $i$ -th element of  $\tilde{\epsilon}$  with  $\zeta_i$

$$\begin{aligned} L(f_{\epsilon, \tilde{\delta}}) &= L(f_{\tilde{\epsilon}, \tilde{\delta}}) + \left. \frac{\partial L(f)}{\partial \epsilon_1} \right|_{\epsilon_1=0} \cdot \epsilon_1 + \left. \frac{\partial L(f)}{\partial \epsilon_2} \right|_{\epsilon_2=0} \cdot \epsilon_2 + \dots \\ &\quad + \left. \frac{\partial L(f)}{\partial \epsilon_{m+m'}} \right|_{\epsilon_{m+m'}=0} \cdot \epsilon_{m+m'} + O(\epsilon^2) \\ &= L(f) + \sum_{i=1}^{m+m'} \left. \frac{\partial L(f)}{\partial \epsilon_i} \right|_{\epsilon_i=0} \cdot \epsilon_i \end{aligned}$$

$$\frac{\partial L(f)}{\partial \epsilon_i} = \frac{\partial L(f_{\tilde{\epsilon}, \tilde{\delta}})}{\partial \epsilon_i} = \frac{1}{\tilde{\epsilon}_i} \int_0^{\tilde{\epsilon}_i} \nabla_{\zeta_i} L(f_{[\tilde{\epsilon}_{-i}, \zeta_i], \tilde{\delta}}) d\zeta_i$$

$[\tilde{\epsilon}_{-i}, \zeta_i]$  denotes replacing the  $i$ -th element of  $\tilde{\epsilon}$  with  $\zeta_i$

$s_i$ : an integrated gradient:  
measures the  
contribution of turning  
on the  $i$ -th new neuron

# Growing Network Width

**Step Two.** Re-Optimize  $\epsilon$  with Taylor approximation on the loss

$$L(f_{\epsilon, \tilde{\delta}}) = L(f) + \sum_{i=1}^{m+m'} \epsilon_i s_i + O(\epsilon^2), \quad \boxed{s_i} = \frac{1}{\tilde{\epsilon}_i} \int_0^{\tilde{\epsilon}_i} \nabla_{\zeta_i} L(f_{[\tilde{\epsilon}_{-i}, \zeta_i], \tilde{\delta}}) d\zeta_i$$

$[\tilde{\epsilon}_{-i}, \zeta_i]$  denotes replacing the  $i$ -th element of  $\tilde{\epsilon}$  with  $\zeta_i$

In practice, approximate the integration in  $s_i$  by discrete sampling:

$$s_i \approx \frac{1}{n} \sum_{z=1}^n \nabla_c c_z L(f_{[\tilde{\epsilon}_{-i}, c_z], \delta})$$
$$\hat{\epsilon} = \arg \min_{\epsilon} \left\{ \sum_{i=1}^{m+m'} \epsilon_i s_i \quad s.t. \quad \|\epsilon\|_0 \leq \eta_t, \quad \|\epsilon\|_{\infty} \leq \epsilon \right\}$$

$s_i$ : an integrated gradient:  
measures the  
contribution of turning  
on the  $i$ -th new neuron

$$f_{t+1} = f_{\hat{\epsilon}, \tilde{\delta}}$$

Goal in Architecture Optimization:

- Instantiate neighborhood  $\partial(f_t, \epsilon)$  ✓
- Solve the optimization in  $f_{t+1}$  ✓



# Experiment Result

**backbone:** VGG-19

**dataset:** CIFAR-10

**method:**

- Net2Net - grows networks uniformly by randomly selecting the existing neurons in each layer
- NASH - a random sampling search method using network morphism
- Firefly Descent

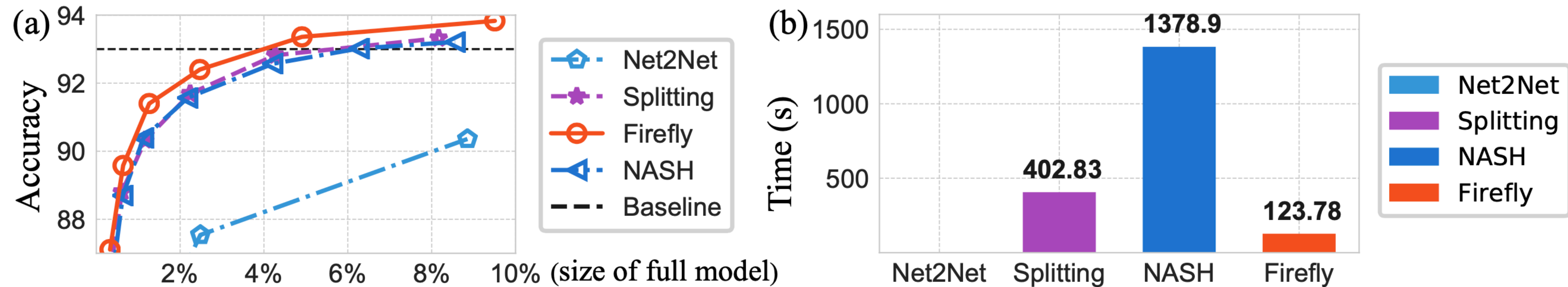


Figure 4: (a) Results of growing increasingly wider networks on CIFAR-10; VGG-19 is used as the backbone. (b) Computation time spent on growing for different methods.

- Outperform baseline test accuracy with less than 10% size of full model, growing fast
- Achieve comparable test accuracy as the full model with only 4% size of full model

# Experiment Result

**backbone:** cell-based architecture searching  
**dataset:** ImageNet

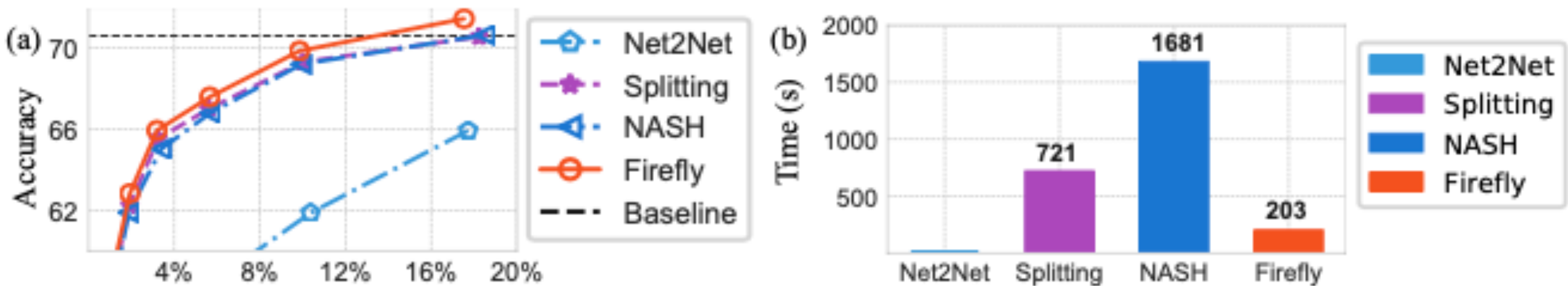
Table 1 reports the results comparing Firefly with several NAS baselines. Our method achieves a similar or better performance comparing with those RL-based and gradient-based methods like ENAS or DARTS, but with higher computational efficiency in terms of the total search time.

Method	Search Time (GPU Days)	Param (M)	Error
NASNet-A (Zoph et al., 2018)	2000	3.1	2.83
ENAS (Pham et al., 2018)	4	4.2	2.91
Random Search	4	3.2	$3.29 \pm 0.15$
DARTS (first order) (Liu et al., 2018b)	1.5	3.3	$3.00 \pm 0.14$
DARTS (second order) (Liu et al., 2018b)	4	3.3	$2.76 \pm 0.09$
Firefly	1.5	3.3	$2.78 \pm 0.05$

Table 1: Performance compared with several NAS baseline

- Achieves a similar or better performance comparing with several NAS methods,
- Much less search time

**backbone:** MobileNetV1  
**dataset:** CIFAR-100



## **Contribution:**

- Smoothly change the network, grow the networks without accuracy drop
- Flexible growing width and depth

## **Drawbacks:**

- 通过gradient optimization将“search”和“grow”的过程combine在一起, continuous growing
- one-layer grow?
- 每一次architecture optimization (grow) 后需要重新训练新网络
- 实验无精确的实验结果, 对比模型较naive
- Formula without proof