

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

## Model selection problem

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

$\Omega$  – family of models

$C(f)$  – complexity of model

# Model selection problem

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

$\partial(f_t, \epsilon)$  – neighbourhood of  $f_t$

---

**Algorithm 1** Firefly Neural Architecture Descent

---

**Input:** Loss function  $L(f)$ ; initial small network  $f_0$ ; search neighborhood  $\partial(f, \epsilon)$ ; maximum increase of size  $\{\eta_t\}$ .

**Repeat:** At the  $t$ -th growing phase:

1. Optimize the parameter of  $f_t$  with fixed structure using a typical optimizer for several epochs.
  2. Minimize  $L(f)$  in  $f \in \partial(f, \epsilon)$  without the complexity constraint (see e.g., (4)) to get a large “over-grown” network  $\tilde{f}_{t+1}$  by performing gradient descent.
  3. Select the top  $\eta_t$  neurons in  $\tilde{f}_{t+1}$  with the highest importance measures to get  $f_{t+1}$  (see (5)).
-

# Neighbourhood

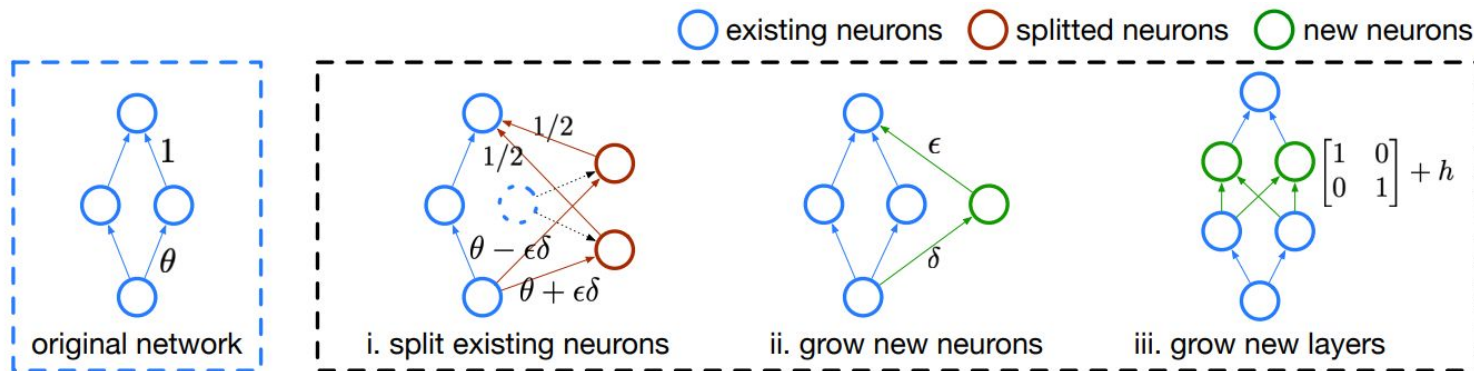
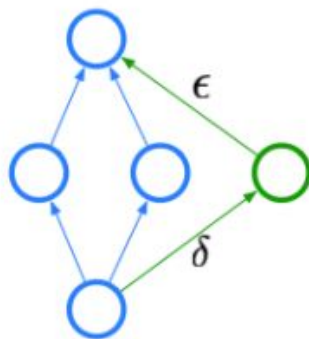
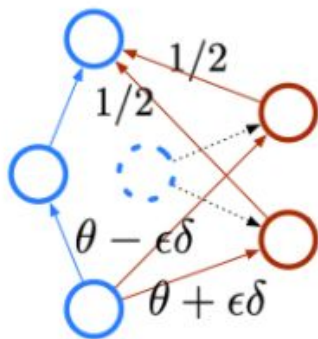


Figure 1: An illustration of three different growing methods within firefly neural architecture descent. Both  $\delta$  and  $h$  are trainable perturbations.

# Splitting and Growing Neuron



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

# Optimization

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

Step 1

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

# Optimization

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

## Step 2

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

$$\hat{\epsilon} = \arg \min_{\epsilon} \left\{ \sum_{i=1}^{m+m'} \epsilon_i s_i \quad \text{s.t.} \quad \|\epsilon\|_0 \leq \eta_t, \quad \|\epsilon\|_{\infty} \leq \epsilon \right\}$$



# Growing new layer

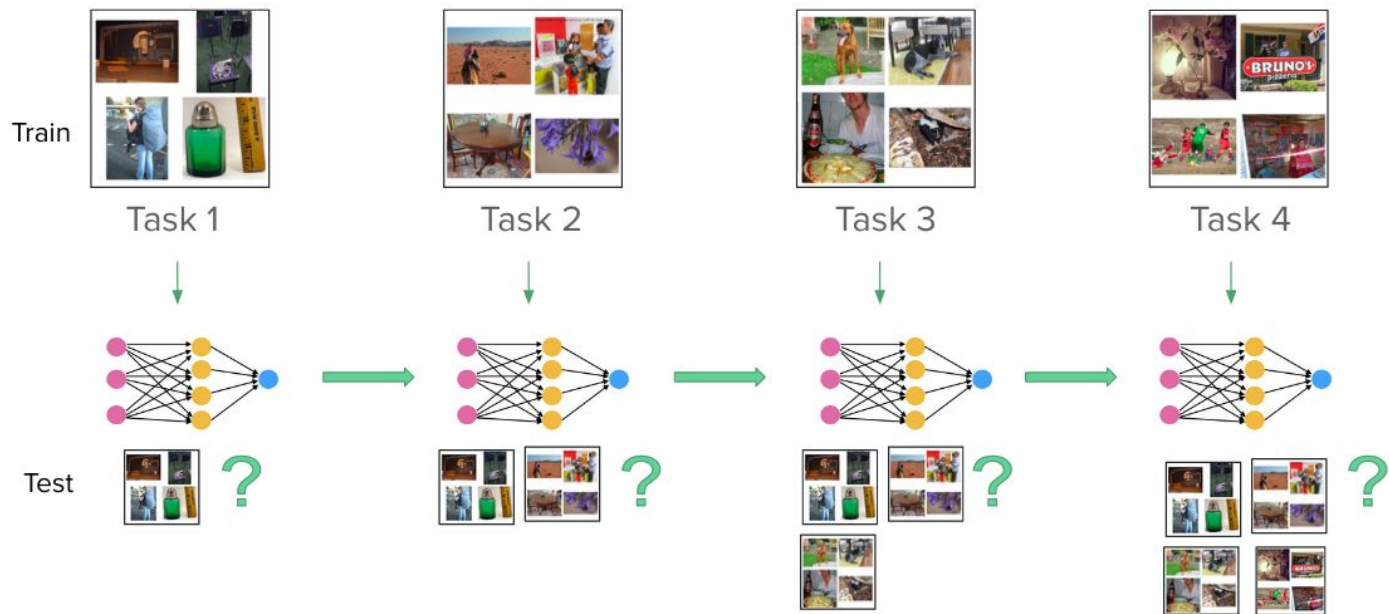
New layer representation

$$f_{\varepsilon, \delta} = g_d \circ (I + h_{d-1}) \cdots (I + h_2) \circ g_2 \circ (I + h_1) \circ g_1, \quad \text{with} \quad h_\ell(\cdot) = \sum_{i=1}^{m'} \varepsilon_{\ell i} \sigma(\cdot, \delta_{\ell i})$$

Optimization problem to solve

$$\min_{\varepsilon, \delta} \{L(f_{\varepsilon, \delta}) \text{ s.t. } \|\varepsilon\|_0 \leq \eta_{t,0}, \quad \underline{\|\varepsilon\|_{\infty,0} \leq \eta_{t,1}}, \quad \|\varepsilon\|_{\infty} \leq \epsilon, \quad \|\delta\|_{2,\infty} \leq 1\},$$

# Continual learning



# Continual learning

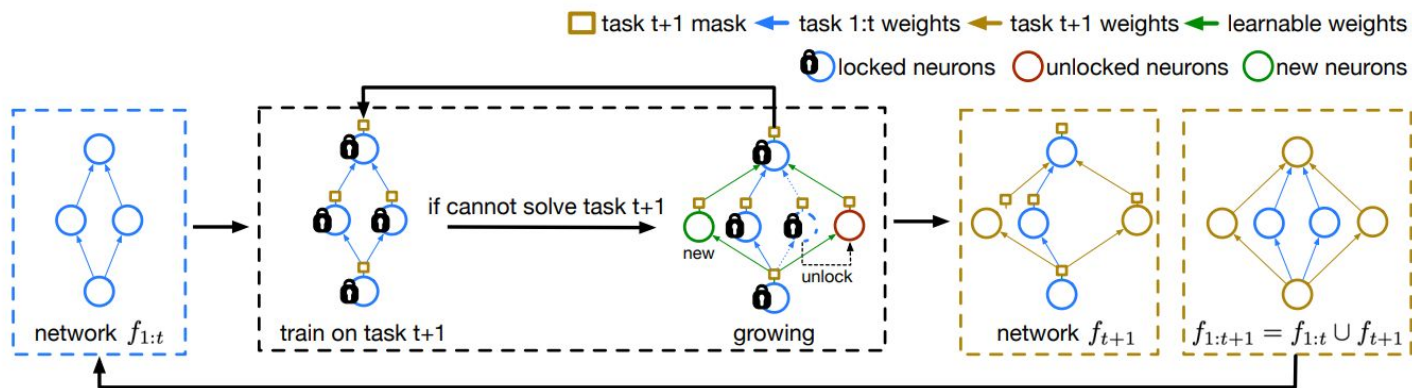


Figure 2: Illustration of how Firefly grows networks in continual learning.

# Firefly vs Random

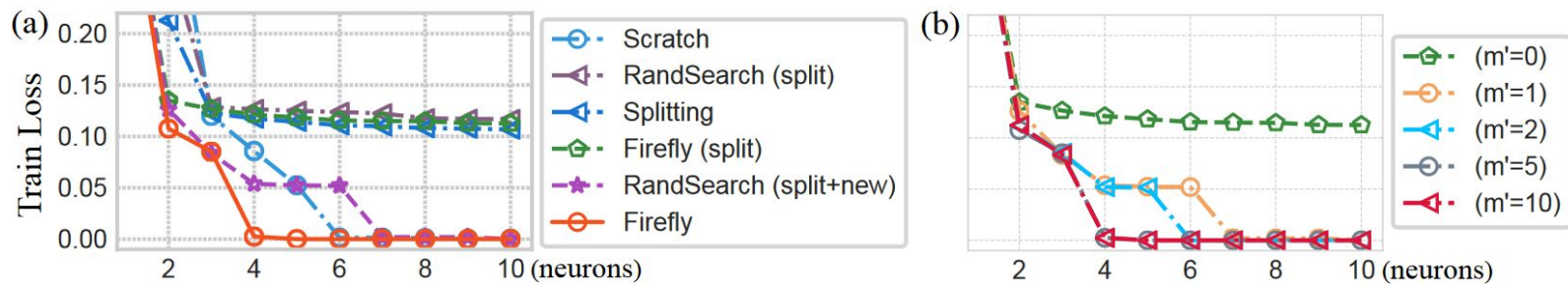


Figure 3: (a) Average training loss of different growing methods versus the number of grown neurons. (b) Firefly descent with different numbers of new neuron candidates.

# Comparison with other methods

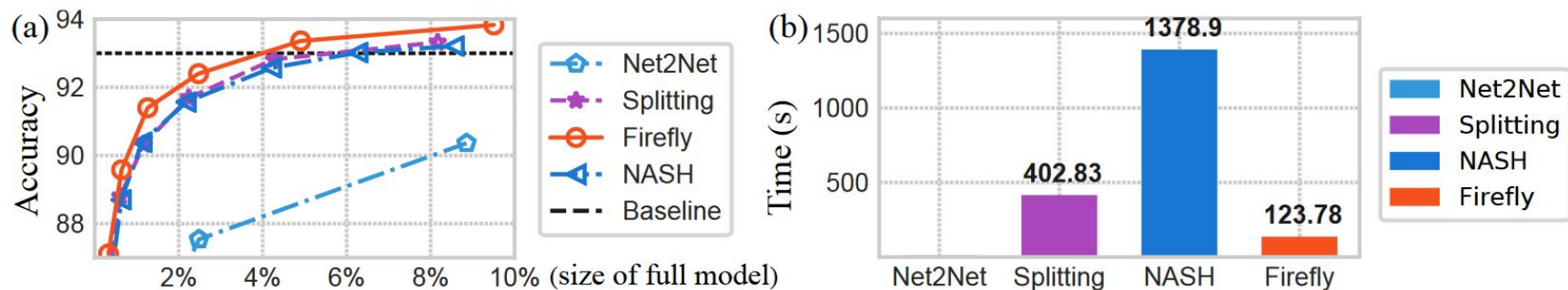


Figure 4: (a) Results of growing increasingly wider networks on CIFAR-10; VGG-19 is used as the backbone. (b) Computation time spent on growing for different methods.

