# Neural Architecture Search

Neural Networks Design And Application

# A framework for NAS
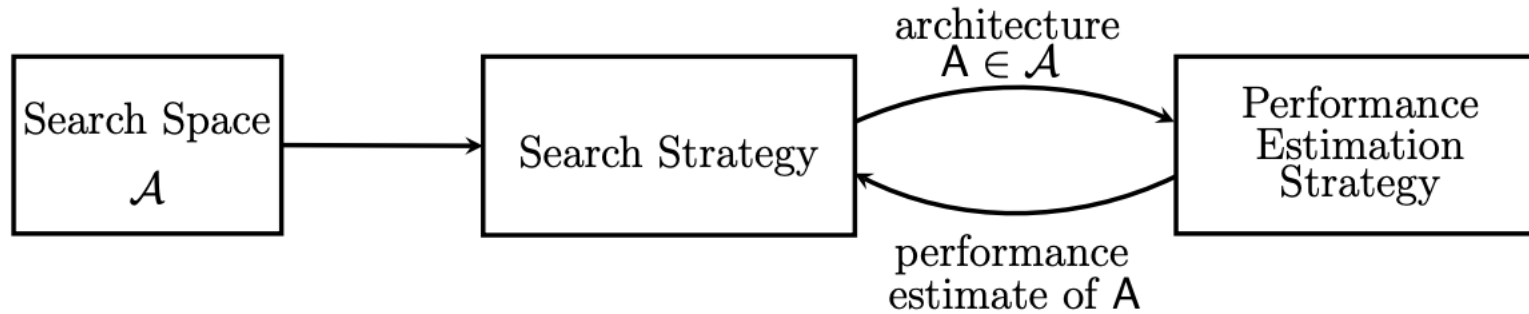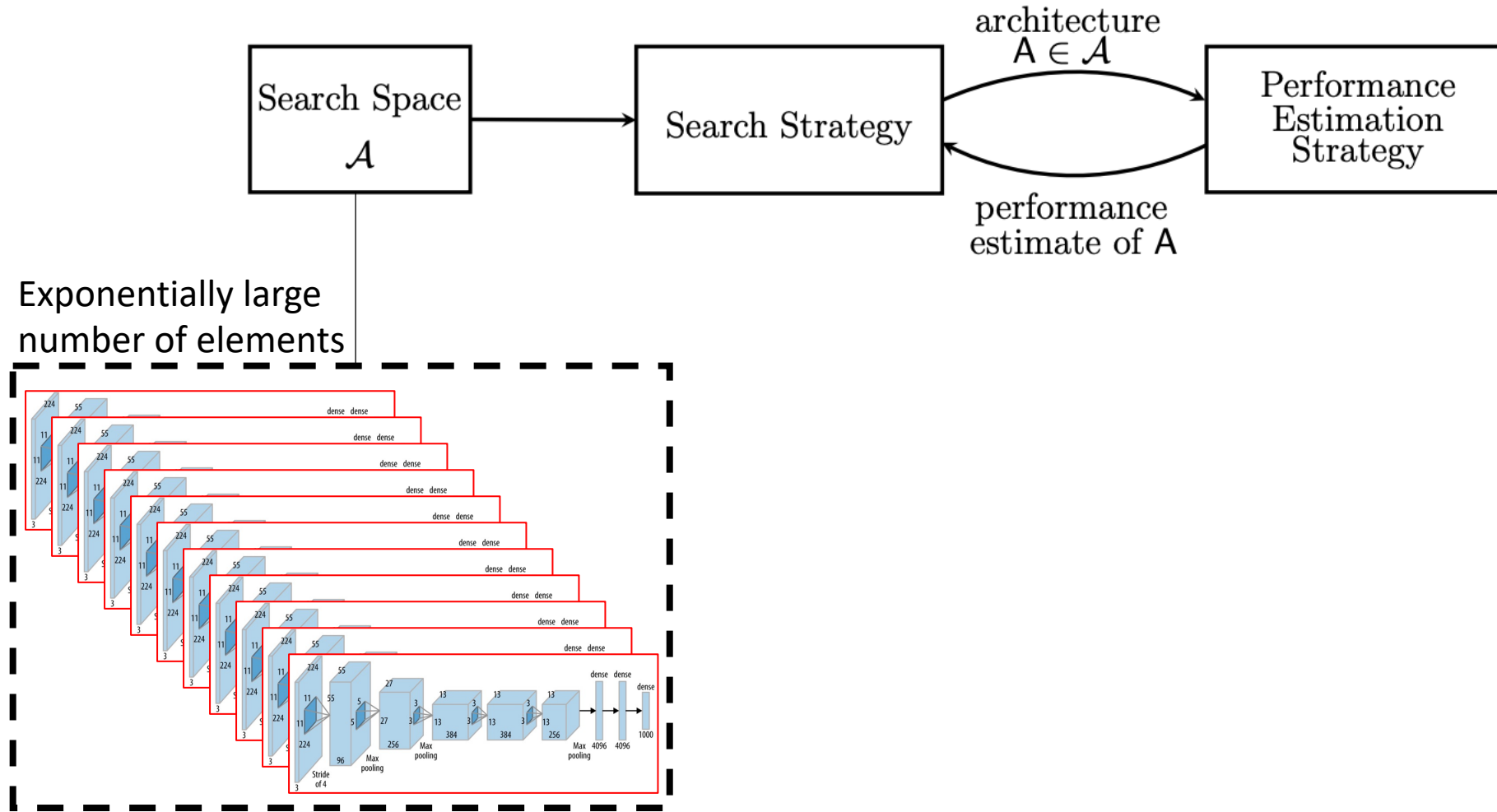
# A framework for NAS
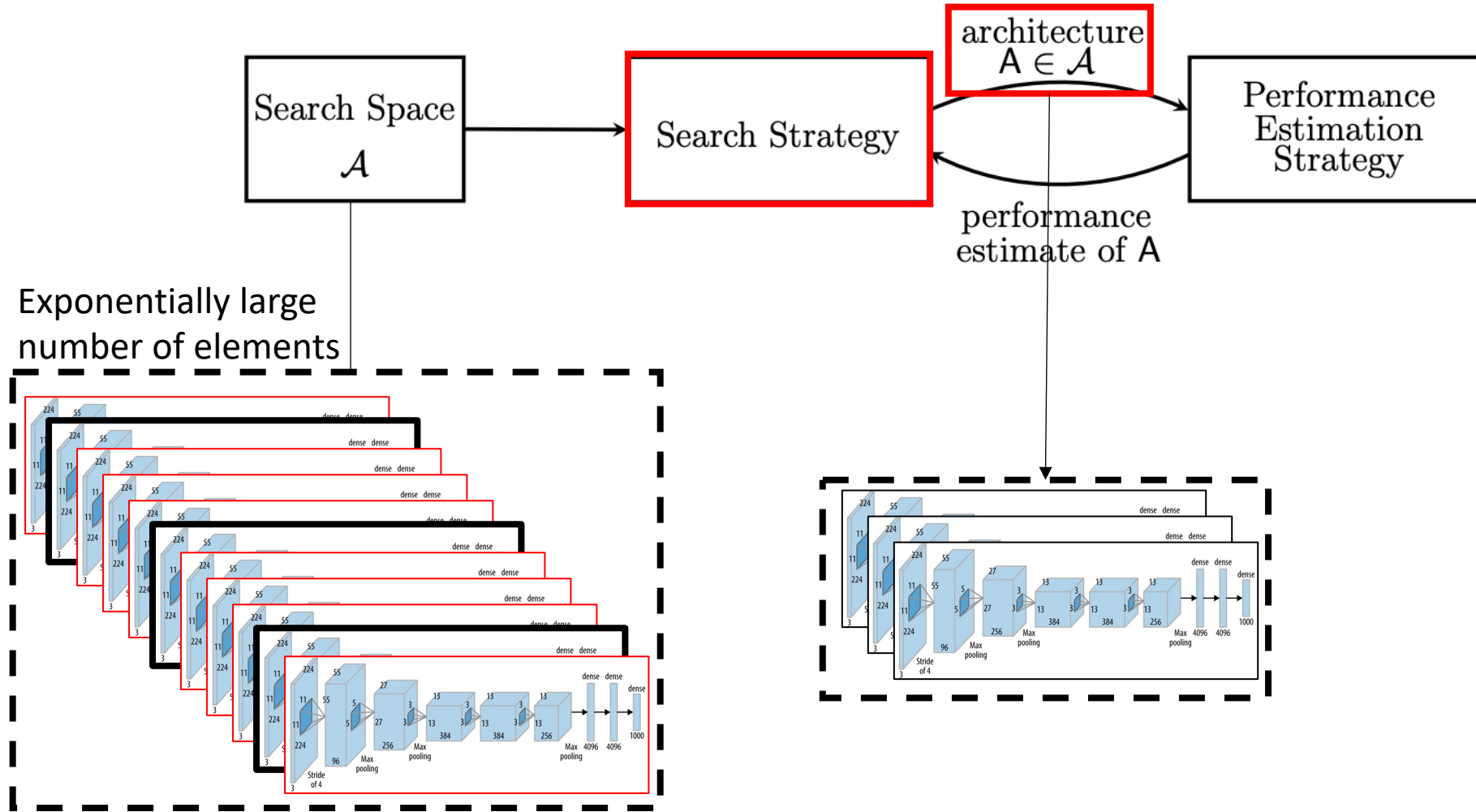
# A framework for NAS



Exponentially large number of elements

# Random search

- Learning rate: {0.0001, 0.001, 0.01, 0.1, 1, 10}
- Weight decay (L2 regularization): {0.0001, 0.001, 0.01, 0.1, 1, 10}

# Random search

- Learning rate: {0.0001, 0.001, 0.01, 0.1, 1, 10}
- Weight decay (L2 regularization): {0.0001, 0.001, 0.01, 0.1, 1, 10}



Grid Layout

Unimportant parameter

Important parameter

https://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf

# Random search

- Learning rate: {0.0001, 0.001, 0.01, 0.1, 1, 10}
- Weight decay (L2 regularization): {0.0001, 0.001, 0.01, 0.1, 1, 10}



Grid Layout

Unimportant parameter

Important parameter

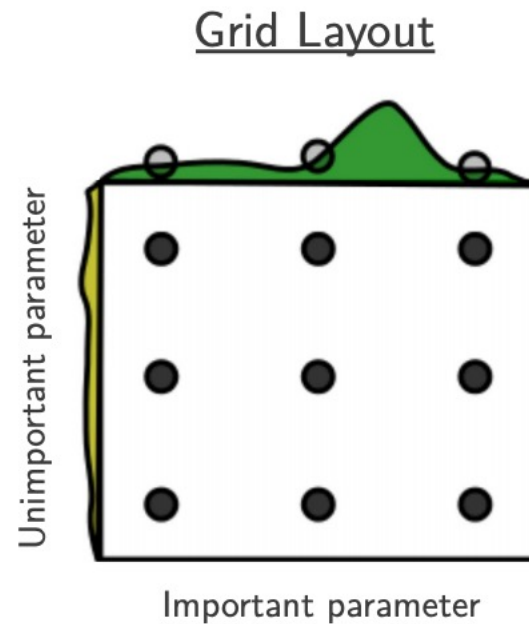https://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf

# Random search

- Learning rate: {0.0001, 0.001, 0.01, 0.1, 1, 10}
- Weight decay (L2 regularization): {0.0001, 0.001, 0.01, 0.1, 1, 10}

Grid Layout



Unimportant parameter

Important parameter

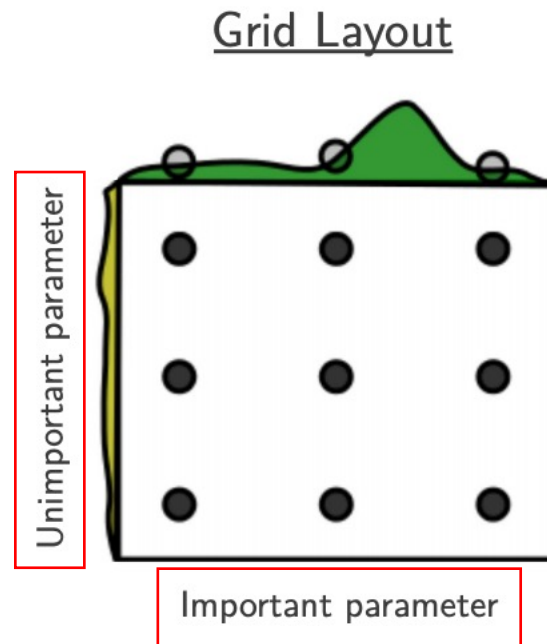https://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf

# Random search

- Learning rate: {0.0001, 0.001, 0.01, 0.1, 1, 10}
- Weight decay (L2 regularization): {0.0001, 0.001, 0.01, 0.1, 1, 10}



Grid Layout

https://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf

# Random search

- Learning rate: {0.0001, 0.001, 0.01, 0.1, 1, 10}
- Weight decay (L2 regularization): {0.0001, 0.001, 0.01, 0.1, 1, 10}
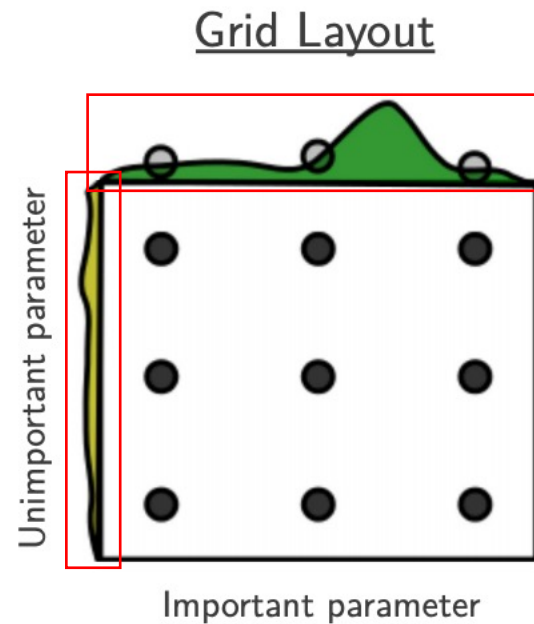
https://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf

# Random search

- Learning rate: {0.0001, 0.001, 0.01, 0.1, 1, 10}
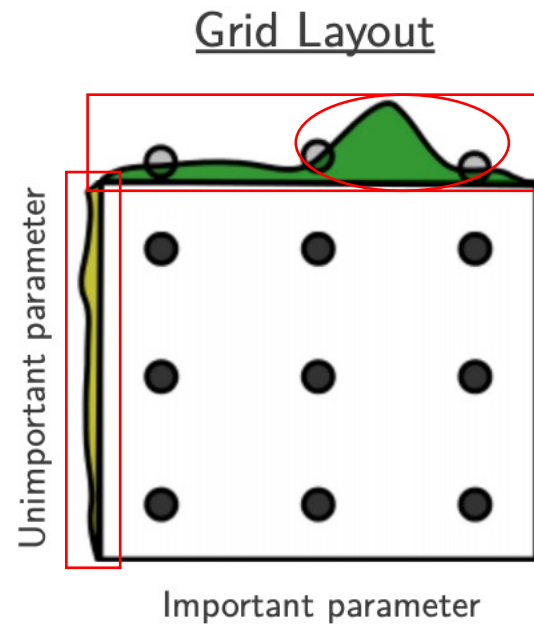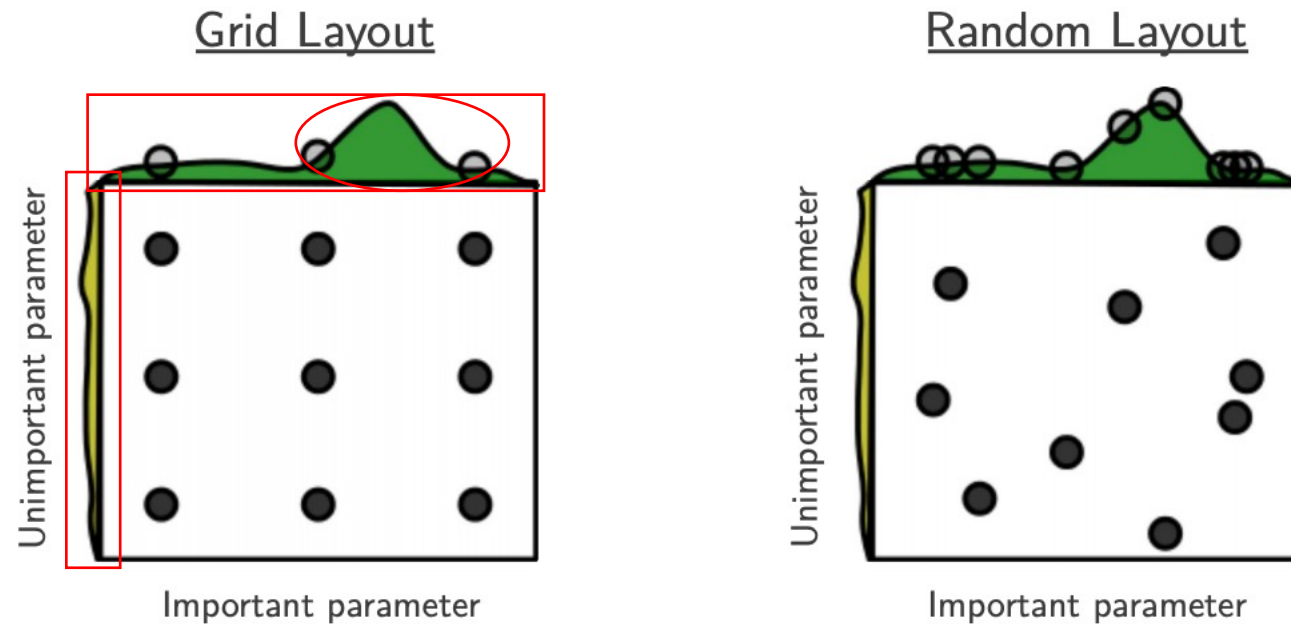- Weight decay (L2 regularization): {0.0001, 0.001, 0.01, 0.1, 1, 10}

https://www.jmlr.org/papers/volume13/bergstra12a/bergstra12a.pdf

# Reinforcement learning



Sample architecture A with probability p

The controller (RNN)

Trains a child network with architecture $A$ to get reward $R$

Compute gradient of p and scale it by R to update the controller

Image credit: https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Reinforcement learning



Sample architecture A with probability p

The controller (RNN)

Trains a child network with architecture $A$ to get reward $R$

Compute gradient of p and scale it by R to update the controller

Image credit: https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Reinforcement learning



Sample architecture A with probability p

The controller (RNN)

Trains a child network with architecture $A$ to get reward $R$

Compute gradient of p and scale it by R to update the controller

Image credit: https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Reinforcement learning



Image credit: https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Reinforcement learning

- **Action space**: The action space is a list of tokens for defining a child network predicted by the controller (See more in the above section). The controller outputs $action$, $a_{1:T}$, where $T$ is the total number of tokens.
- **Reward**: The accuracy of a child network that can be achieved at convergence is the reward for training the controller, $R$.
- **Loss**: NAS optimizes the controller parameters $\theta$ with a REINFORCE loss. We want to maximize the expected reward (high accuracy) with the gradient as follows. The nice thing here with policy gradient is that it works even when the reward is non-differentiable.

$$\nabla_\theta J(\theta) = \sum_{t=1}^{T} \mathbb{E}[\nabla_\theta \log P(a_t|a_{1:(t-1)}; \theta)R]$$

https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html
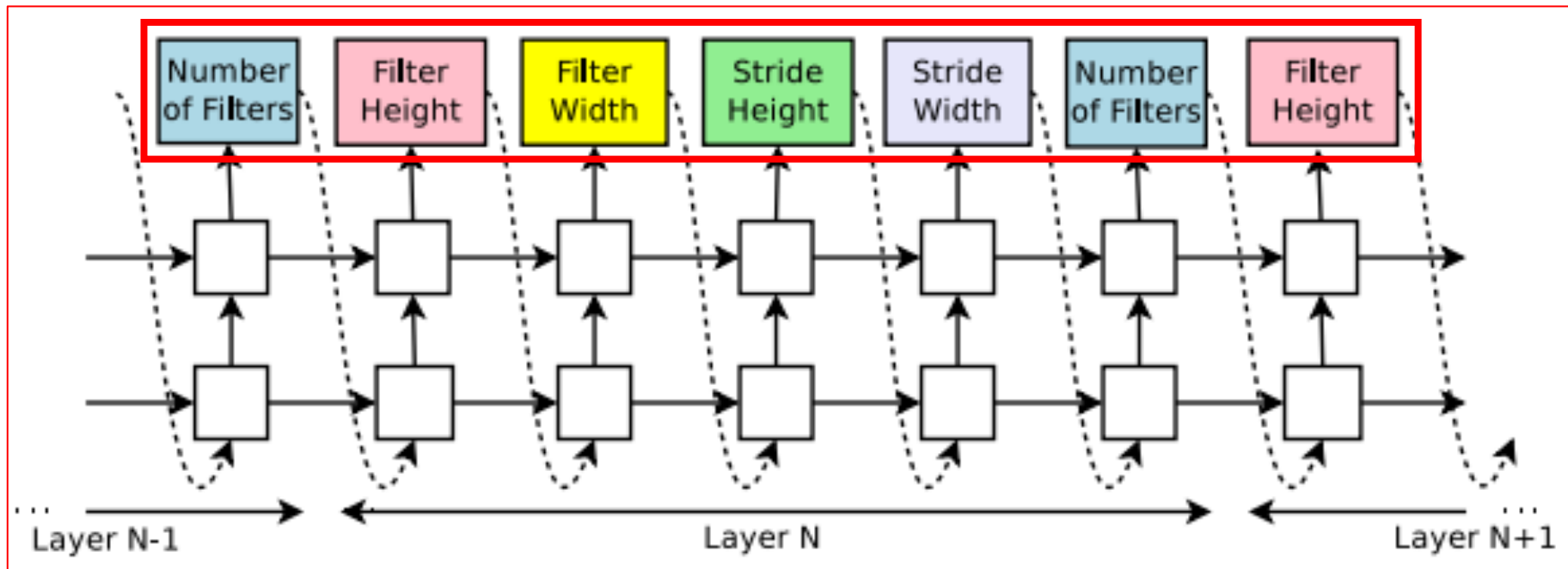
# Reinforcement learning

- **Action space** The action space is a list of tokens for defining a child network predicted by the controller (See more in the above section). The controller outputs $action$, $a_{1:T}$, where $T$ is the total number of tokens.

- **Reward**: The accuracy of a child network that can be achieved at convergence is the reward for training the controller, $R$.

- **Loss**: NAS optimizes the controller parameters $\theta$ with a REINFORCE loss. We want to maximize the expected reward (high accuracy) with the gradient as follows. The nice thing here with policy gradient is that it works even when the reward is non-differentiable.

$$\nabla_\theta J(\theta) = \sum_{t=1}^{T} \mathbb{E}[\nabla_\theta \log P(a_t|a_{1:(t-1)}; \theta)R]$$

https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Defining search space



A RNN network

Zoph, Barret, and Quoc V. Le. "Neural architecture search with reinforcement learning." *arXiv preprint arXiv:1611.01578* (2016). https://arxiv.org/pdf/1611.01578.pdf

# Reinforcement learning

- **Action space** The action space is a list of tokens for defining a child network predicted by the controller (See more in the above section). The controller outputs $action, a_{1:T}$ where $T$ is the total number of tokens.

- **Reward**: The accuracy of a child network that can be achieved at convergence is the reward for training the controller, $R$.

- **Loss**: NAS optimizes the controller parameters $\theta$ with a REINFORCE loss. We want to maximize the expected reward (high accuracy) with the gradient as follows. The nice thing here with policy gradient is that it works even when the reward is non-differentiable.

$$\nabla_\theta J(\theta) = \sum_{t=1}^{T} \mathbb{E}[\nabla_\theta \log P(a_t|a_{1:(t-1)}; \theta)R]$$

https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html
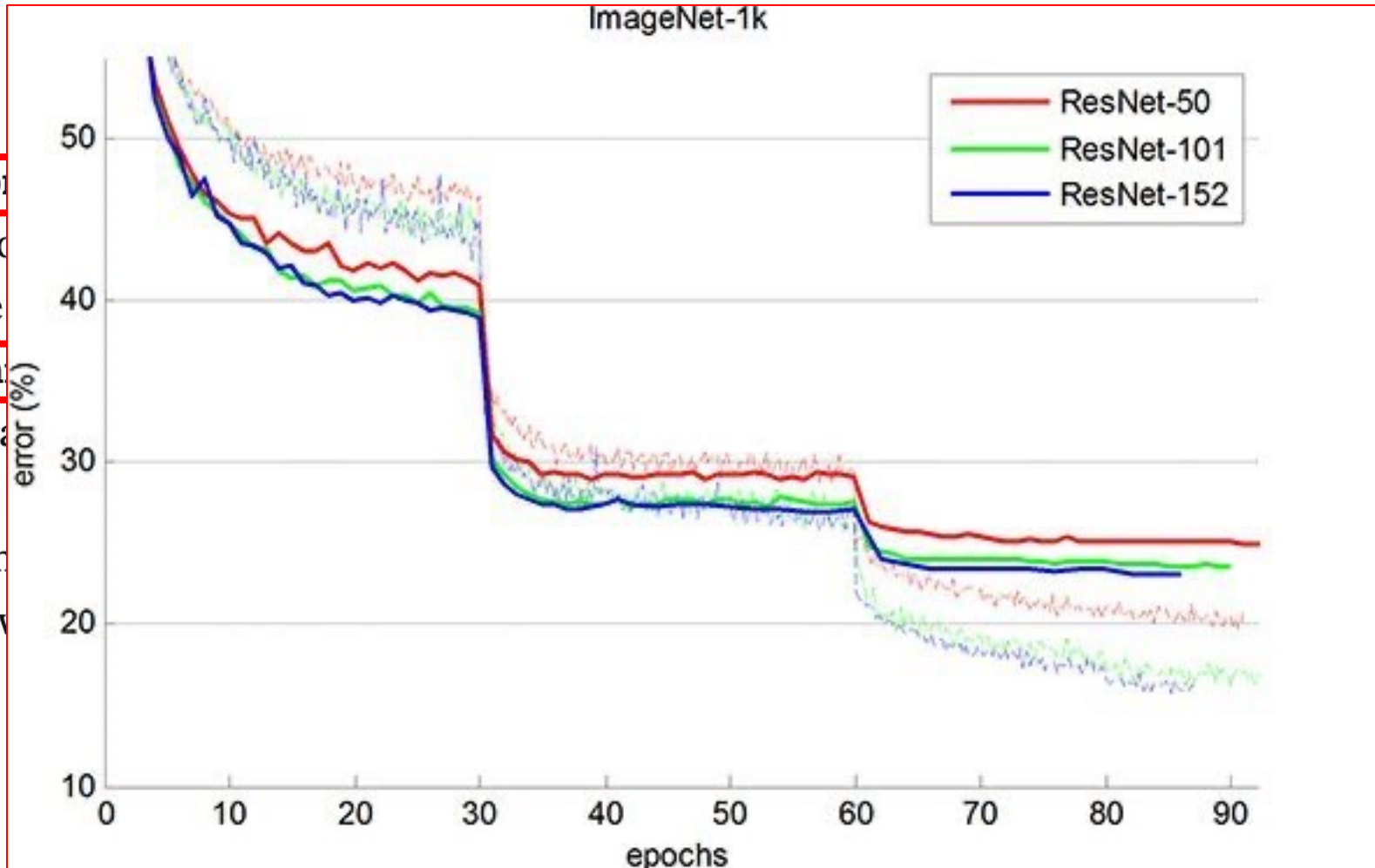
# Reinforcement learning

- **Action space** The action space is a list of tokens for defining a child network predicted by the controller (See more in the above section). The controller outputs $action, a_{1:T}$ where $T$ is the total number of tokens.
- **Reward:** The accuracy of a child network that can be achieved at convergence is the reward for training the controller, $R$.
- **Loss:** NAS optimizes the controller parameters $\theta$ with a REINFORCE loss. We want to maximize the expected reward (high accuracy) with the gradient as follows. The nice thing here with policy gradient is that it works even when the reward is non-differentiable.

$$\nabla_\theta J(\theta) = \sum_{t=1}^{T} \mathbb{E}[\nabla_\theta \log P(a_t | a_{1:(t-1)}; \theta) R]$$

20

https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Reinforcement learning



- **Action**: ... ted by the co... is the ... where $T$
- **Reward**: ... reward for tra...
- **Loss**: ... maxim... thing here ... e.

https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Reinforcement learning

- **Action space** The action space is a list of tokens for defining a child network predicted by the controller (See more in the above section). The controller outputs $action, a_{1:T}$ where $T$ is the total number of tokens.
- **Reward**: The accuracy of a child network that can be achieved at convergence is the reward for training the controller, $R$.
- **Loss**: NAS optimizes the controller parameters $\theta$ with a REINFORCE loss. We want to maximize the expected reward (high accuracy) with the gradient as follows. The nice thing here with policy gradient is that it works even when the reward is non-differentiable.

$$\nabla_\theta J(\theta) = \sum_{t=1}^{T} \mathbb{E}[\nabla_\theta \log P(a_t | a_{1:(t-1)}; \theta) R]$$

https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Reinforcement learning

- **Action space** The action space is a list of tokens for defining a child network predicted by the controller (See more in the above section). The controller outputs $action, a_{1:T}$ where $T$ is the total number of tokens.
- **Reward**: The accuracy of a child network that can be achieved at convergence is the reward for training the controller, $R$.
- **Loss**: NAS optimizes the controller parameters $\theta$ with a REINFORCE loss. We want to maximize the expected reward (high accuracy) with the gradient as follows. The nice thing here with policy gradient is that it works even when the reward is non-differentiable.

$$\nabla_{\theta} J(\theta) = \sum_{t=1}^{T} \mathbb{E}[\nabla_{\theta} \log P(a_t | a_{1:(t-1)}; \theta) R]$$

https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Reinforcement learning

- **Action space** The action space is a list of tokens for defining a child network predicted by the controller (See more in the above section). The controller outputs action, $a_{1:T}$ where $T$ is the total number of tokens.
- **Reward**: The accuracy of a child network that can be achieved at convergence is the reward for training the controller, $R$.
- **Loss**: NAS optimizes the controller parameters $\theta$ with a REINFORCE loss. We want to maximize the expected reward (high accuracy) with the gradient as follows. The nice thing here with policy gradient is that it works even when the reward is non-differentiable.

$$\nabla_\theta J(\theta) = \sum_{t=1}^{T} \mathbb{E}[\nabla_\theta \log P(a_t|a_{1:(t-1)}; \theta)R]$$

https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# Reinforcement learning

- **Action space** The action space is a list of tokens for defining a child network predicted by the controller (See more in the above section). The controller outputs action, $a_{1:T}$ where $T$ is the total number of tokens.
- **Reward**: The accuracy of a child network that can be achieved at convergence is the reward for training the controller, $R$.
- **Loss**: NAS optimizes the controller parameters $\theta$ with a REINFORCE loss. We want to maximize the expected reward (high accuracy) with the gradient as follows. The nice thing here with policy gradient is that it works even when the reward is non-differentiable.

$$\nabla_\theta J(\theta) = \sum_{t=1}^{T} \mathbb{E}[\nabla_\theta \log P(a_t | a_{1:(t-1)}; \theta) R]$$

https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html
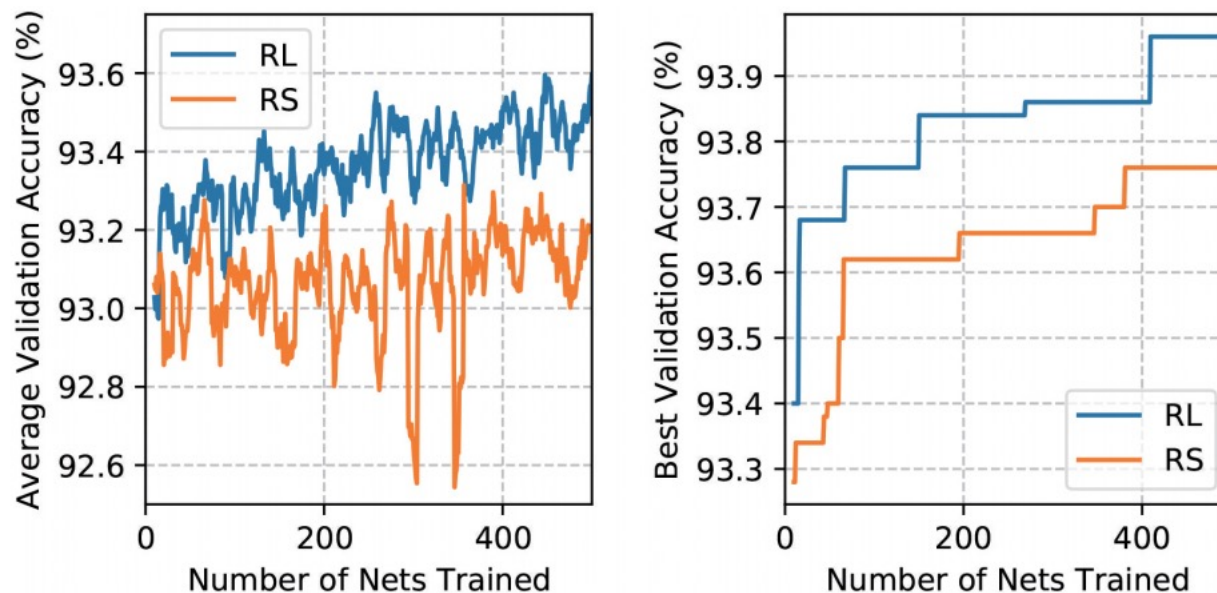
# Reinforcement learning

- **Action space** The action space is a list of tokens for defining a child network predicted by the controller (See more in the above section). The controller outputs $action, a_{1:T}$ where $T$ is the total number of tokens.
- **Reward**: The accuracy of a child network that can be achieved at convergence is the reward for training the controller, $R$.
- **Loss**: NAS optimizes the controller parameters $\theta$ with a REINFORCE loss. We want to maximize the expected reward (high accuracy) with the gradient as follows. The nice thing here with policy gradient is that it works even when the reward is non-differentiable.

$$\nabla_\theta J(\theta) = \sum_{t=1}^{T} \mathbb{E}[\nabla_\theta \log P(a_t | a_{1:(t-1)}; \theta) R]$$

$$\mathrm{E}[X] = \sum_{i=1}^{k} x_i \, p_i = x_1 p_1 + x_2 p_2 + \cdots + x_k p_k.$$

$$L(\theta \,|\, x) = \prod_{i=1}^{n} p_X(x_i \,|\, \theta)$$

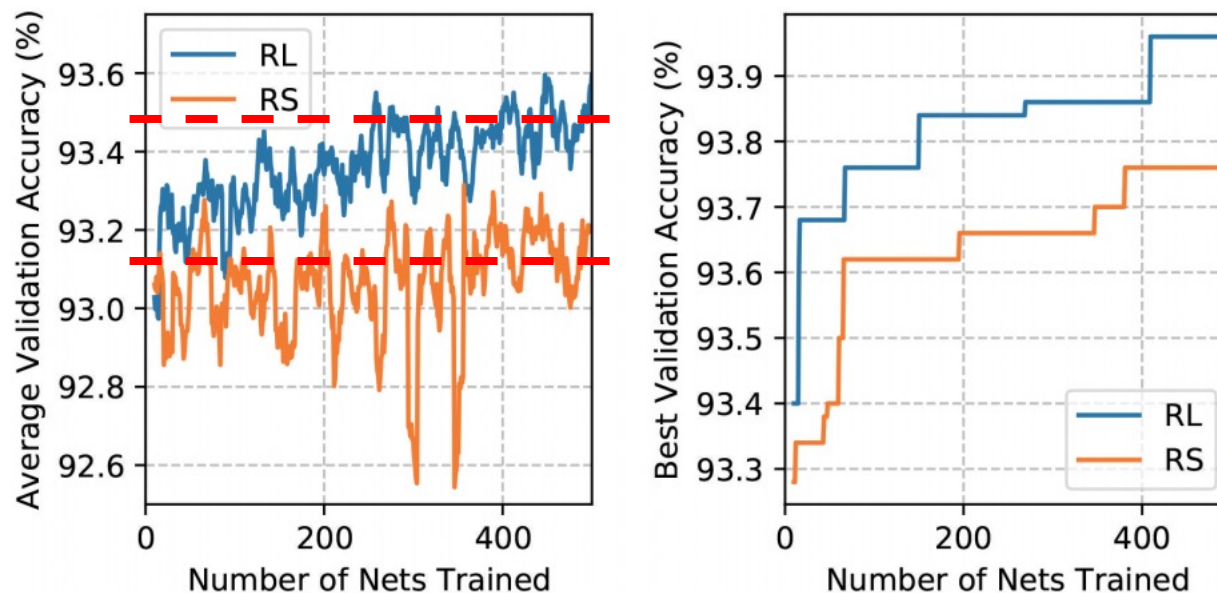https://lilianweng.github.io/lil-log/2020/08/06/neural-architecture-search.html

# RS vs RL



Cai, Han, Jiacheng Yang, Weinan Zhang, Song Han, and Yong Yu. "Path-level network transformation for efficient architecture search." In *International Conference on Machine Learning*, pp. 678-687. PMLR, 2018.
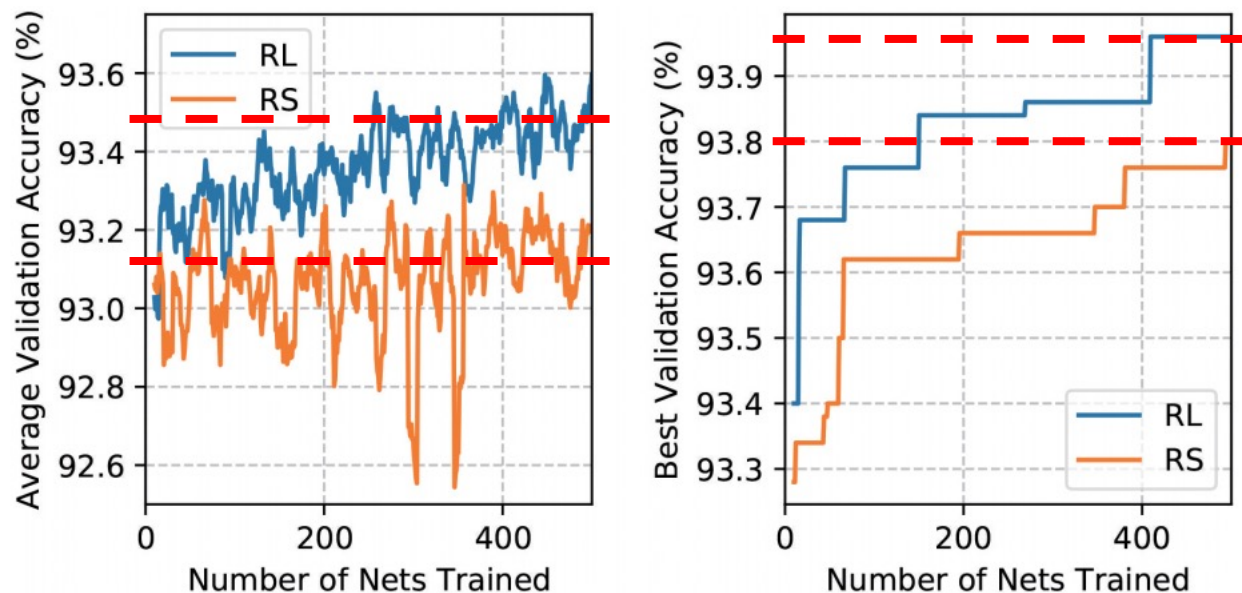https://arxiv.org/pdf/1806.02639.pdf

27

# RS vs RL



Cai, Han, Jiacheng Yang, Weinan Zhang, Song Han, and Yong Yu. "Path-level network transformation for efficient architecture search." In *International Conference on Machine Learning*, pp. 678-687. PMLR, 2018.
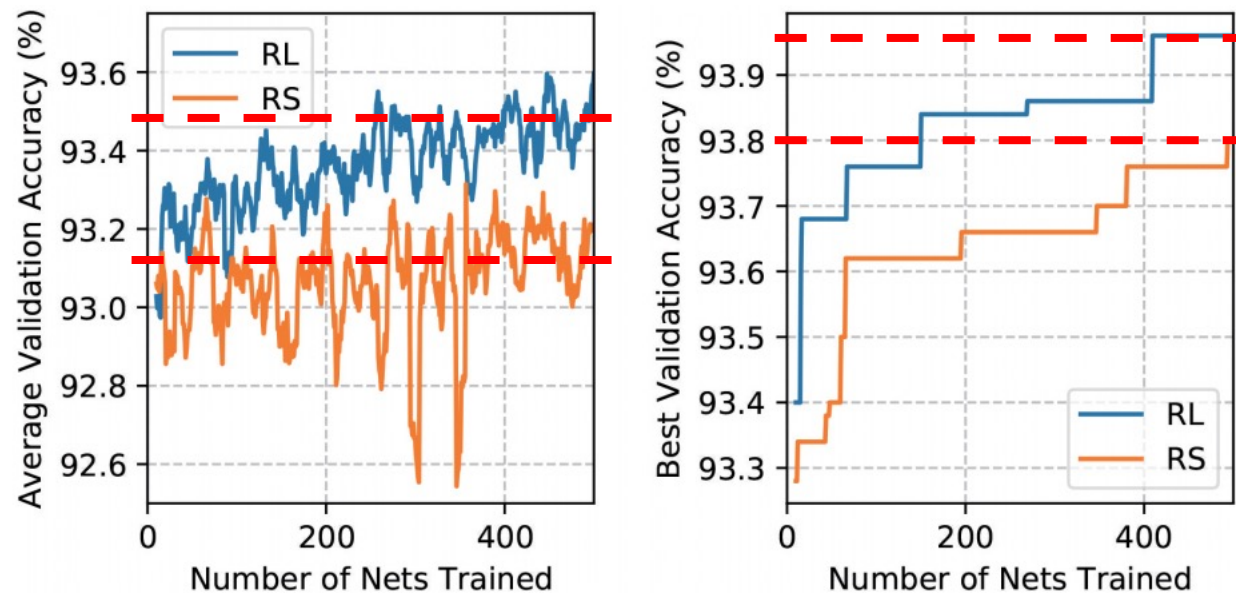https://arxiv.org/pdf/1806.02639.pdf

# RS vs RL



Cai, Han, Jiacheng Yang, Weinan Zhang, Song Han, and Yong Yu. "Path-level network transformation for efficient architecture search." In *International Conference on Machine Learning*, pp. 678-687. PMLR, 2018.
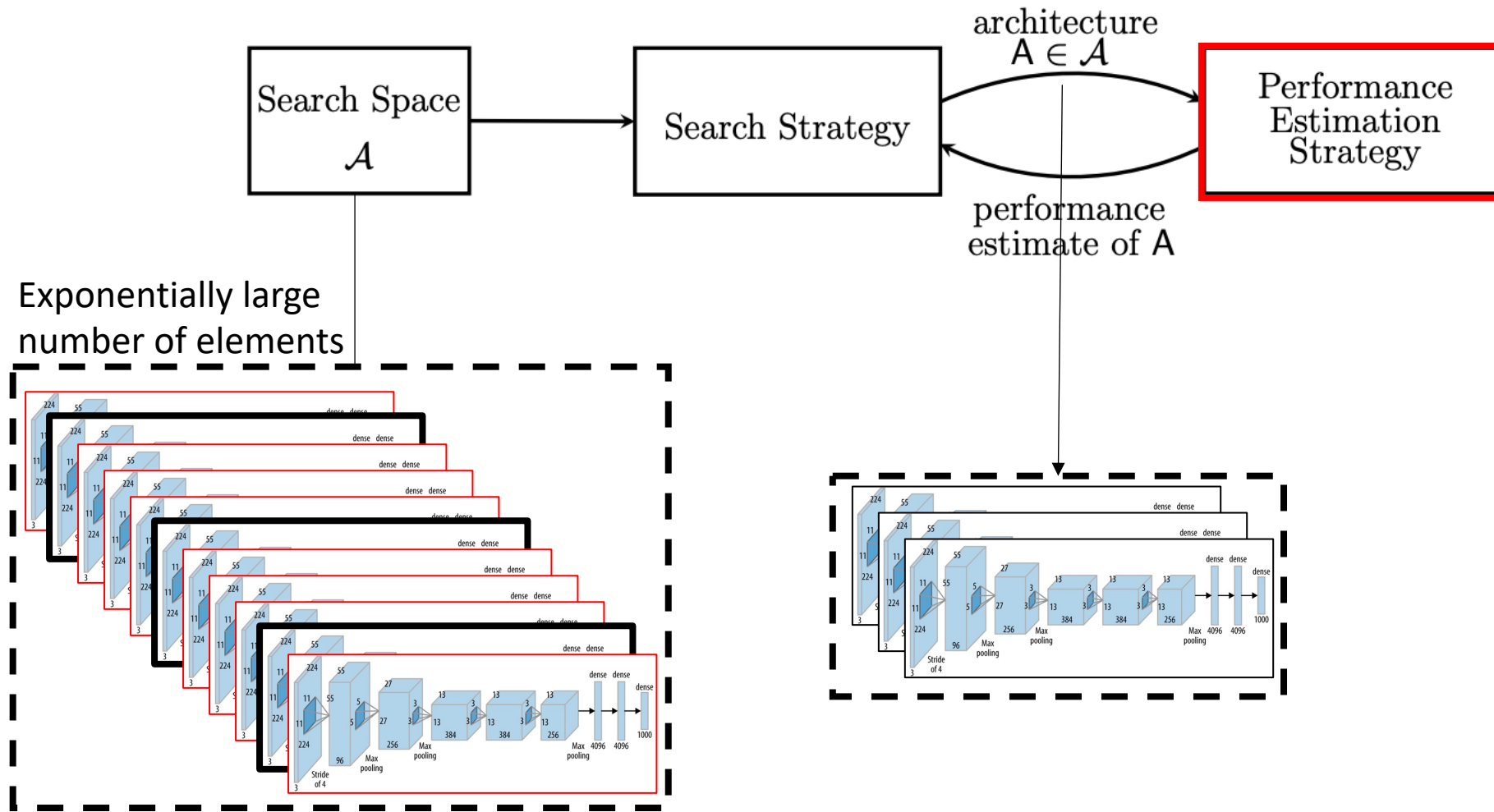https://arxiv.org/pdf/1806.02639.pdf

# RS vs RL

Q: what conclusion can we make?



Cai, Han, Jiacheng Yang, Weinan Zhang, Song Han, and Yong Yu. "Path-level network transformation for efficient architecture search." In *International Conference on Machine Learning*, pp. 678-687. PMLR, 2018.
https://arxiv.org/pdf/1806.02639.pdf

30

# A framework for NAS

# Evaluation Strategy

- Training from scratch (until convergence)

# Evaluation Strategy

- Training from scratch (until convergence)
- Proxy Task Performance

# Evaluation Strategy

- Training from scratch (until convergence)
- Proxy Task Performance
  - Train on a smaller dataset.

# Evaluation Strategy

- Training from scratch (until convergence)
- Proxy Task Performance
  - Train on a smaller dataset.
  - Train for fewer epochs.

# Evaluation Strategy

- Training from scratch (until convergence)
- Proxy Task Performance
  - Train on a smaller dataset.
  - Train for fewer epochs.
  - Train and evaluate a down-scaled model in the search stage
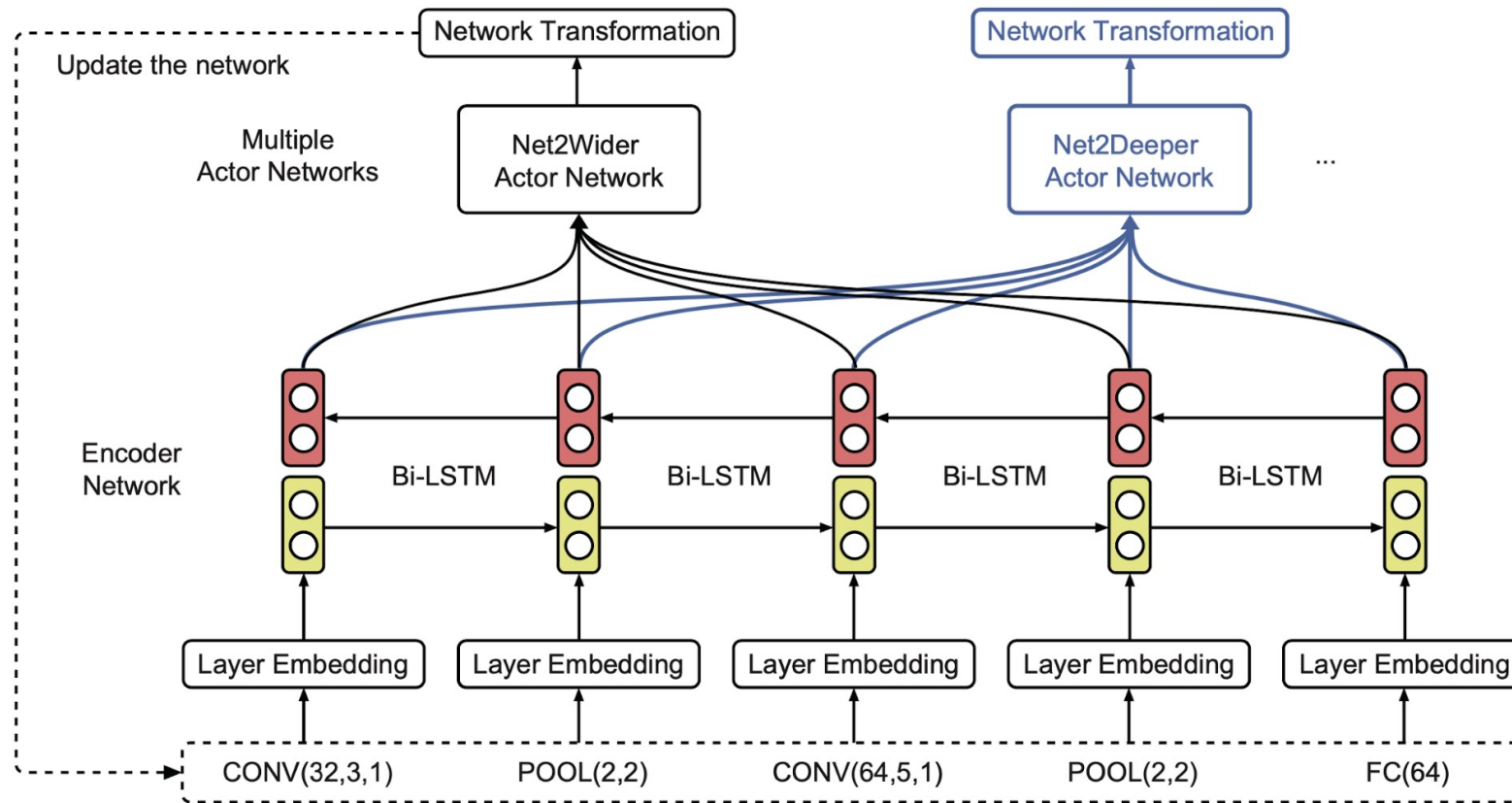
# Evaluation Strategy

- Training from scratch (until convergence)
- Proxy Task Performance
  - Train on a smaller dataset.
  - Train for fewer epochs.
  - Train and evaluate a down-scaled model in the search stage
    - E.g., once a cell structure is learned, the final architecture is determined by selecting a number of cell repeats https://arxiv.org/pdf/1707.07012.pdf

putational budget. After having learned the convolutional cells, several hyper-parameters may be explored to build a final network for a given task: (1) the number of cell repeats $N$ and (2) the number of filters in the initial convolutional cell. After selecting the number of initial filters, we use a

# Evaluation Strategy

- Training from scratch (until convergence)
- Proxy Task Performance
  - Train on a smaller dataset.
  - Train for fewer epochs.
  - Train and evaluate a down-scaled model in the search stage
    - E.g., once a cell structure is learned, the final architecture is determined by selecting a number of cell repeats https://arxiv.org/pdf/1707.07012.pdf
- Parameter Sharing

# Parameter Sharing



Cai, Han, Tianyao Chen, Weinan Zhang, Yong Yu, and Jun Wang. "Efficient architecture search by network transformation."
In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1. 2018.
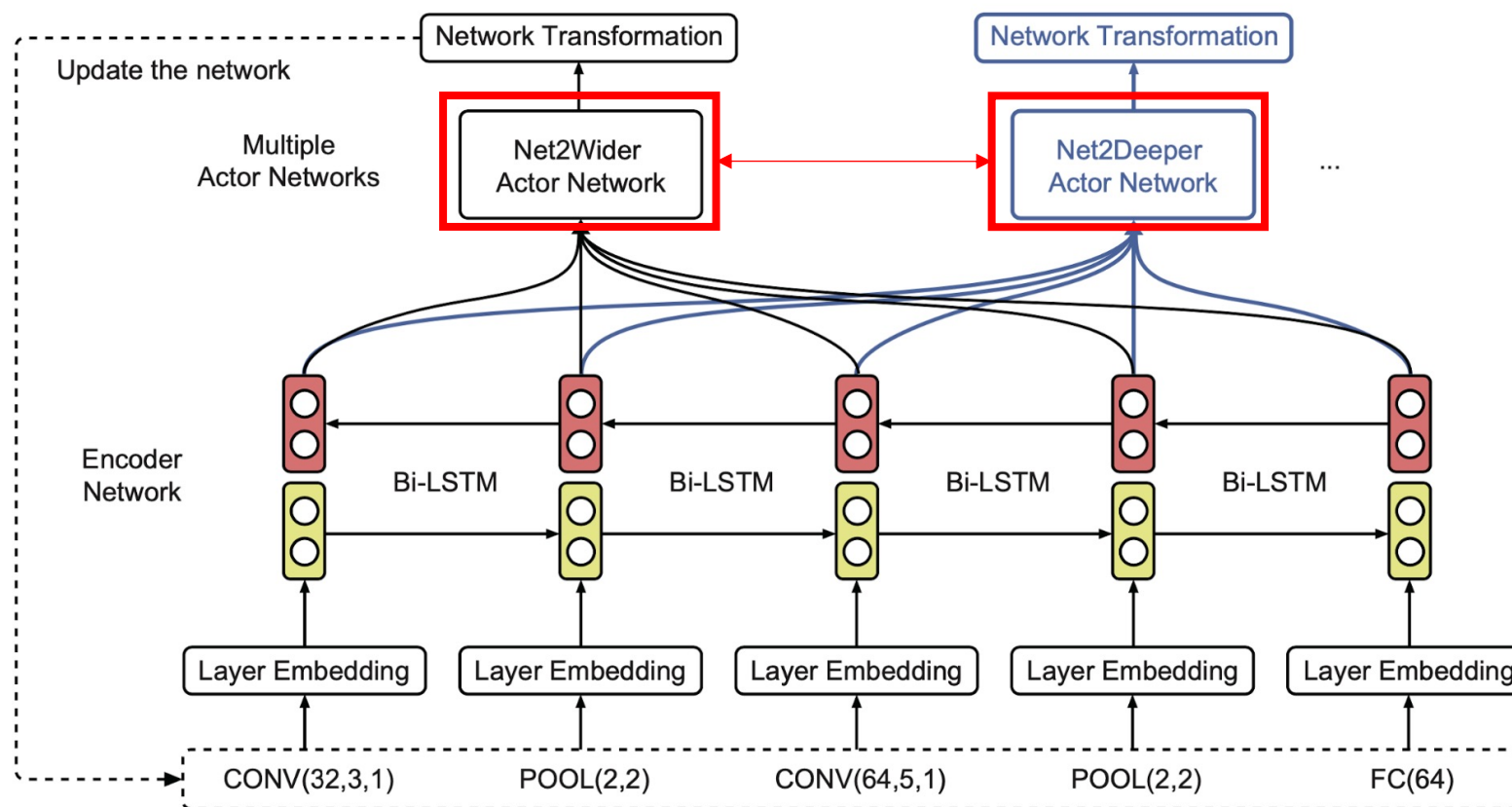https://arxiv.org/pdf/1707.04873.pdf

39

# Parameter Sharing



Cai, Han, Tianyao Chen, Weinan Zhang, Yong Yu, and Jun Wang. "Efficient architecture search by network transformation."
In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1. 2018.
https://arxiv.org/pdf/1707.04873.pdf

40

# Parameter Sharing

more units for fully-connected layers
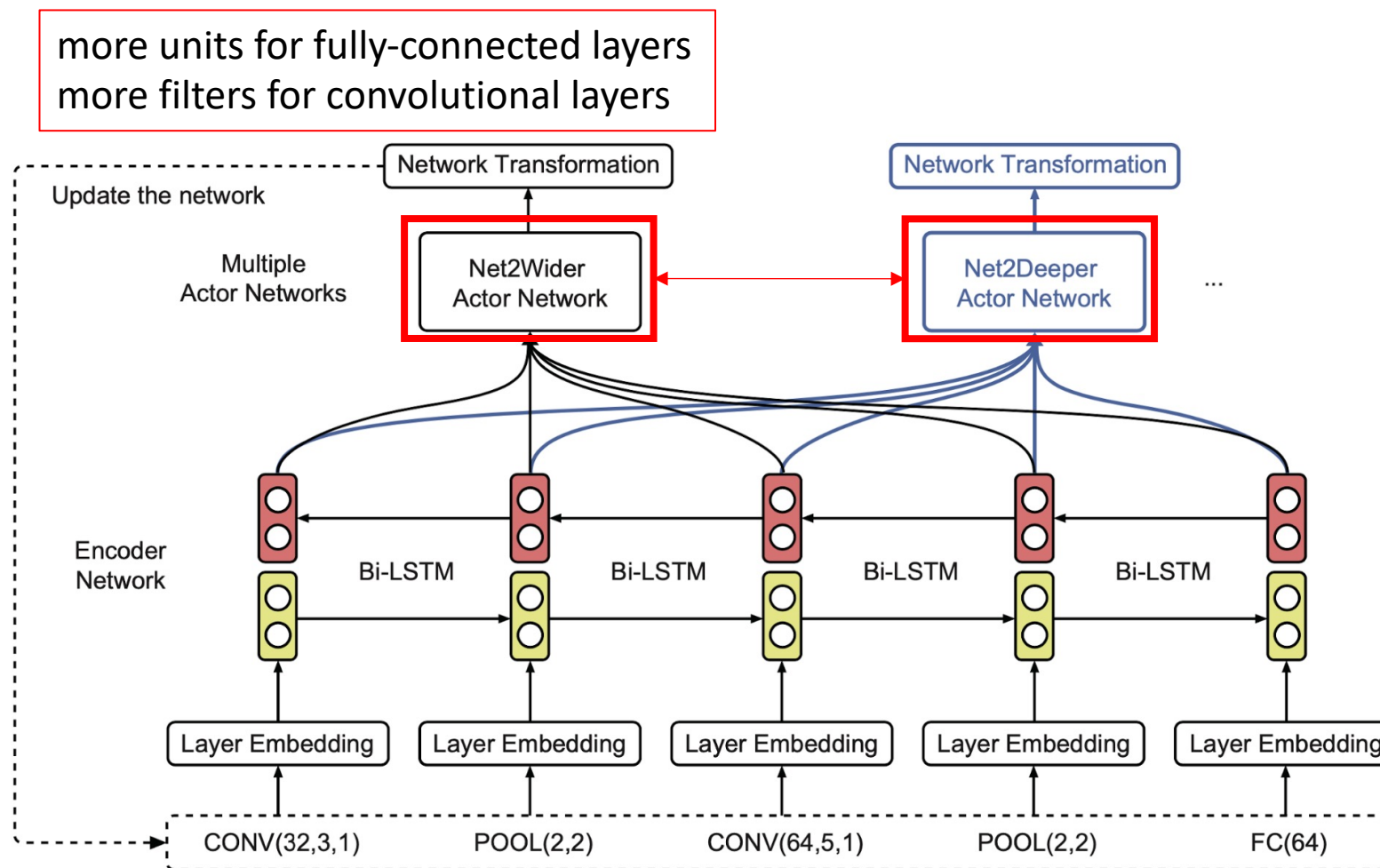more filters for convolutional layers



Cai, Han, Tianyao Chen, Weinan Zhang, Yong Yu, and Jun Wang. "Efficient architecture search by network transformation."
In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1. 2018.
https://arxiv.org/pdf/1707.04873.pdf

41

# Parameter Sharing

more units for fully-connected layers
more filters for convolutional layers

insert a new layer
convolutional layer: the kernel is set to be identity filters
fully-connected layer: the weight matrix is set to be identity matrix



Cai, Han, Tianyao Chen, Weinan Zhang, Yong Yu, and Jun Wang. "Efficient architecture search by network transformation."
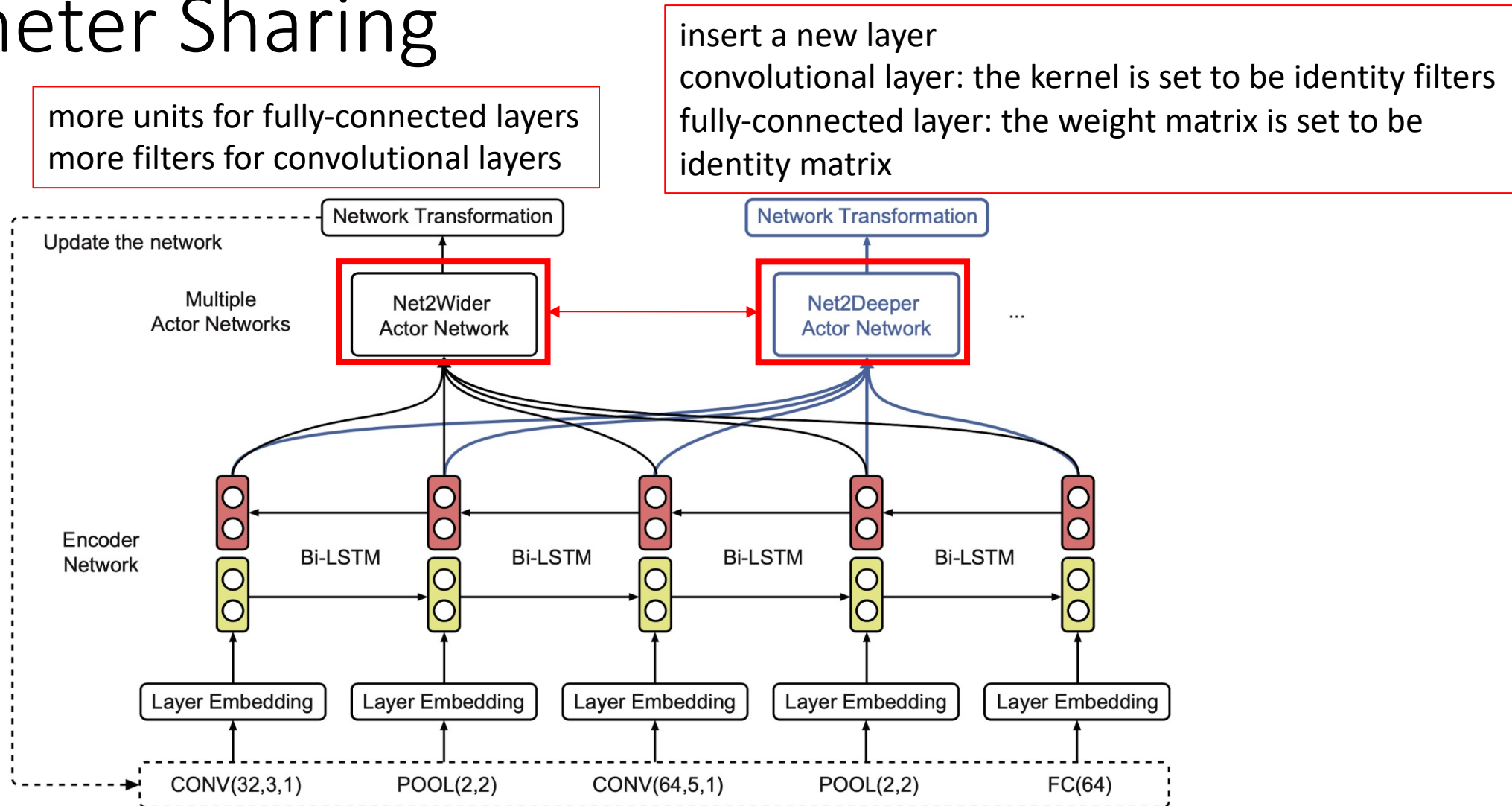In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1. 2018.
https://arxiv.org/pdf/1707.04873.pdf

42

# DARTS: differentiable architecture search

---

**Algorithm 1:** DARTS – Differentiable Architecture Search

---

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$

**while** *not converged* **do**

    1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$

    2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}} \alpha_o^{(i,j)}$ for each edge $(i,j)$

---

$$\min_\alpha \mathcal{L}_{\text{validate}}(w^*(\alpha), \alpha)$$

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

# DARTS: differentiable architecture search

---

**Algorithm 1:** DARTS – Differentiable Architecture Search

---

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$

**while** *not converged* **do**

  1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$
  2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}} \alpha_o^{(i,j)}$ for each edge $(i,j)$

---

$$\min_\alpha \mathcal{L}_{\text{validate}}(w^*(\alpha), \alpha)$$

Architecture

Corresponding weights

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

# DARTS: differentiable architecture search

---

**Algorithm 1:** DARTS – Differentiable Architecture Search

---

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$

**while** *not converged* **do**

    1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$

    2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}} \, \alpha_o^{(i,j)}$ for each edge $(i,j)$

---

$$\min_\alpha \mathcal{L}_{\text{validate}}(w^*(\alpha), \alpha)$$

Architecture

Corresponding weights

$$\text{s.t.} \; w^*(\alpha) = \arg\min_w \mathcal{L}_{\text{train}}(w, \alpha)$$

# DARTS: differentiable architecture search

---

**Algorithm 1:** DARTS – Differentiable Architecture Search

---

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$

**while** *not converged* **do**

    1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$

    2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}} \, \alpha_o^{(i,j)}$ for each edge $(i,j)$

---

$$\min_\alpha \mathcal{L}_{\text{validate}}(w^*(\alpha), \alpha)$$

Architecture

Corresponding weights

$$\text{s.t.} \; w^*(\alpha) = \arg\min_w \mathcal{L}_{\text{train}}(w, \alpha)$$

46

# DARTS: differentiable architecture search

---

**Algorithm 1:** DARTS – Differentiable Architecture Search

---

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$
**while** *not converged* **do**
    1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$
    2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi\nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}} \, \alpha_o^{(i,j)}$ for each edge $(i,j)$

---

$$\min_\alpha \mathcal{L}_{\text{validate}}(w^*(\alpha), \alpha)$$

Architecture

Corresponding weights

$$\text{s.t.} \; w^*(\alpha) = \arg\min_w \mathcal{L}_{\text{train}}(w, \alpha)$$

47

# DARTS: differentiable architecture search

---

**Algorithm 1:** DARTS – Differentiable Architecture Search

---

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$

**while** *not converged* **do**

    1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$

    2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}} \alpha_o^{(i,j)}$ for each edge $(i,j)$

---

$$\min_\alpha \mathcal{L}_{\text{validate}}(w^*(\alpha), \alpha)$$

Architecture

Corresponding weights

$$\text{s.t.} \quad w^*(\alpha) = \arg\min_w \mathcal{L}_{\text{train}}(w, \alpha)$$

# DARTS: differentiable architecture search

**Algorithm 1:** DARTS – Differentiable Architecture Search

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$

**while** *not converged* **do**

    1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$

    2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}} \alpha_o^{(i,j)}$ for each edge $(i,j)$



(a) Initially unknown operations on the edges.

(b) Continuous relaxation by placing a mixture of operations on each edge.

(c) Bilevel optimization to jointly train mixing probabilities and weights.

(d) Finalized the model based on the learned mixing probabilities.

# DARTS: differentiable architecture search

**Algorithm 1:** DARTS – Differentiable Architecture Search

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$
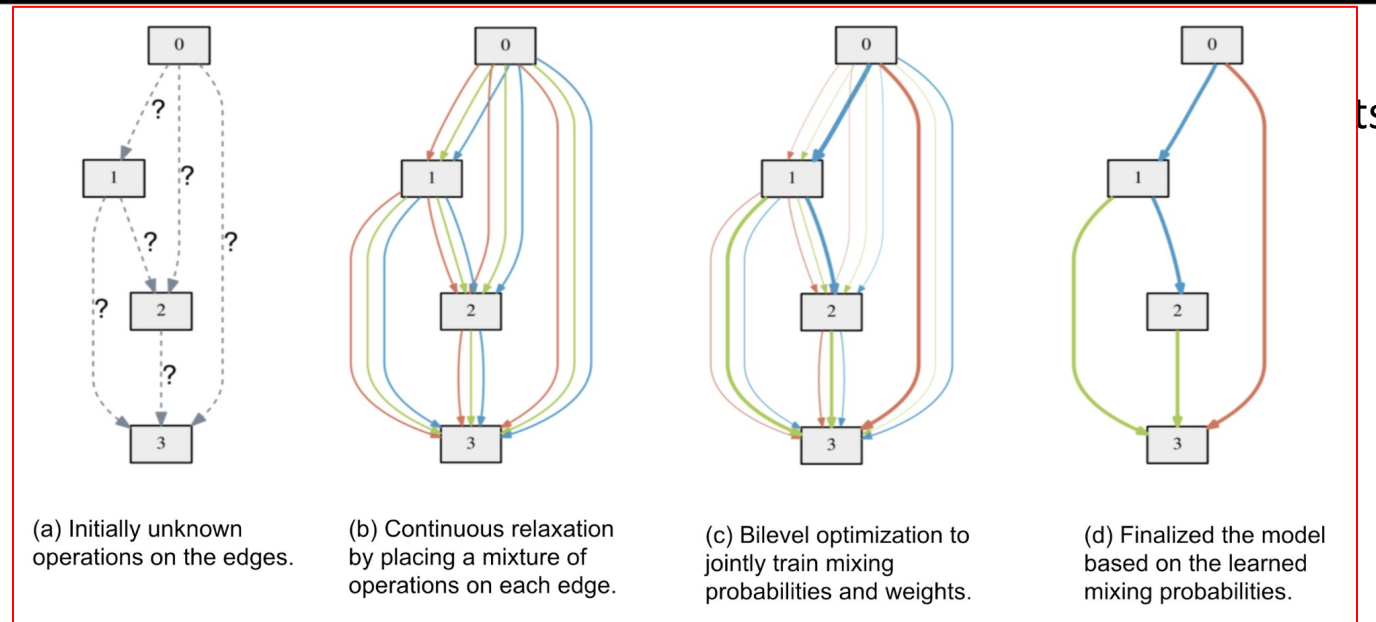
**while** *not converged* **do**

    1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$
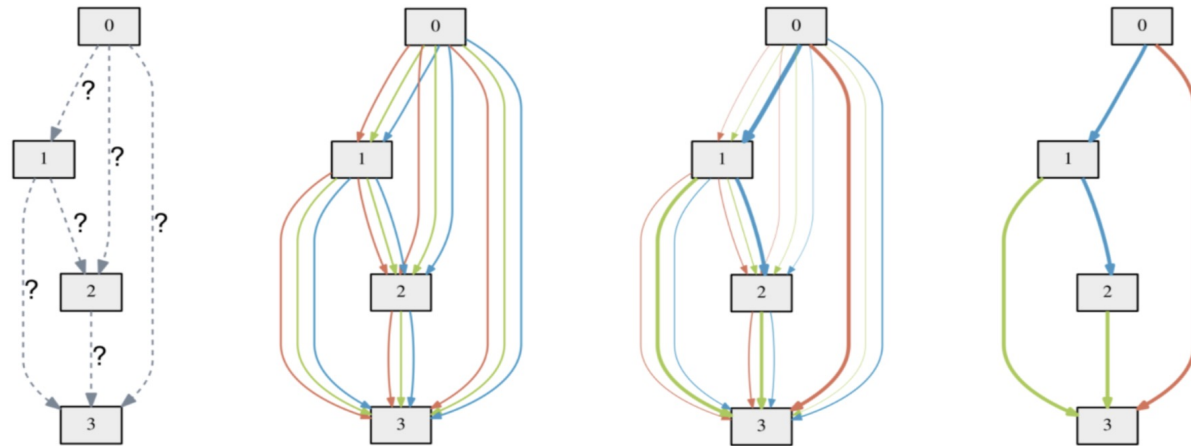
    2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi\nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}}\ \alpha_o^{(i,j)}$ for each edge $(i,j)$

$$x_i = \sum_{j<i} o^{(i,j)}(x_j)$$

$$\bar{o}^{(i,j)}(x) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_{ij}^o)}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{ij}^{o'})} o(x)$$



(a) Initially unknown operations on the edges.

(b) Continuous relaxation by placing a mixture of operations on each edge.

(c) Bilevel optimization to jointly train mixing probabilities and weights.

(d) Finalized the model based on the learned mixing probabilities.

50

# DARTS: differentiable architecture search

**Algorithm 1:** DARTS – Differentiable Architecture Search

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$

**while** *not converged* **do**

    1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$

    2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

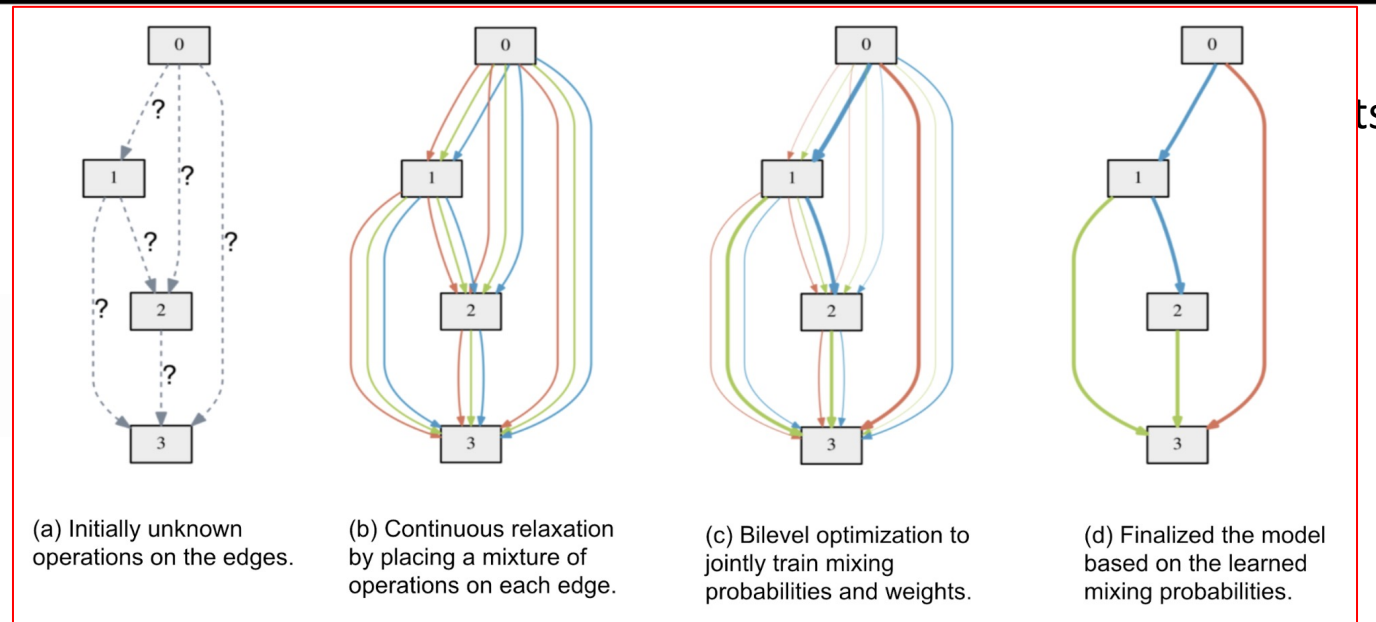Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}} \, \alpha_o^{(i,j)}$ for each edge $(i,j)$

$$x_i = \sum_{j<i} o^{(i,j)}(x_j)$$

Mix of operations

$$\bar{o}^{(i,j)}(x) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_{ij}^o)}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{ij}^{o'})} o(x)$$

differentiable

ts



(a) Initially unknown operations on the edges.

(b) Continuous relaxation by placing a mixture of operations on each edge.

(c) Bilevel optimization to jointly train mixing probabilities and weights.

(d) Finalized the model based on the learned mixing probabilities.

51

# DARTS: differentiable architecture search

**Algorithm 1:** DARTS – Differentiable Architecture Search

Create a mixed operation $\bar{o}^{(i,j)}$ parametrized by $\alpha^{(i,j)}$ for each edge $(i,j)$

**while** *not converged* **do**

1. Update weights $w$ by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$
2. Update architecture $\alpha$ by descending $\nabla_\alpha \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$

Replace $\bar{o}^{(i,j)}$ with $o^{(i,j)} = \text{argmax}_{o \in \mathcal{O}} \alpha_o^{(i,j)}$ for each edge $(i,j)$

$$x_i = \sum_{j < i} o^{(i,j)}(x_j)$$

Mix of operations

$$\bar{o}^{(i,j)}(x) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_{ij}^o)}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{ij}^{o'})} o(x)$$

differentiable



(a) Initially unknown operations on the edges.

(b) Continuous relaxation by placing a mixture of operations on each edge.

(c) Bilevel optimization to jointly train mixing probabilities and weights.

(d) Finalized the model based on the learned mixing probabilities.

52