

# NAT: Neural Architecture Transformer for Accurate and Compact Architectures

Yong Guo, Yin Zheng, Mingkui Tan, Qi Chen,  
Jian Chen, Peilin Zhao, Junzhou Huang

Published in NeurIPS 2019

# Contents

1. Background
2. Proposed Method
3. Experimental Results
4. Conclusion

# Contents

1. Background

2. Proposed Method

3. Experimental Results

4. Conclusion

# Background

Deep neural networks have achieved great success in many computer vision tasks, such as [image classification](#), [face recognition](#), [object detection](#), etc.

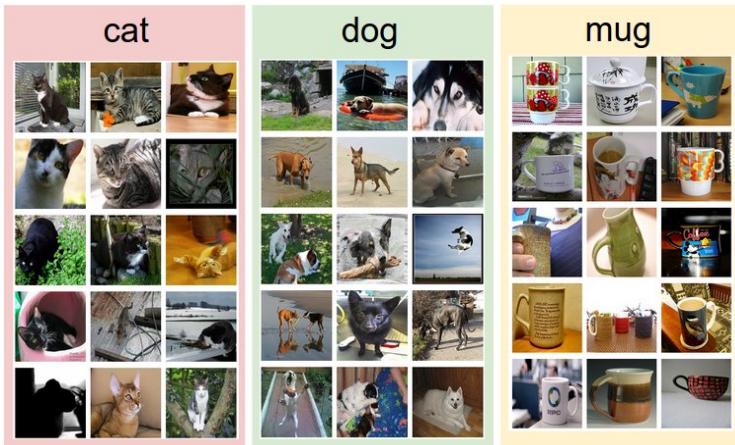
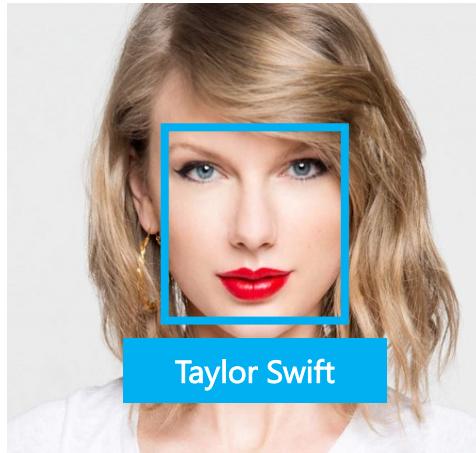
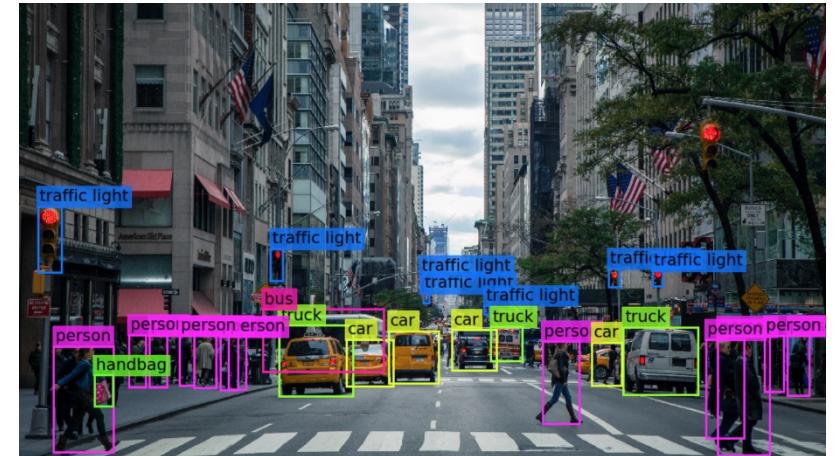


Image Classification



Face Recognition



Object Detection

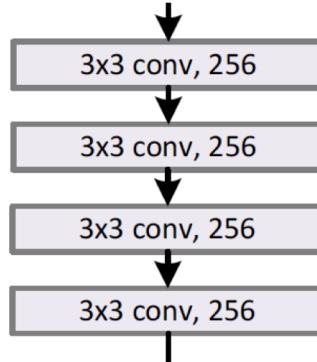
Figure: Applications of deep neural networks.

# Neural Architecture Design

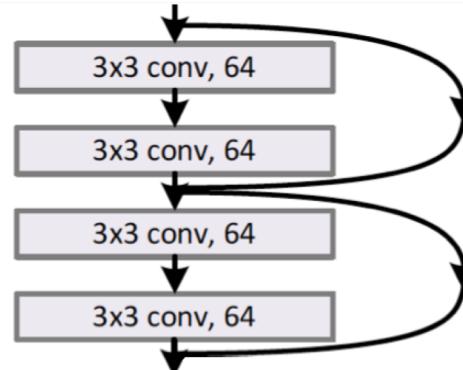
- Neural architecture design is one of the key factors behind the success of deep neural networks.
- Existing architectures can be divided into two categories:
  1. Hand-crafted architectures
  2. Automatically searched architectures

# Hand-crafted Architectures

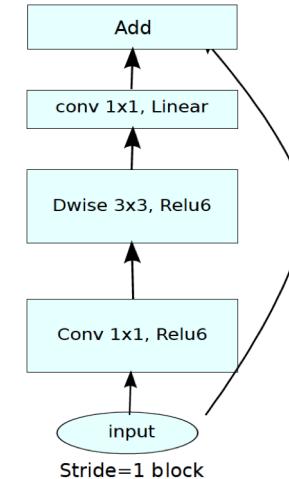
Several widely used hand-crafted architectures:



VGG



ResNet



MobileNetV2

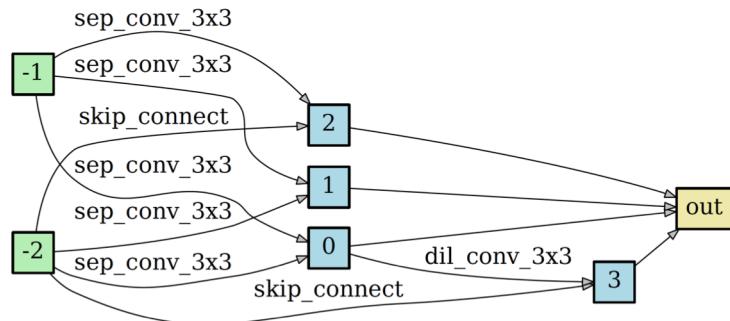
Figure: Examples of hand-crafted architectures.

## Limitations of hand-crafted architecture design process

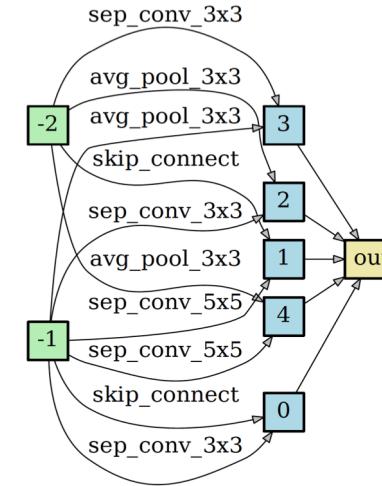
- Hand-crafted methods rely on **substantial human expertise**.
- Hand-crafted methods cannot fully **explore the whole architecture space**.

# Automatically Searched Architectures

- There is a growing interest to **automate the manual process** of architecture design by Neural Architecture Search (NAS).



DARTS normal cell



ENAS normal cell

Figure: Examples of NAS based architectures.

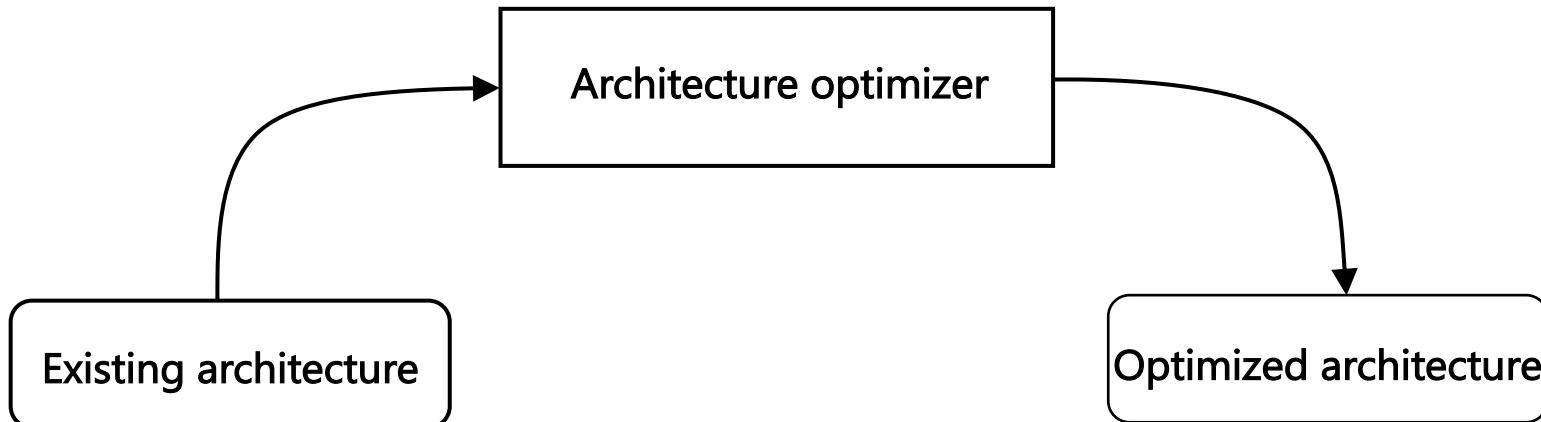
## Limitations of NAS methods

- Search space is extremely large, *e.g.*, billions of candidate architectures.
- NAS methods may find suboptimal architectures with limited performance.

# Architecture Optimization

Since both the hand-crafted and NAS based architectures are not optimal, can we **optimize architectures to obtain the better ones?**

- One can design architecture optimization methods to optimize existing architectures for **better performance**.



**Figure:** Architecture optimization scheme.

# Existing Architecture Optimization Methods

## ■ Neural Architecture Optimization (NAO)

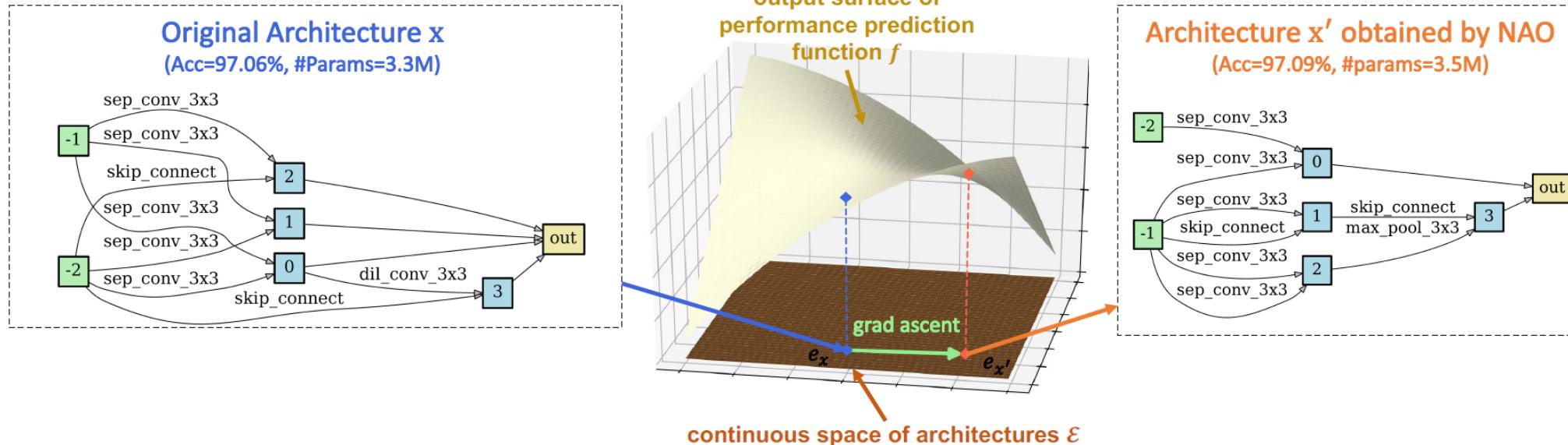


Figure: Framework of Neural Architecture Optimization (NAO).

## Limitations of NAO

- NAO often produces a **totally different architecture** from the input one.
- NAO may **introduce extra parameters** or **additional computational cost**.

# Contents

1. Background
2. Proposed Method
3. Experimental Results
4. Conclusion

# Motivation

- Both hand-crafted architectures and NAS based architectures may contain non-significant or redundant operations.
- Existing architecture optimization methods may introduce extra parameters or additional computational cost into the architectures.

How to transform the redundant operations in **any arbitrary architecture** to improve the performance without introducing extra computational cost?

# Problem Definition

**Our goal:** Transforming any arbitrary architecture for **better performance and less computational cost**.

**One solution:** Replacing the redundant operations with the **more efficient ones**.

We divide the operations into three categories  $\{S, N, O\}$ .  $S$  denotes **skip connection**,  $N$  denotes **null connection**,  $O$  denotes the **other operations**.

- We have  $c(O) > c(S) > c(N)$ , where  $c(\cdot)$  is a function to evaluate the computational cost.
- To **reduce the computational cost**, we allow the transitions:  $O \rightarrow S$ ,  $O \rightarrow N$ ,  $S \rightarrow N$ .
- Since skip connection has negligible cost but often can significantly improve the performance, we also allow  $N \rightarrow S$ .

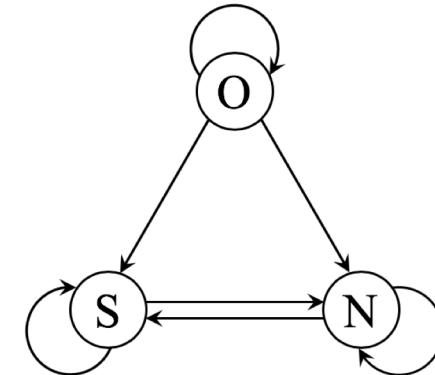


Figure: Operation transformation scheme.

# Optimization for Single Architecture

Given a specific architecture  $\hat{\beta}$ , we seek to find the optimal architecture  $\alpha$ . Then, the optimization problem can be formulated as

$$\max_{\alpha} R(\alpha | \hat{\beta}), \text{ s.t. } c(\alpha) \leq \kappa$$

- $R(\alpha | \hat{\beta}) = R(\alpha, w_\alpha) - R(\hat{\beta}, w_{\hat{\beta}})$  denotes the performance difference between the optimized architectures  $\alpha$  and the given architectures  $\hat{\beta}$ .  $w_\alpha$  and  $w_{\hat{\beta}}$  are the parameters of  $\alpha$  and  $\hat{\beta}$ .
- $c(\cdot)$  is a function to measure the computation cost of architectures.
- $\kappa$  is an upper bound of the computational cost.

# Optimization for Arbitrary Architecture

- The previous optimization problem is only concentrated on one specific  $\hat{\beta}$ .
- We hope to find a transformer to optimize any arbitrary architecture  $\beta \sim p(\cdot)$ .

To this end, we seek to train the transformer parameterized  $\theta$  by optimizing the objective

$$\max_{\theta} \mathbb{E}_{\beta \sim p(\cdot)} [\mathbb{E}_{\alpha \sim \pi(\cdot | \beta; \theta)} R(\alpha | \beta)], \text{ s.t. } c(\alpha) \leq \kappa, \alpha \sim \pi(\cdot | \beta; \theta)$$

where  $\mathbb{E}_{\beta \sim p(\cdot)} [\mathbb{E}_{\alpha \sim \pi(\cdot | \beta; \theta)} R(\alpha | \beta)]$  denotes the expectation of  $R(\alpha | \beta)$  over the distribution of  $\beta \sim p(\cdot)$  and the distribution of  $\alpha \sim \pi(\cdot | \beta; \theta)$ .

# Optimization for Arbitrary Architecture

$$\max_{\theta} \underbrace{\mathbb{E}_{\beta \sim p(\cdot)} [\mathbb{E}_{\alpha \sim \pi(\cdot | \beta; \theta)} R(\alpha | \beta)]}_{\text{s.t. } c(\alpha) \leq \kappa, \alpha \sim \pi(\cdot | \beta; \theta)}$$

## Several challenges regarding the optimization problem

- It is hard to find a comprehensive measure to accurately evaluate the cost.
- The upper bound of computational cost  $\kappa$  is hard to determine.
- How to compute  $\mathbb{E}_{\beta \sim p(\cdot)} [\mathbb{E}_{\alpha \sim \pi(\cdot | \beta; \theta)} R(\alpha | \beta)]$  still remains a question.

# Markov Decision Process for Learning NAT

## Our solution

- We cast the optimization problem into a Markov decision process (MDP).
- We seek to make a series of decisions to replace redundant operations with the more computationally efficient operations.

**Benefits:** We do not have to evaluate the cost  $c(\alpha)$  or determine the upper bound  $\kappa$  to obtain an architecture with less computational cost.

# Markov Decision Process for Learning NAT

## Details of MDP

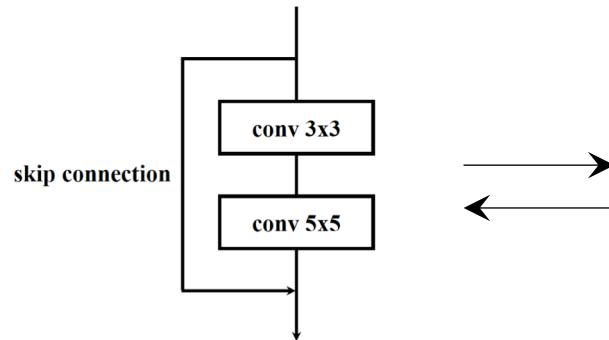
- An architecture is defined as a **state**.
- A transformation mapping  $\beta \rightarrow \alpha$  is defined as an **action**.
- The accuracy improvement on validation set is regarded as **reward**.
- The policy  $\pi(\cdot | \beta; \theta)$  parameterized by  $\theta$  is the **probability distribution of the action**.

Based on MDP, how to build a model to learn the **optimal policy**  $\pi$  ?

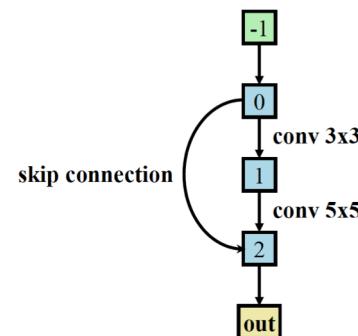
# Policy Learning by Graph Convolution Networks

**Graph Representation of Architectures:** We use the [directed acyclic graph \(DAG\)](#) to represent both hand-crafted architectures and NAS based architectures.

- Node: feature maps of a specific layer in deep networks
- Edge: a computational module or operation, e.g., convolution or max pooling



(a) network view



(b) graph view

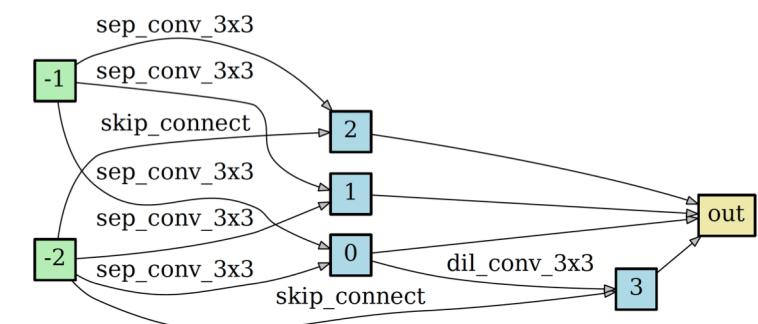


Figure: Graph representation examples of NAS based architectures.

Figure: Graph representation examples of hand-crafted architectures.

# Policy Learning by Graph Convolution Networks

To better exploit the **adjacency information** of the operations in an architecture, we use a two-layer **graph convolutional network (GCN)** to build the controller.

$$\mathbf{Z} = f(\mathbf{X}, \mathbf{A}) = \text{Softmax} \left( \mathbf{A} \sigma \left( \mathbf{A} \mathbf{X} \mathbf{W}^{(0)} \right) \mathbf{W}^{(1)} \mathbf{W}^{\text{FC}} \right)$$

## Notations

- $\mathbf{A}$  : the adjacency matrix of the architecture graph.
- $\mathbf{X}$  : the attributes of the nodes together with their two input edges in the graph.
- $\mathbf{W}^{(0)}$  and  $\mathbf{W}^{(1)}$  : the weights of two graph convolution layers.
- $\mathbf{W}^{\text{FC}}$  : the weight of the fully-connected layer,  $\sigma$  is a non-linear activation function.
- $\mathbf{Z}$  : the probability distribution of different candidate operations, *i.e.*, the learned policy.

# Training Method

We train the transformer parameters  $\theta$  and the model parameter  $w$  in an **alternative** way.

- Training the model parameters  $w$  :

$$w \leftarrow w - \eta \frac{1}{m} \sum_{i=1}^m \nabla_w \mathcal{L}(\beta_i, w)$$

where  $\mathcal{L}(\cdot)$  is the cross-entropy loss,  $\eta$  is the learning rate.

- Training the transformer parameters  $\theta$  :

To encourage exploration, we introduce an **entropy regularization term**:

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\beta \sim p(\cdot)} \left[ \mathbb{E}_{\alpha \sim \pi(\cdot | \beta; \theta)} [R(\alpha, w) - R(\beta, w)] + \underline{\lambda H(\pi(\cdot | \beta; \theta))} \right] \\ &= \sum_{\beta} p(\beta) \left[ \sum_{\alpha} \pi(\alpha | \beta; \theta) (R(\alpha, w) - R(\beta, w)) + \underline{\lambda H(\pi(\cdot | \beta; \theta))} \right] \end{aligned}$$

where  $H(\cdot)$  evaluates the entropy of the policy, and  $\lambda$  controls the strength of the entropy regularization term.

# Training Method

---

**Algorithm 1** Training method for Neural Architecture Transformer (NAT).

**Require:** The number of sampled input architectures in an iteration  $m$ , the number of sampled optimized architectures for each input architecture  $n$ , learning rate  $\eta$ , regularizer parameter  $\lambda$ , input architecture distribution  $p(\cdot)$ , shared model parameters  $w$ , transformer parameters  $\theta$ .

```
1: Initiate  $w$  and  $\theta$ .
2: while not convergent do
3:   for each iteration on training data do
4:     // Fix  $\theta$  and update  $w$ .
5:     Sample  $\beta_i \sim p(\cdot)$  to construct a batch  $\{\beta_i\}_{i=1}^m$ .
6:     Update the model parameters  $w$  by descending the gradient:
7:        $w \leftarrow w - \eta \frac{1}{m} \sum_{i=1}^m \nabla_w \mathcal{L}(\beta_i, w)$ .
8:   end for
9:   for each iteration on validation data do
10:    // Fix  $w$  and update  $\theta$ .
11:    Sample  $\beta_i \sim p(\cdot)$  to construct a batch  $\{\beta_i\}_{i=1}^m$ .
12:    Obtain  $\{\alpha_j\}_{j=1}^n$  according to the policy learned by GCN.
13:    Update the parameters  $\theta$  by descending the gradient:
14:       $\theta \leftarrow \theta - \eta \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n [\nabla_\theta \log \pi(\alpha_j | \beta_i; \theta) (R(\alpha_j, w) - R(\beta_i, w)) + \lambda \nabla_\theta H(\pi(\cdot | \beta_i; \theta))]$ .
15:  end for
16: end while
```

# Contents

1. Background

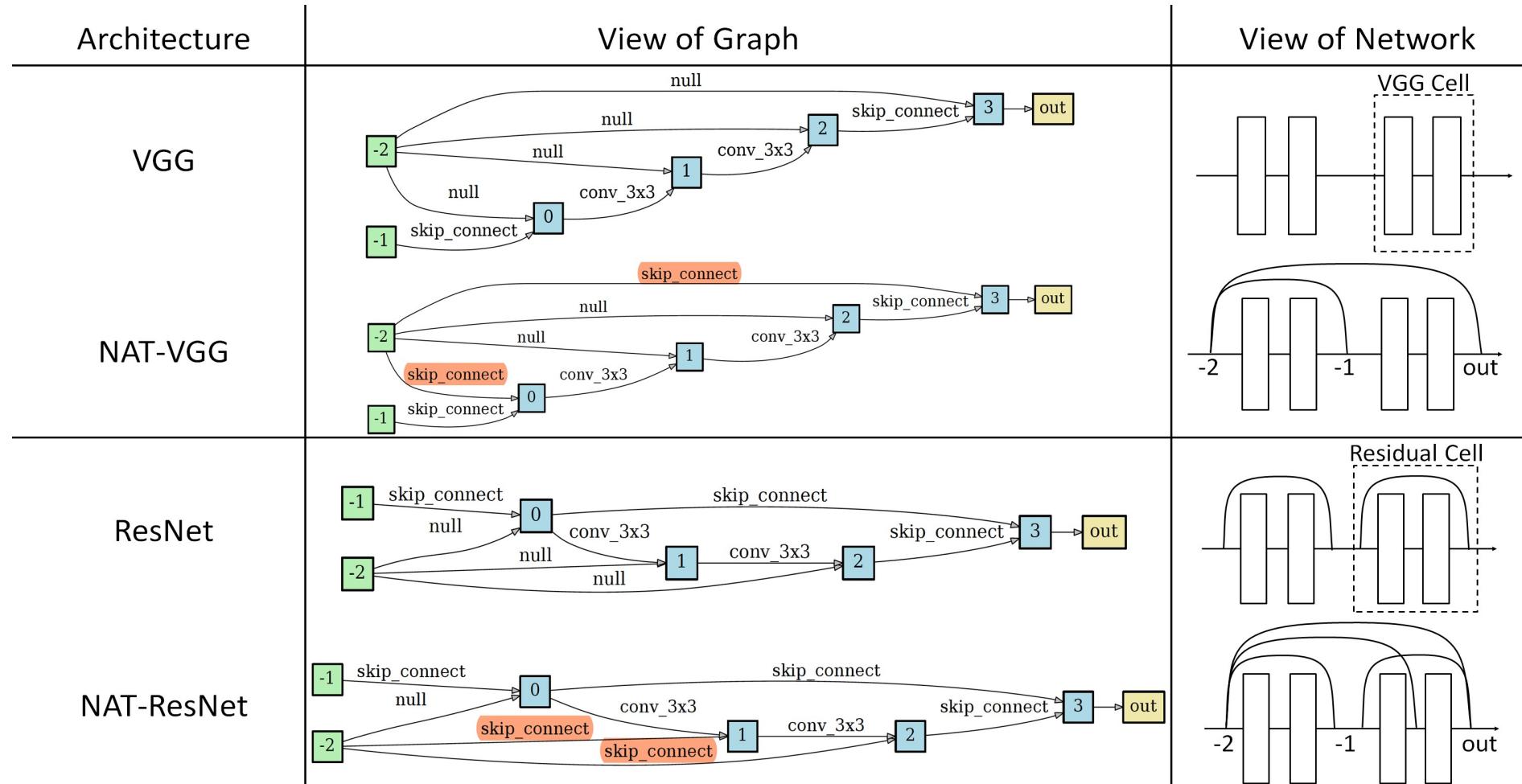
2. Proposed Method

3. Experimental Results

4. Conclusion

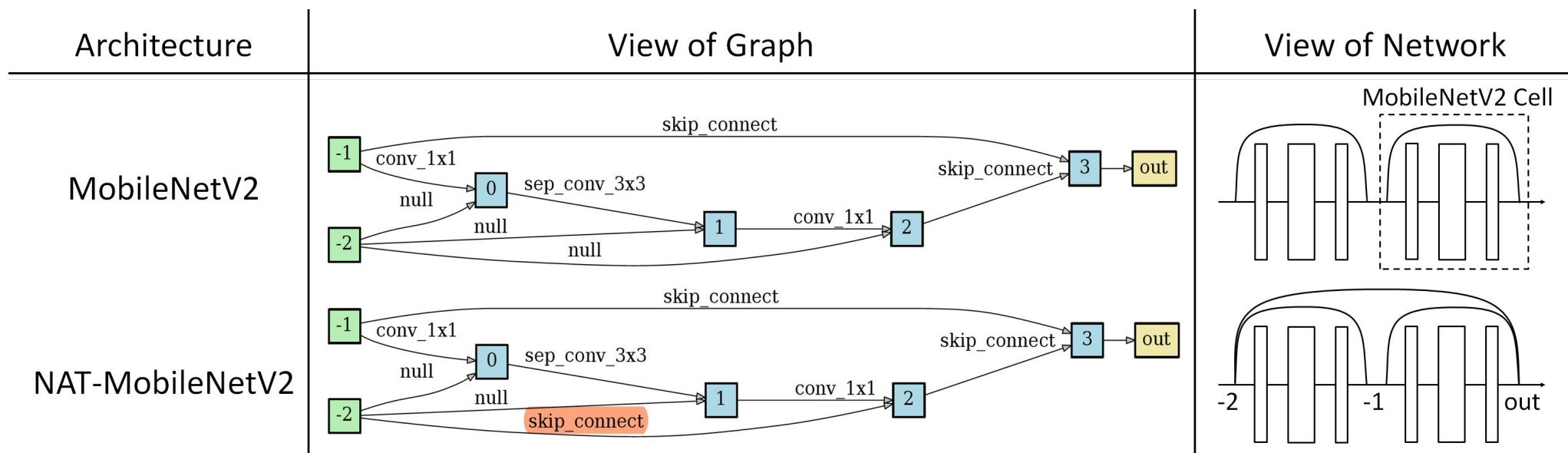
# Visual Results of Hand-crafted Architectures

- Results on several hand-crafted architectures, including VGG, ResNet, and MobileNet.



# Visual Results of Hand-crafted Architectures

- Results on several hand-crafted architectures, including VGG, ResNet, and MobileNet.



- NAT introduces additional skip connections to improve the performance.

# Comparison on Hand-crafted Architectures

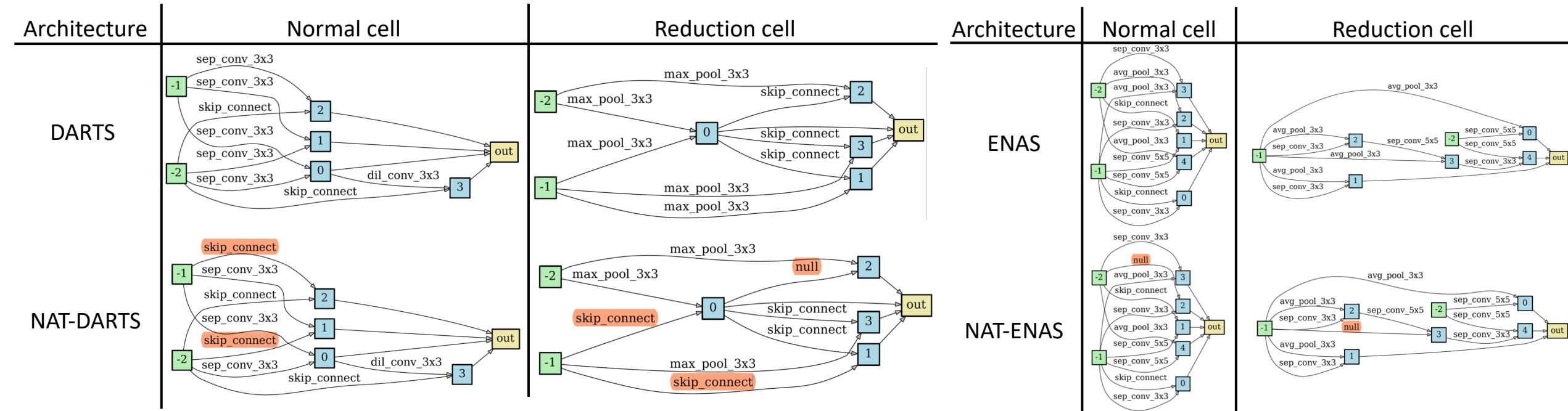
- Results on several hand-crafted architectures, including VGG, ResNet, and MobileNet.

CIFAR-10					ImageNet					
Model	Method	#Params (M)	#MAdds (M)	Acc. (%)	Model	Method	#Params (M)	#MAdds (M)	Acc. (%)	
									Top-1	Top-5
VGG16	/	15.2	313	93.56	VGG16	/	138.4	15620	71.6	90.4
	NAO[31]	19.5	548	95.72		NAO [31]	147.7	18896	72.9	91.3
	NAT	15.2	313	<b>96.04</b>		NAT	138.4	15620	<b>74.3</b>	<b>92.0</b>
ResNet20	/	0.3	41	91.37	ResNet18	/	11.7	1580	69.8	89.1
	NAO [31]	0.4	61	92.44		NAO [31]	17.9	2246	70.8	89.7
	NAT	0.3	41	<b>92.95</b>		NAT	11.7	1580	<b>71.1</b>	<b>90.0</b>
ResNet56	/	0.9	127	93.21	ResNet50	/	25.6	3530	76.2	92.9
	NAO [31]	1.3	199	95.27		NAO [31]	34.8	4505	77.4	93.2
	NAT	0.9	127	<b>95.40</b>		NAT	25.6	3530	<b>77.7</b>	<b>93.5</b>
MobileNetV2	/	2.3	91	94.47	MobileNetV2	/	3.4	300	72.0	90.3
	NAO [31]	2.9	131	94.75		NAO [31]	4.5	513	72.2	90.6
	NAT	2.3	91	<b>95.17</b>		NAT	3.4	300	<b>72.5</b>	<b>91.0</b>

- NAT based models yield **significantly better performance** with approximately **the same computational cost** as the baseline models.

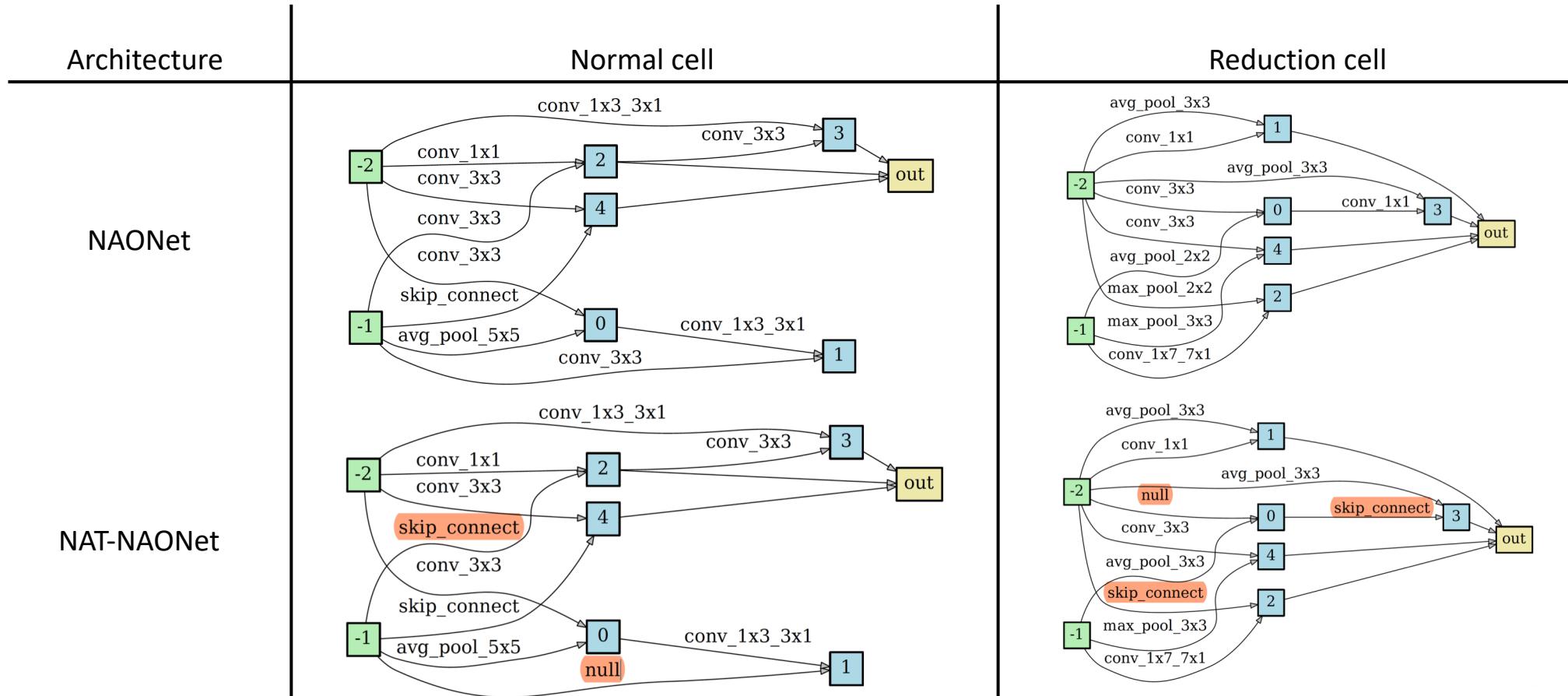
# Visual Results on NAS based Architectures

- Results on several NAS based architectures, including ENAS, DARTS, and NAONet.



# Visual Results on NAS based Architectures

- Results on several NAS based architectures, including ENAS, DARTS, and NAONet.



# Comparison on NAS based Architectures

- Results on several NAS based architectures, including ENAS, DARTS, and NAONet.

CIFAR-10					ImageNet				
Model	Method	#Params (M)	#MAdds (M)	Acc. (%)	Model	Method	#Params (M)	#MAdds (M)	Acc. (%)
AmoebaNet <sup>†</sup> [34]	/	3.2	-	96.73	AmoebaNet [34]	/	5.1	555	74.5 92.0
PNAS <sup>†</sup> [28]	/	3.2	-	96.67	PNAS [28]	/	5.1	588	74.2 91.9
SNAS <sup>†</sup> [48]	/	2.9	-	97.08	SNAS [48]	/	4.3	522	72.7 90.8
GHN <sup>†</sup> [52]	/	5.7	-	97.22	GHN [52]	/	6.1	569	73.0 91.3
ENAS <sup>†</sup> [33]	/	4.6	804	97.11	ENAS [33]	/	5.6	679	73.8 91.7
	NAO [31]	4.5	763	97.05		NAO [31]	5.5	656	73.7 91.7
	NAT	4.6	804	<b>97.24</b>		NAT	5.6	679	<b>73.9</b> <b>91.8</b>
DARTS <sup>†</sup> [29]	/	3.3	533	97.06	DARTS [29]	/	5.9	595	73.1 91.0
	NAO [31]	3.5	577	97.09		NAO [31]	6.1	627	73.3 91.1
	NAT	3.0	483	<b>97.28</b>		NAT	3.9	515	<b>74.4</b> <b>92.2</b>
NAONet <sup>†</sup> [31]	/	128	66016	97.89	NAONet [31]	/	11.35	1360	74.3 91.8
	NAO [31]	143	73705	97.91		NAO [31]	11.83	1417	74.5 92.0
	NAT	113	58326	<b>98.01</b>		NAT	8.36	1025	<b>74.8</b> <b>92.3</b>

- NAT based models yield significantly better performance with less or comparable computational cost as the baseline models.

# Comparison of Different Policy Learners

- We compare several policy learners, including Random Search, LSTM, and two GCN based methods.

Method	VGG16	ResNet20	MobileNetV2	ENAS <sup>†</sup>	DARTS <sup>†</sup>	NAONet <sup>†</sup>
/	93.56	91.37	94.47	97.11	97.06	97.89
Random Search	93.17	91.56	94.38	96.58	95.17	96.31
LSTM	94.45	92.19	95.01	97.05	97.05	97.93
Maximum-GCN	94.37	92.57	94.87	96.92	97.00	97.90
Sampling-GCN (Ours)	<b>95.93</b>	<b>92.97</b>	<b>95.13</b>	<b>97.21</b>	<b>97.26</b>	<b>97.99</b>

- Our Sampling-GCN method significantly outperforms the other methods.

# Contents

1. Background
2. Proposed Method
3. Experimental Results
4. Conclusion

# Conclusion

- We propose a novel Neural Architecture Transformers (NAT) to optimize any arbitrary architectures for better performance without extra computational cost.
- We cast the problem into a Markov decision process (MDP) and employ graph convolutional network (GCN) to learn the optimal policy.
- Extensive experiments show the effectiveness of NAT on both hand-crafted and NAS based architectures.

Thanks!  
Q & A