

# NAT: Neural Architecture Transformer for Accurate and Compact Architectures

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





### BACKGROUND AND MOTIVATION

Limitations of existing architecture design methods:

- Hand-crafted architecture design methods: rely on substantial human expertise and cannot fully explore the whole architecture design space.
- Neural architecture search (NAS) methods: often produce subopimal architectures due to the extremely large search space.

Both hand-crafted architectures and NAS based architectures may contain non-significant or redundant operations.

#### CONTRIBUTIONS

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

# PROBLEM DEFINITION

- Operations can be categorized into  $\{S, N, O\}$ .
- . S denotes skip connection
- 2. N denotes null connection
- The computational cost follows:

3. O denotes other connections



• We only allow the transitions:  $O \rightarrow S, O \rightarrow N, S \rightarrow N, S \rightarrow N$ .

Goal: Transform any arbitrary architecture for better performance and less computational cost.

Solution: Replace redundant operations with more efficient ones, such as S and N.

# OPTIMIZATION FOR SINGLE ARCHITECTURE

Given a specific architecture  $\widehat{\beta}$ , we seek to find an optimal architecture  $\alpha$ . The optimization problem can be formulated as:

$$\max_{\alpha} R\left(\alpha | \widehat{\beta}\right), \text{ s.t. } c(\alpha) \leq \kappa. \tag{1}$$

- $R(\alpha|\beta) = R(\alpha, w_{\alpha}) R(\beta, w_{\beta})$  denotes the performance difference between the optimized architectures  $\alpha$  and the given architectures  $\beta$ .  $w_{\alpha}$  and  $w_{\beta}$  are the parameters of  $\alpha$  and  $\beta$ , respectively.
- ullet  $c(\cdot)$  measures the computation cost of an architecture.
- $\kappa$  is an upper bound of the cost.

#### OPTIMIZATION FOR ARBITRARY ARCHITECTURE

To optimize arbitrary architecture, we solve the following problem:

$$\max_{\theta} \mathbb{E}_{\beta \sim p(\cdot)} \left[ \mathbb{E}_{\alpha \sim \pi(\cdot|\beta;\theta)} R(\alpha|\beta) \right], \text{ s.t. } c(\alpha) \leq \kappa, \ \alpha \sim \pi(\cdot|\beta;\theta). \tag{2}$$

We further cast the problem into an Markov Decision Process.

- An architecture is defined as a state.
- A transformation mapping  $\beta \to \alpha$  is defined as an action.
- The accuracy improvement on validation set is regraded as reward.

# POLICY LEARNING BY GCN

To better exploit the adjacency information of architecture graph, we use a two-layer graph convolution network to build the controller:

$$\mathbf{Z} = f(\mathbf{X}, \mathbf{A}) = \text{Softmax}\left(\mathbf{A}\sigma\left(\mathbf{A}\mathbf{X}\mathbf{W}^{(0)}\right)\mathbf{W}^{(1)}\mathbf{W}^{\text{FC}}\right).$$
 (3)

- A: adjacency matrix of the architecture graph.
- X: attributes of the nodes in the graph.
- $W^{(0)}$  and  $W^{(1)}$ : weights of graph convolution layers.
- W<sup>FC</sup>: weight of the fully-connected layer.
- $\sigma$ : non-linear activation function.
- $\bullet$  **Z**: probability distribution of different operations, *i.e.*, the learned policy  $\pi(\cdot|\beta;\theta)$ .

# TRAINING AND INFERENCE METHOD

Training method for NAT

Algorithm 1 Training method for Neural Architecture Transformer.

- Initiate w and  $\theta$ .
- : while not convergent do
- for each iteration on training data do
- Sample  $\beta_i \sim p(\cdot)$  to construct a batch  $\{\beta_i\}_{i=1}^m$ . Update the model parameters w by descending the gradient.
- for each iteration on validation data do
- Sample  $\beta_i \sim p(\cdot)$  to construct a batch  $\{\beta_i\}_{i=1}^m$ .
- Obtain  $\{\alpha_j\}_{j=1}^n$  according to the policy learned by GCN.
- Update the parameters  $\theta$  by descending the gradient.
- end while
- Inferring the optimized architectures
  - 1. Sample candidate architectures from the learned policy  $\pi(\cdot|\beta;\theta)$ .
  - 2. Select the architectures with the highest validation accuracy.

## RESULTS ON DIFFERENT ARCHITECTURES

• Results on hand-crafted architectures (comparisons on ImageNet)

VGG16     NAO     147.7     18896     72.9     91.3       NAT     138.4     15620     74.3     92.0       NAT     138.4     15620     74.3     92.0       ResNet18     NAO     17.9     2246     70.8     89.7       NAT     11.7     1580     71.1     90.0       ResNet50     NAO     34.8     4505     77.4     93.2       NAT     25.6     3530     77.7     93.5       MobileNetV2     NAO     4.5     513     72.2     90.6       NAT     3.4     300     72.5     91.0	Model	Method	#Params (M)	#MAdds (M)	Acc. (%)	
VGG16     NAO     147.7     18896     72.9     91.3       NAT     138.4     15620 <b>74.3 92.0</b> Image: ResNet18     Image: NAO     11.7     1580     69.8     89.1       ResNet18     NAO     17.9     2246     70.8     89.7       NAT     11.7     1580 <b>71.1 90.0</b> ResNet50     NAO     34.8     4505     77.4     93.2       NAT     25.6     3530 <b>77.7 93.5</b> Image: Note of the content of the co					Top-1	Top-5
NAT     138.4     15620     74.3     92.0       Image: ResNet18 NAO ResNet18 NAO NAT	VGG16	/	138.4	15620	71.6	90.4
ResNet18   / NAO NAT   11.7 11.7 1580 170.8 89.7 11.7 1580 11.7 1580 11.7 1580 11.7 1580 11.7 1580 11.1 11.7 1580 11.1 11.7 1580 11.1 1580 11.1 11.7 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1 1580 11.1		NAO	147.7	18896	72.9	91.3
ResNet18   NAO   17.9   2246   70.8   89.7     NAT   11.7   1580 <b>71.1 90.0</b> ResNet50   NAO   34.8   4505   77.4   93.2     NAT   25.6   3530 <b>77.7 93.5</b> NobileNetV2   NAO   4.5   513   72.2   90.6		NAT	138.4	15620	<b>74.3</b>	92.0
NAT   11.7   1580   71.1   90.0     Image: ResNet50 ResNet50 NAO NAO NAT	ResNet18	/	11.7	1580	69.8	89.1
/   25.6   3530   76.2   92.9     ResNet50   NAO   34.8   4505   77.4   93.2     NAT   25.6   3530   77.7   93.5     /   3.4   300   72.0   90.3     MobileNetV2   NAO   4.5   513   72.2   90.6		NAO	17.9	2246	70.8	89.7
ResNet50   NAO   34.8   4505   77.4   93.2     NAT   25.6   3530   77.7   93.5     J   3.4   300   72.0   90.3     MobileNetV2   NAO   4.5   513   72.2   90.6		NAT	11.7	1580	<b>71.1</b>	90.0
NAT   25.6   3530   77.7   93.5     /   3.4   300   72.0   90.3     MobileNetV2   NAO   4.5   513   72.2   90.6	ResNet50	/	25.6	3530	76.2	92.9
/ 3.4 300 72.0 90.3   MobileNetV2 NAO 4.5 513 72.2 90.6		NAO	34.8	4505	77.4	93.2
MobileNetV2 NAO 4.5 513 72.2 90.6		NAT	25.6	3530	77.7	93.5
	MobileNetV2	/	3.4	300	72.0	90.3
NAT 3.4 300 <b>72.5 91.0</b>		NAO	4.5	513	72.2	90.6
		NAT	3.4	300	72.5	91.0

Results on NAS based architectures (comparisons on ImageNet)

$\mathbf{N}I \circ A \circ 1$	Method	#Params (M)	#MAdds (M)	Acc. (%)	
Model				Top-1	Top-5
AmoebaNet [34]		5.1	555	74.5	92.0
PNAS [28]		5.1	588	74.2	91.9
SNAS [48]		4.3	522	72.7	90.8
GHN [52]		6.1	569	73.0	91.3
	/	5.6	679	73.8	91.7
ENAS [33]	NAO	5.5	656	73.7	91.7
	NAT	5.6	679	<b>73.9</b>	91.8
	/	5.9	595	73.1	91.0
<b>DARTS</b> [29]	NAO	6.1	627	73.3	91.1
	NAT	3.9	515	<b>74.4</b>	92.2
	/	11.35	1360	74.3	91.8
NAONet [31]	NAO	11.83	1417	74.5	92.0
	NAT	8.36	1025	<b>74.8</b>	92.3

### VISUALIZATION OF ARCHITECTURES

• Architecture optimization results on hand-crafted architectures

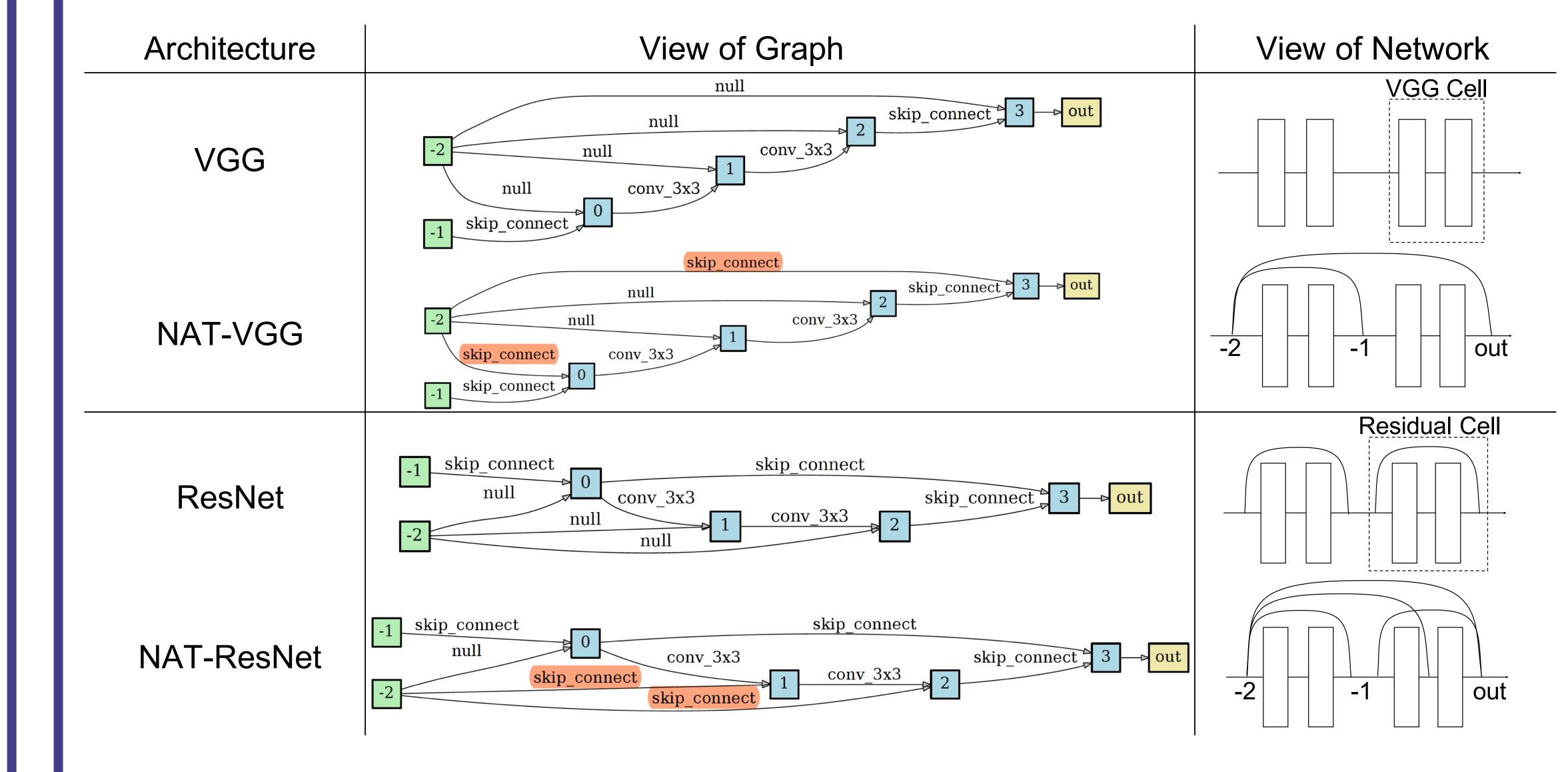


Figure 1: Visualization of some optimized hand-crafted architectures

• Architecture optimization results on NAS based architectures

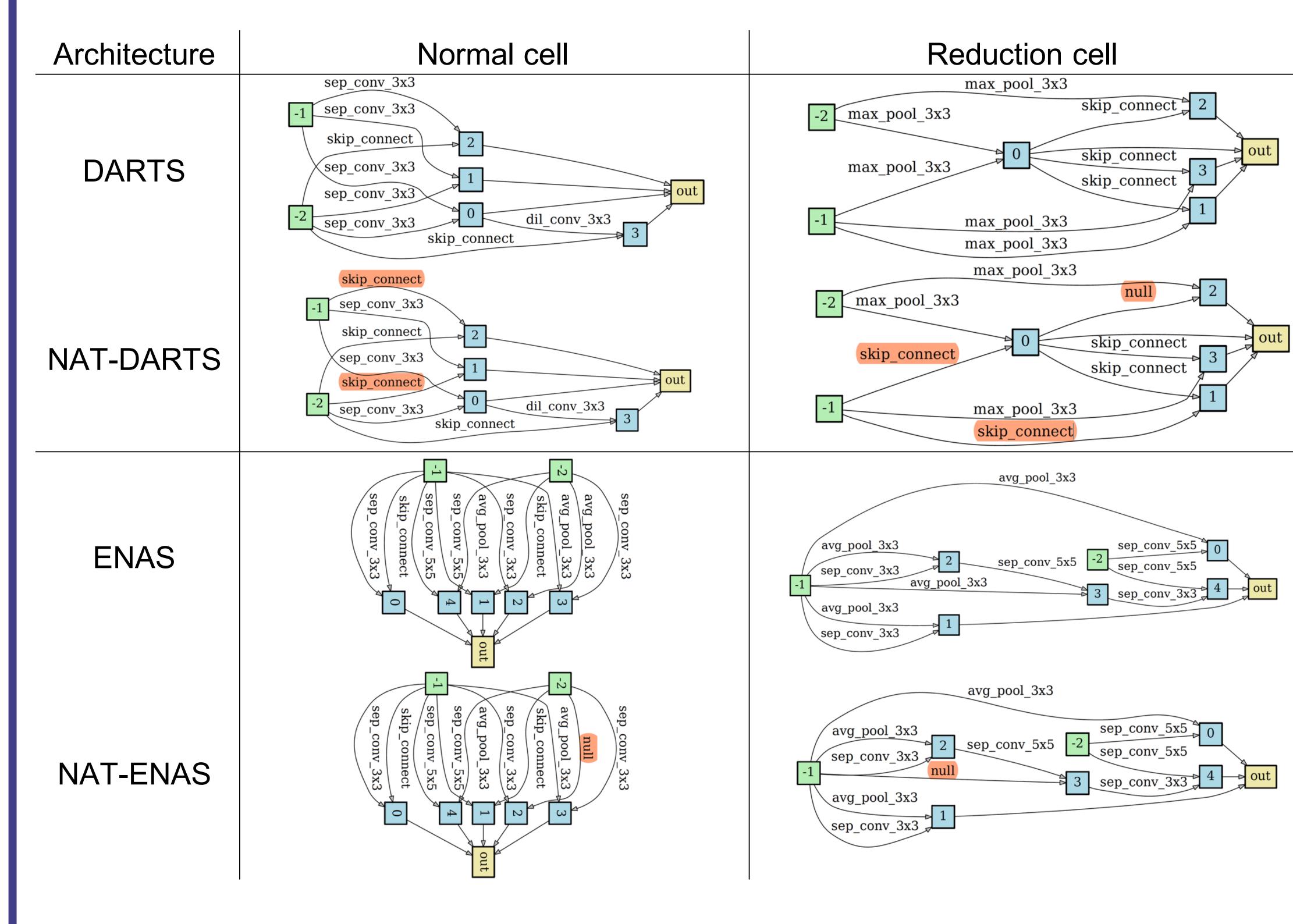


Figure 2: Visualization of some optimized NAS-based architectures

# CONTACT INFORMATION AND CODE

- Email: mingkuitan@scut.edu.cn
- Code: https://github.com/guoyongcs/NAT

