

# CATE: Computation-aware Neural Architecture Encoding with Transformers

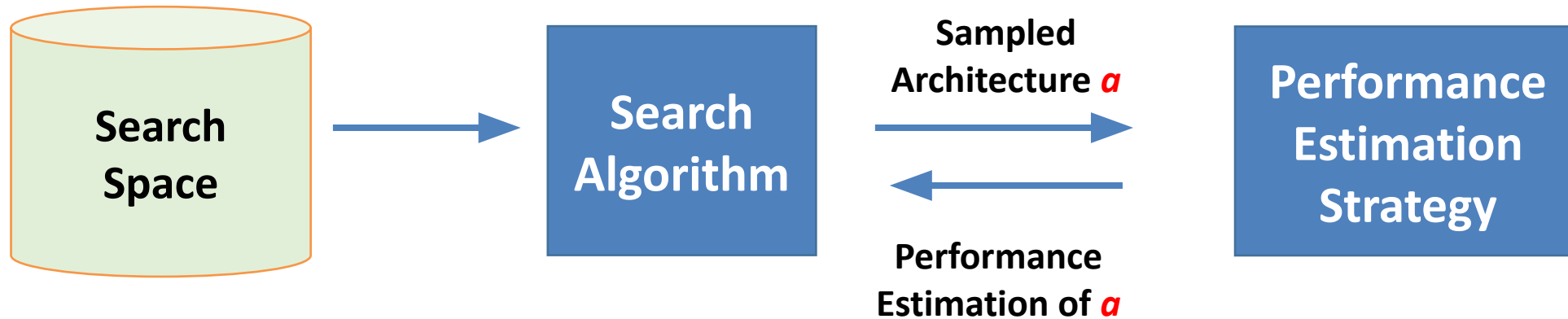
Shen Yan, Kaiqiang Song, Fei Liu, Mi Zhang

Michigan State University, Tencent AI Lab, University of Central Florida

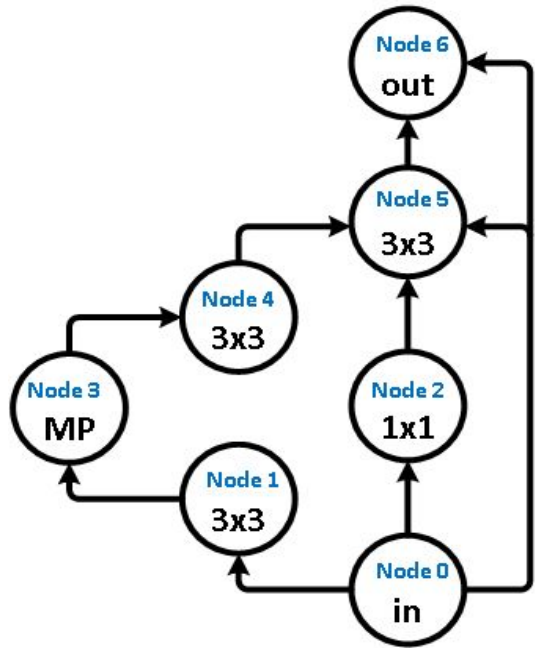
<https://arxiv.org/abs/2102.07108>

June 20<sup>th</sup>, 2021

# Neural Architecture Search (NAS) Pipeline



# Structure-aware encodings: Adjacency matrix-based

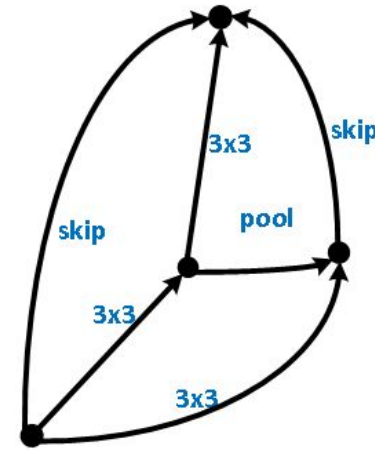


One-hot adjacency encoding of NAS-Bench-101

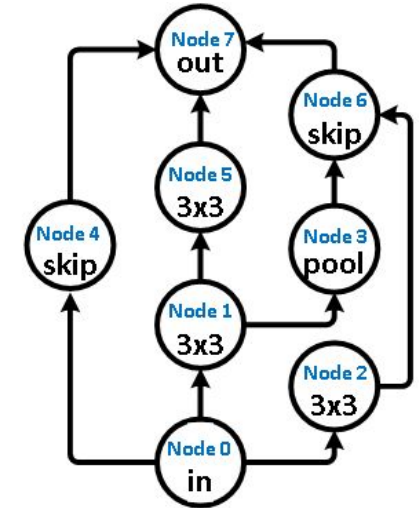
		Node						
		[0	1	2	3	4	5	6]
Node	[0	0	1	1	0	0	1	1
	1	0	0	0	1	0	0	0
	2	0	0	0	0	0	1	0
	3	0	0	0	0	1	0	0
	4	0	0	0	0	0	1	0
	5	0	0	0	0	0	0	1
	6	0	0	0	0	0	0	0

		Operation				
		[in	1x1	3x3	MP	out]
Node	[0	1	0	0	0	0
	1	0	0	1	0	0
	2	0	1	0	0	0
	3	0	0	0	1	0
	4	0	0	1	0	0
	5	0	0	1	0	0
	6	0	0	0	0	1



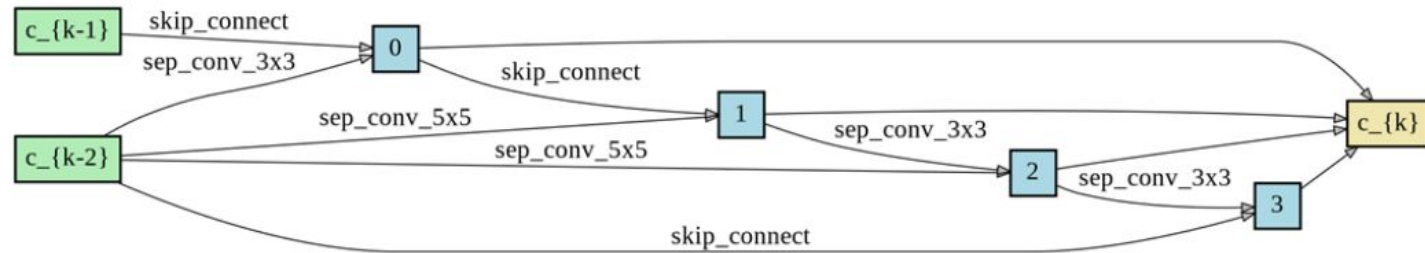
		Node							
		[0	1	2	3	4	5	6	7]
Node	[0	0	1	1	0	1	0	0	0
	1	0	0	0	1	0	1	0	0
	2	0	0	0	0	0	0	1	0
	3	0	0	0	0	0	0	1	0
	4	0	0	0	0	0	0	0	1
	5	0	0	0	0	0	0	0	1
	6	0	0	0	0	0	0	0	1
	7	0	0	0	0	0	0	0	0



		Operation							
		[in	1x1	3x3	pool	skip	none	out]	
Node	[0	1	0	0	0	0	0	0	
	1	0	0	1	0	0	0	0	
	2	0	0	1	0	0	0	0	
	3	0	0	0	1	0	0	0	
	4	0	0	0	0	1	0	0	
	5	0	0	1	0	0	0	0	
	6	0	0	0	0	1	0	0	
	7	0	0	0	0	0	0	1	

One-hot adjacency encoding of NAS-Bench-201

# Structure-aware encodings: Adjacency matrix-based



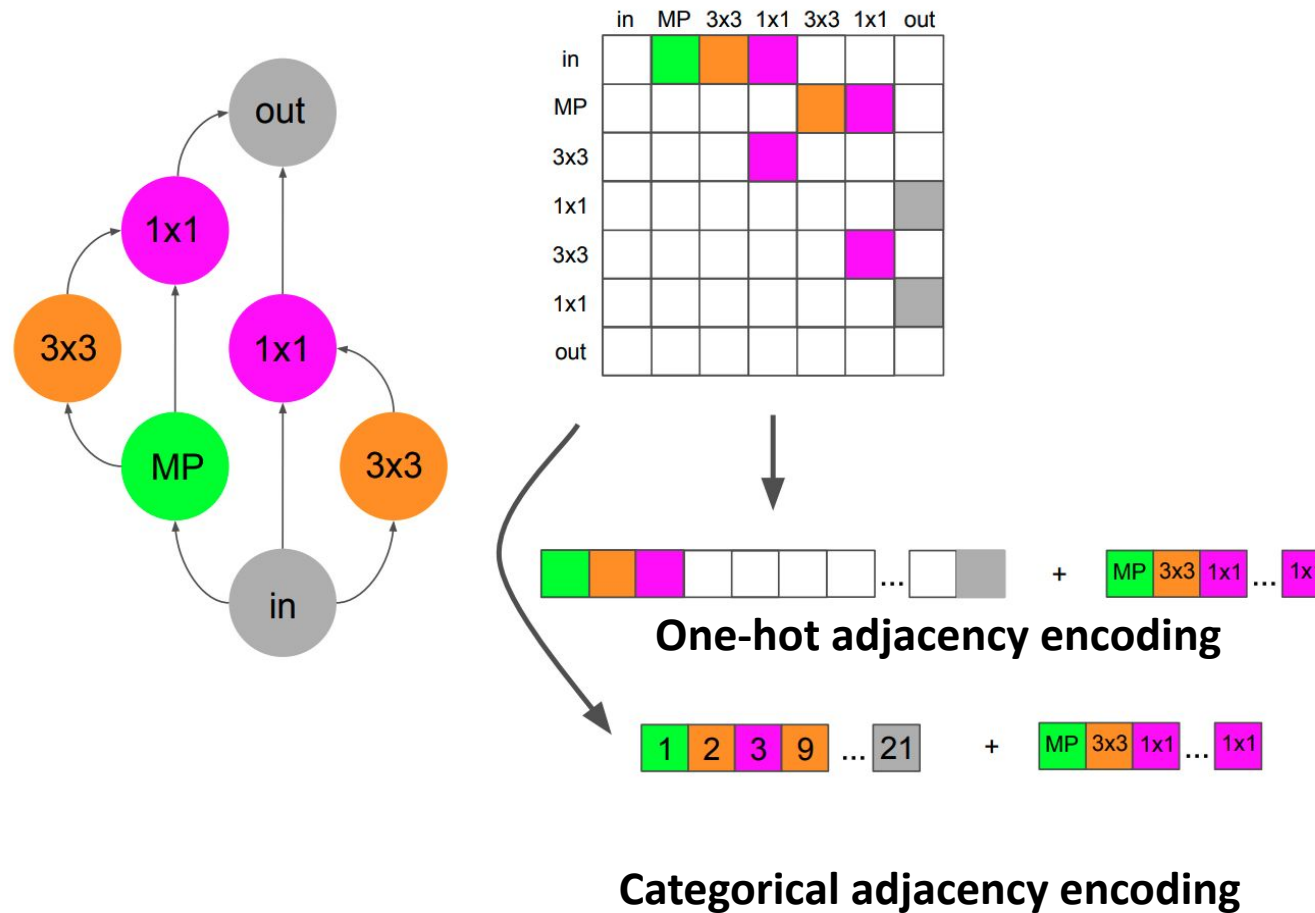
	$c_{k-2}$	$c_{k-1}$	node0	node1	node2	node3	$c_k$
$c_{k-2}$	0	0	1	0	1	0	0
$c_{k-1}$	0	0	0	1	0	0	0
node0	0	0	0	0	1	0	0
node1	0	0	0	0	0	1	0
node2	0	0	0	0	0	0	1
node3	0	0	0	0	0	0	0
$c_k$	0	0	0	0	0	0	0

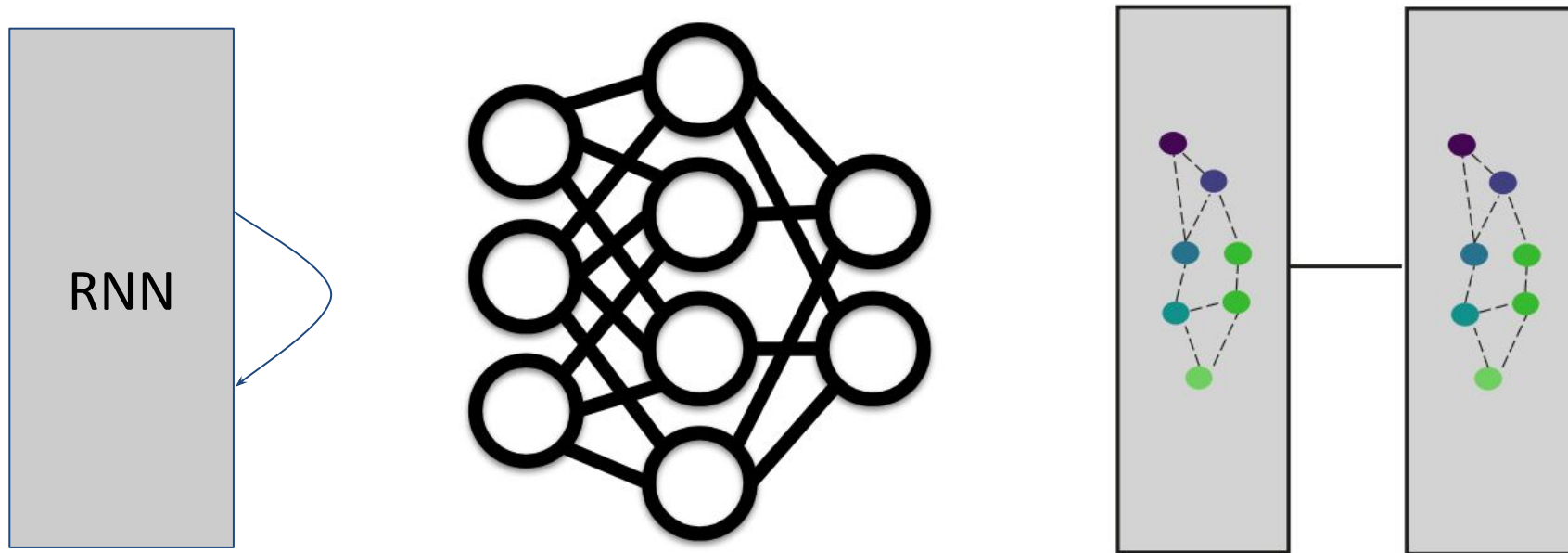
	$c_{k-2}$	$c_{k-1}$	max_pool_3x3	avg_pool_3x3	skip_connect	sep_conv_3x3	sep_conv_5x5	dil_conv_3x3	dil_conv_5x5	$c_k$
$c_{k-2}$	1	0	0	0	0	0	0	0	0	0
$c_{k-1}$	0	1	0	0	0	0	0	0	0	0
max_pool_3x3	0	0	0	0	0	0	1	0	0	0
avg_pool_3x3	0	0	0	0	0	0	0	0	0	0
skip_connect	0	0	0	0	0	1	0	0	0	0
sep_conv_3x3	0	0	0	0	0	0	1	0	0	0
sep_conv_5x5	0	0	0	0	0	0	0	1	0	0
dil_conv_3x3	0	0	0	0	0	0	0	0	1	0
dil_conv_5x5	0	0	0	0	0	0	0	0	0	1
$c_k$	0	0	0	0	0	0	0	0	0	0

One-hot adjacency encoding of a DARTS cell

# Structure-aware encodings: Adjacency matrix-based

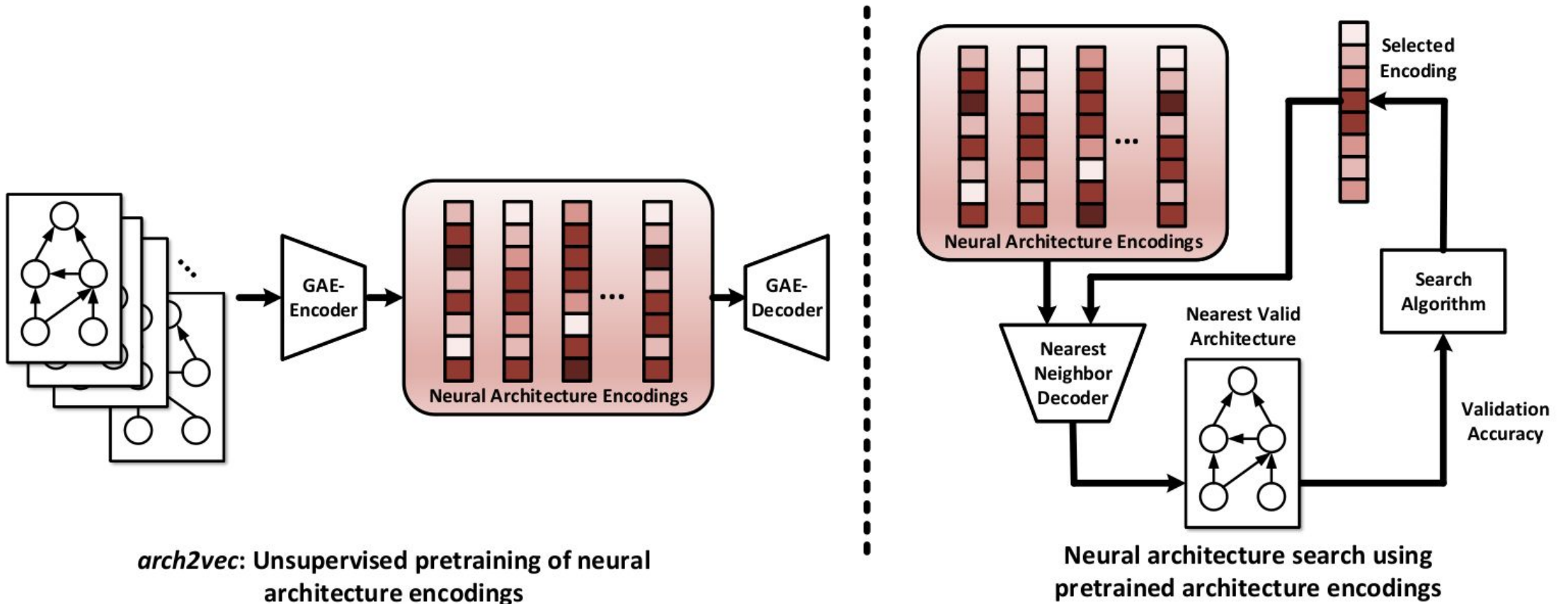


## Structure-aware encodings: LSTM/MLP/GCN

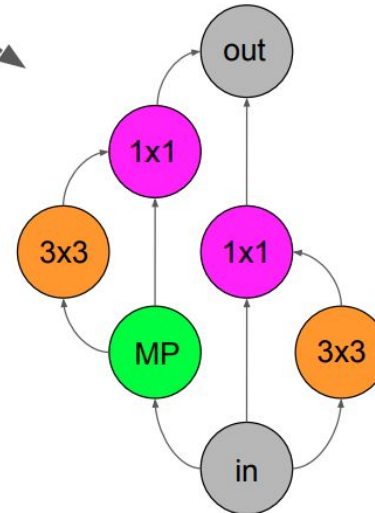
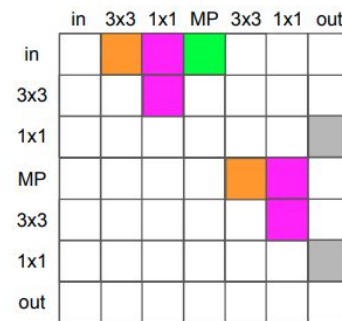
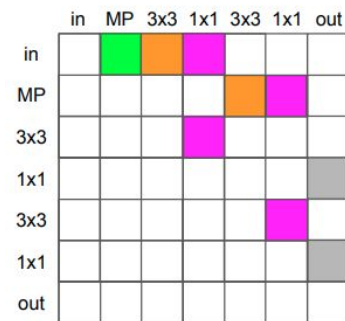
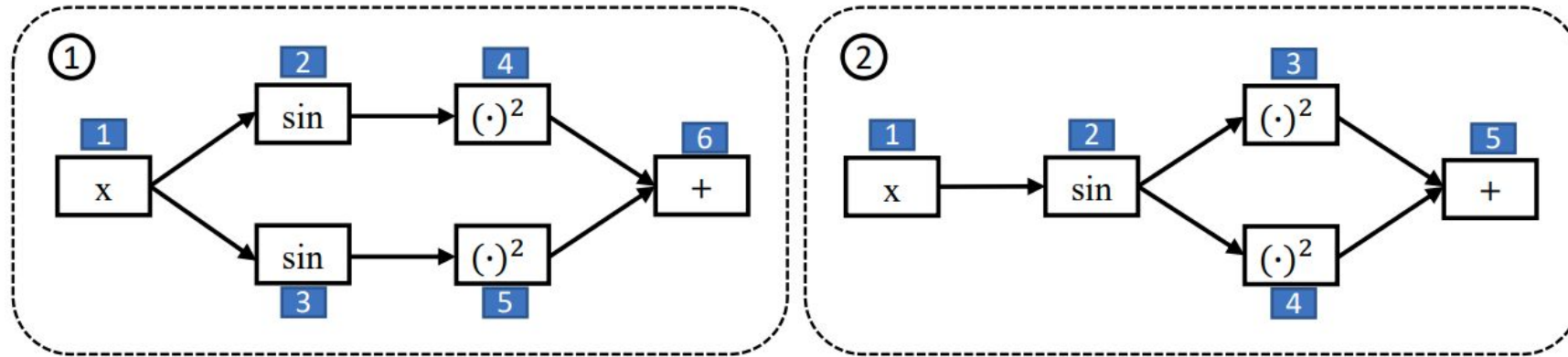


Different types of architecture encoders

# Structure-aware encodings: Unsupervised Pre-training

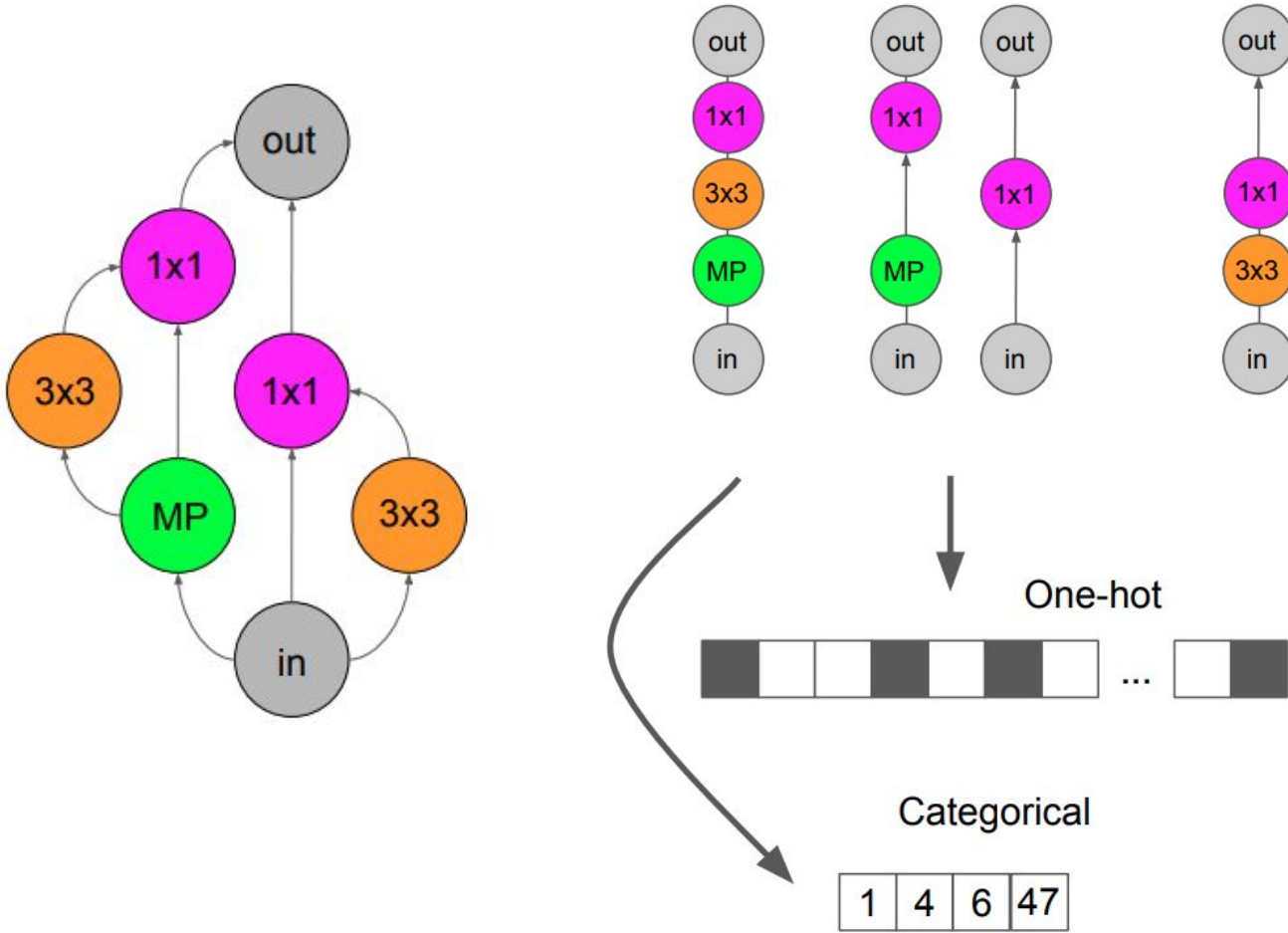


## Drawbacks of Structure-aware encodings



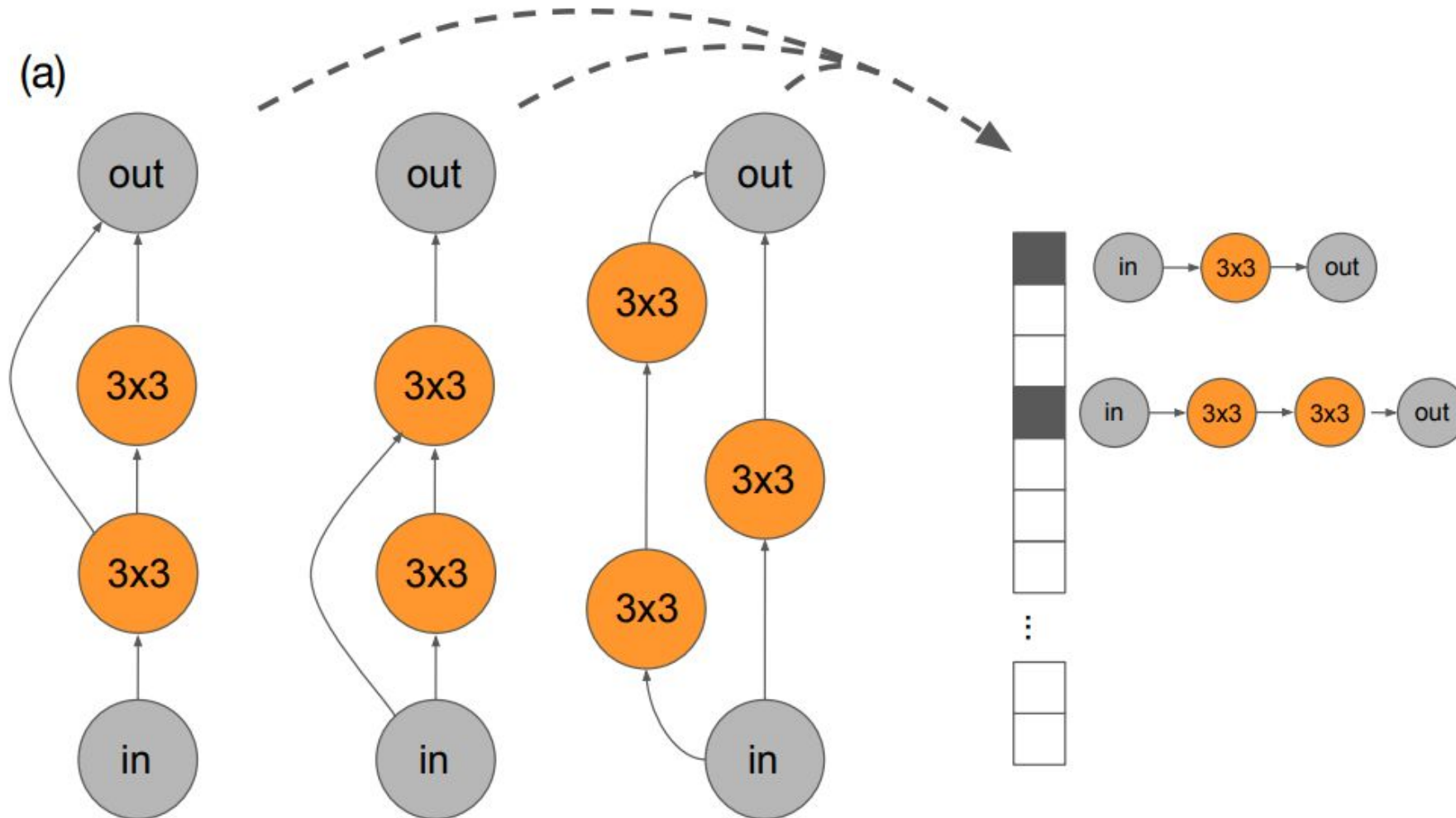


# Computation-aware encodings: Path Encodings

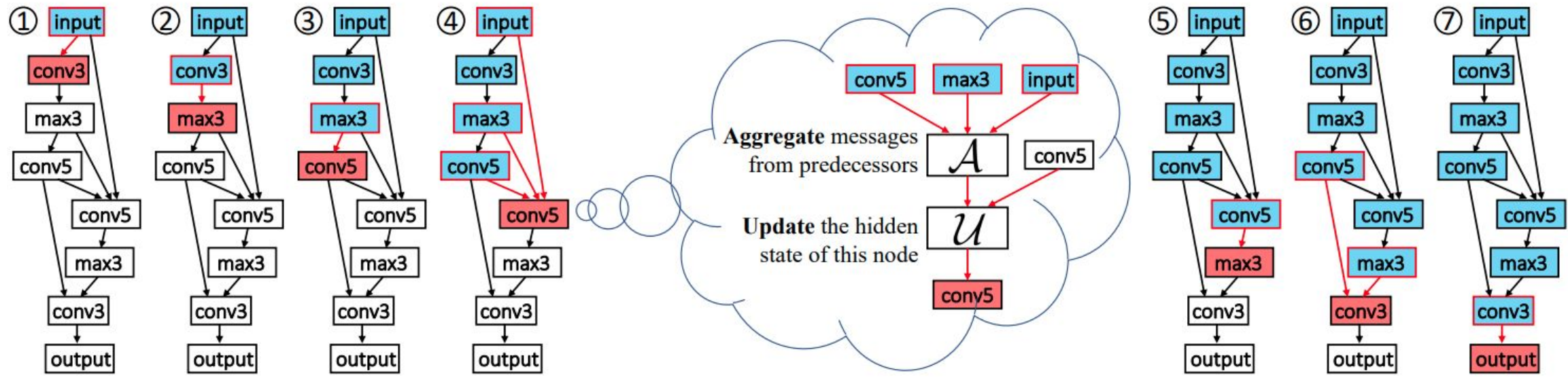


**One-hot / Categorical path encoding**

# Advantages of Computation-aware Encodings

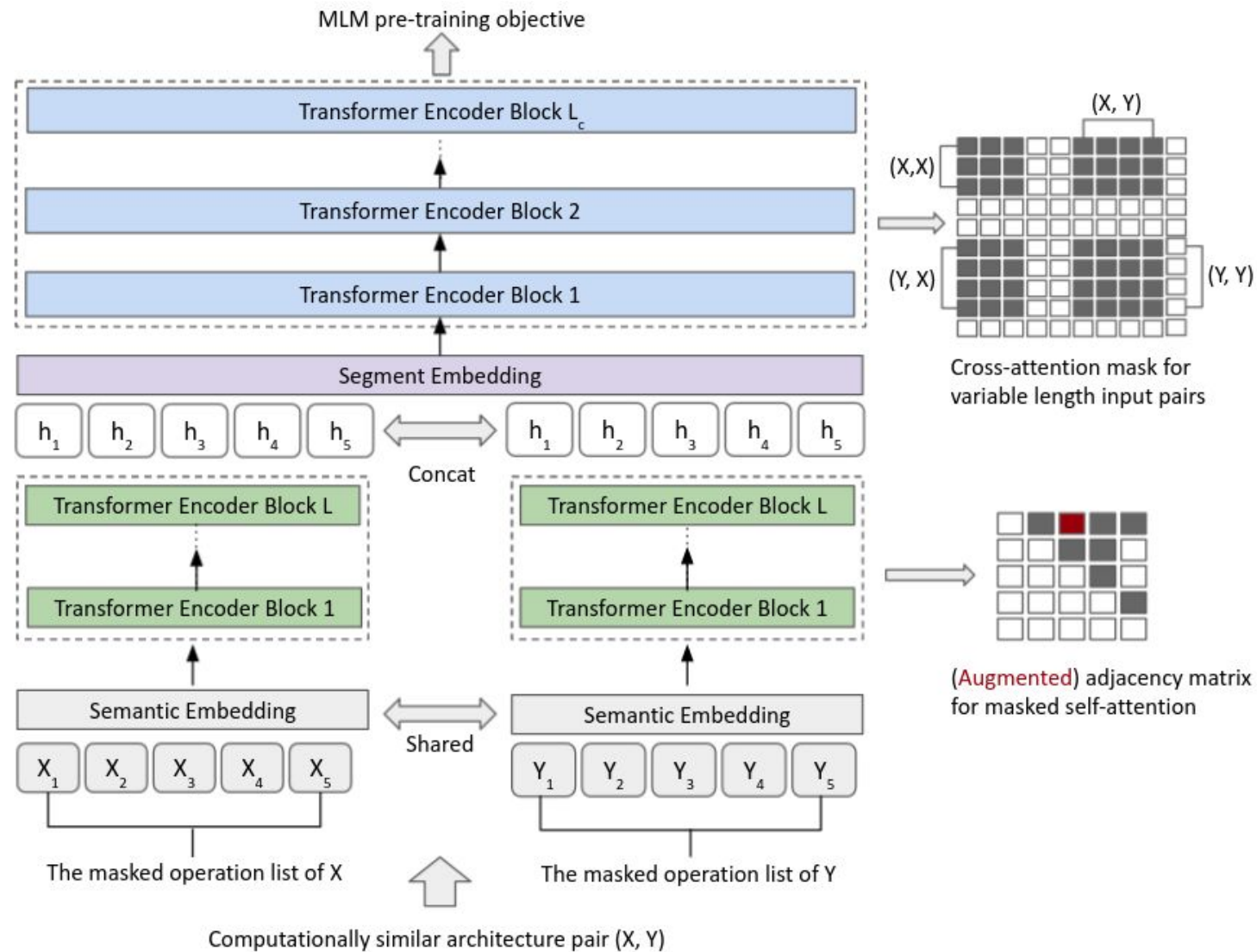


# Computation-aware encodings: D-VAE



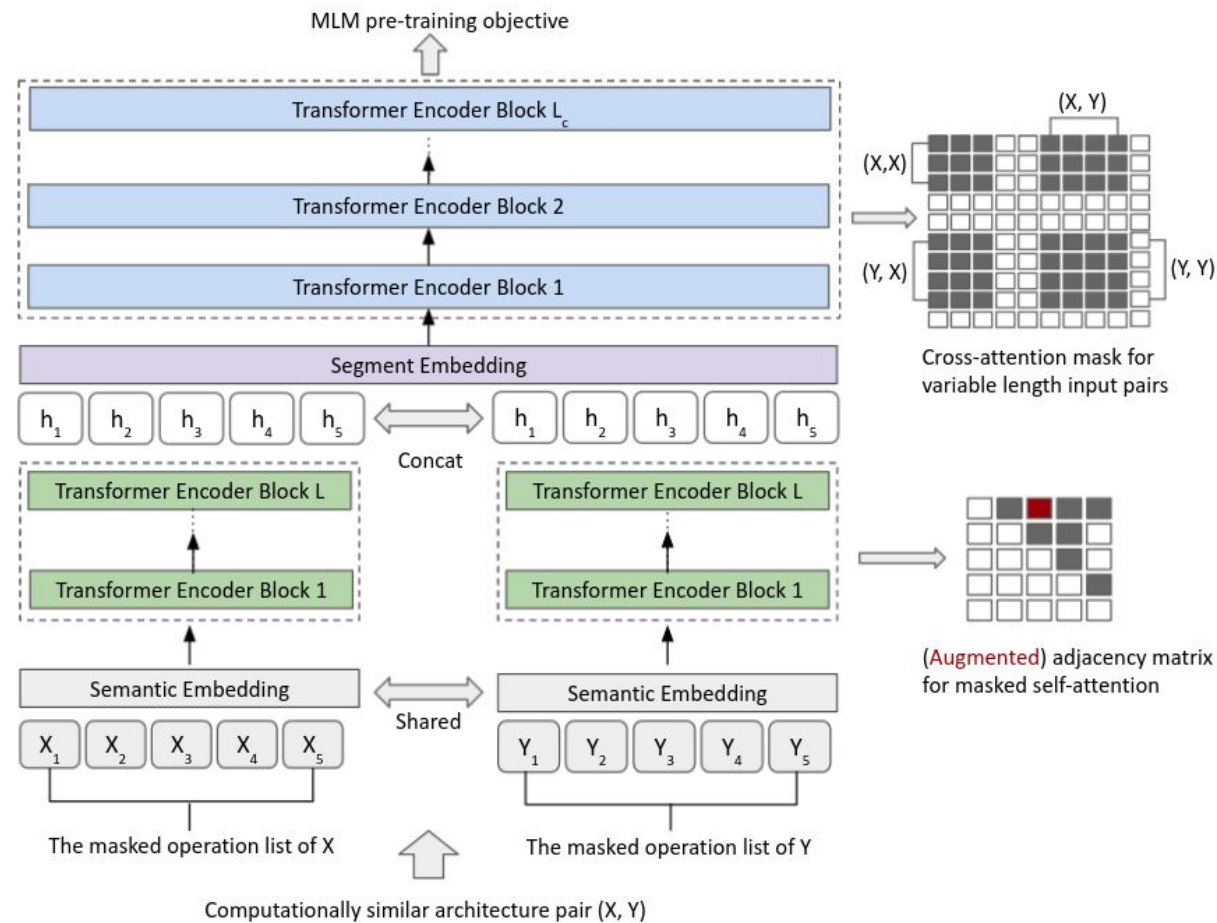
Muhan Zhang, Shali Jiang, Zhicheng Cui, Roman Garnett, Yixin Chen., D-VAE: A Variational Autoencoder for Directed Acyclic Graphs, NeurIPS 2019

# Our Proposed Method: CATE



# Why Pairwise Sampling

Conv 3x3 -> Conv 1x1 -> Conv 5x5 -> ~~Conv 1x1~~ -> Max Pool -> ...  
 Conv 3x3 -> Conv 1x1 -> Conv 5x5 -> ? -> Max Pool -> ...



# Attention Mask

---

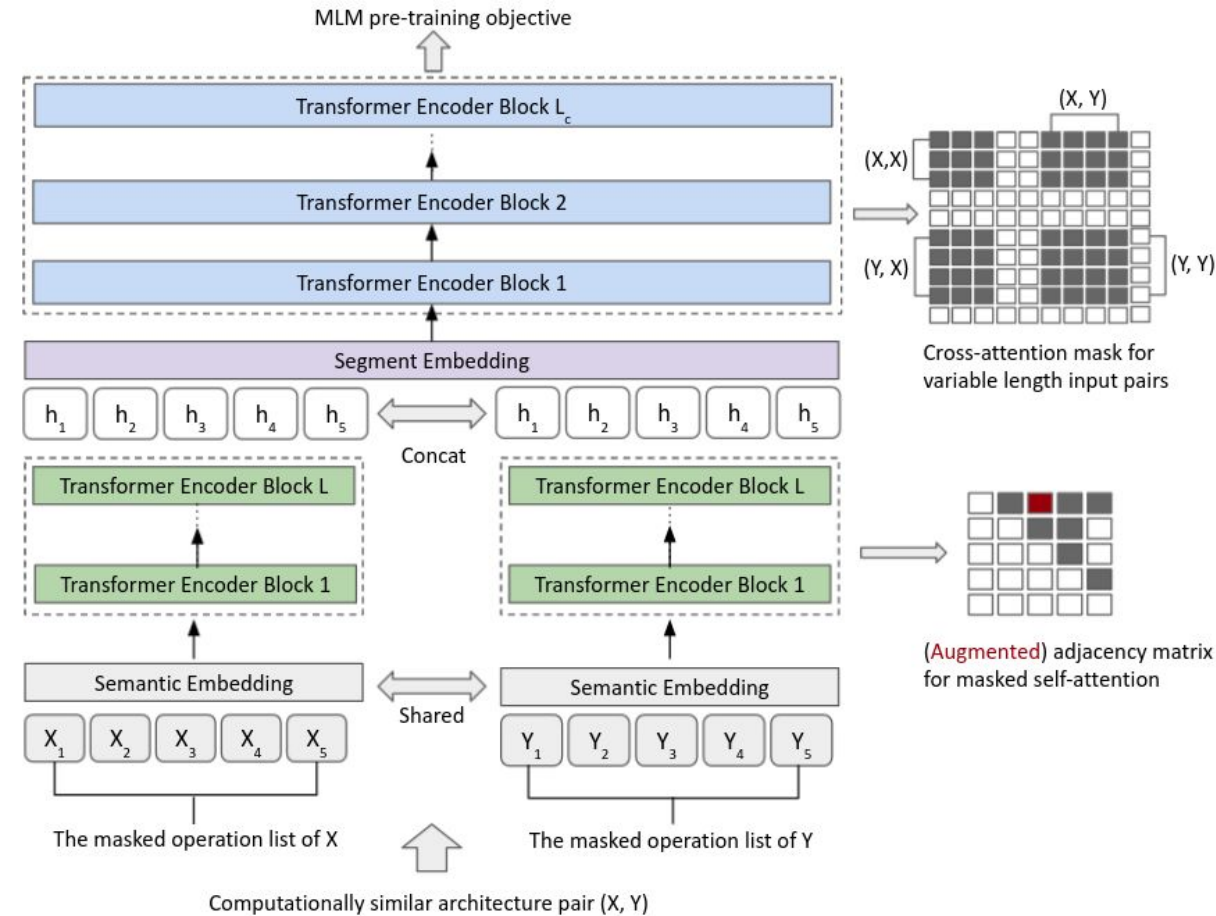
## Algorithm 1 Floyd Algorithm

---

- 1: **Input:** the node set  $\mathcal{V}$ , the adjacent matrix  $\mathbf{A}$
  - 2:  $\tilde{\mathbf{A}} \leftarrow \mathbf{A}$
  - 3: **for**  $k \in \mathcal{V}$  **do**
  - 4:   **for**  $i \in \mathcal{V}$  **do**
  - 5:     **for**  $j \in \mathcal{V}$  **do**
  - 6:        $\tilde{\mathbf{A}}_{i,j} \mid = \tilde{\mathbf{A}}_{i,k} \ \& \ \tilde{\mathbf{A}}_{k,j}$
  - 7: **Output:**  $\tilde{\mathbf{A}}$
- 

$$\mathbf{M}_{i,j}^{Direct} = \begin{cases} 0, & \text{if } A_{i,j} = 1 \\ -\infty, & \text{if } A_{i,j} = 0 \end{cases}$$

$$\mathbf{M}_{i,j}^{Indirect} = \begin{cases} 0, & \text{if } \tilde{A}_{i,j} = 1 \\ -\infty, & \text{if } \tilde{A}_{i,j} = 0 \end{cases}$$



# Encoding-dependent NAS Subroutines

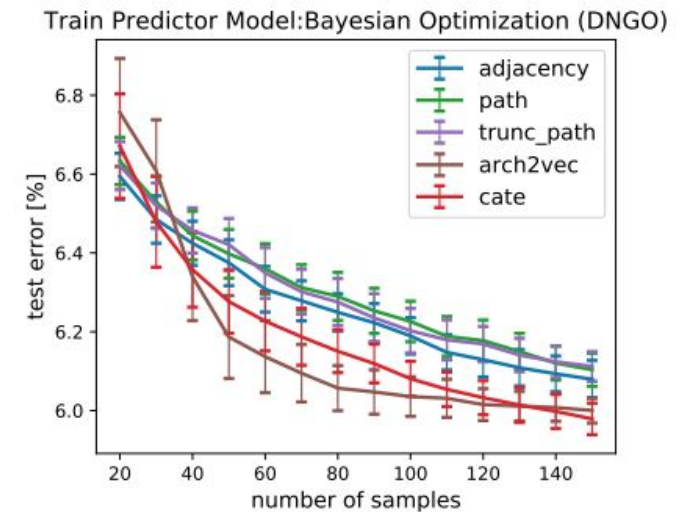
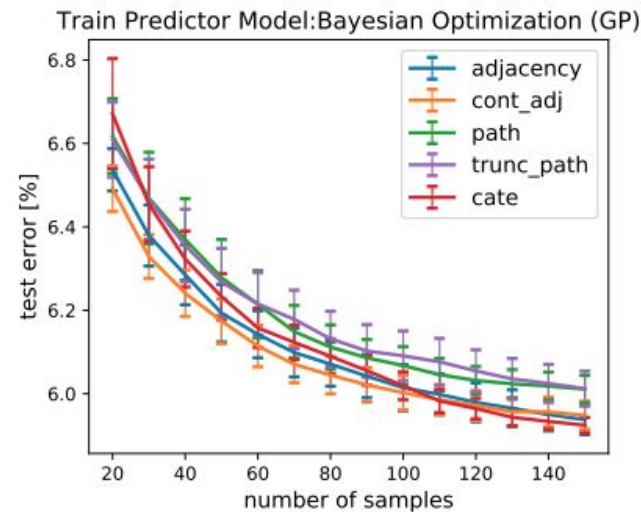
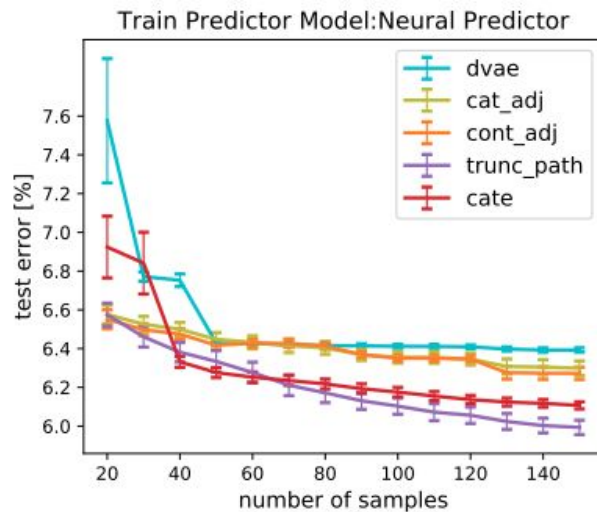
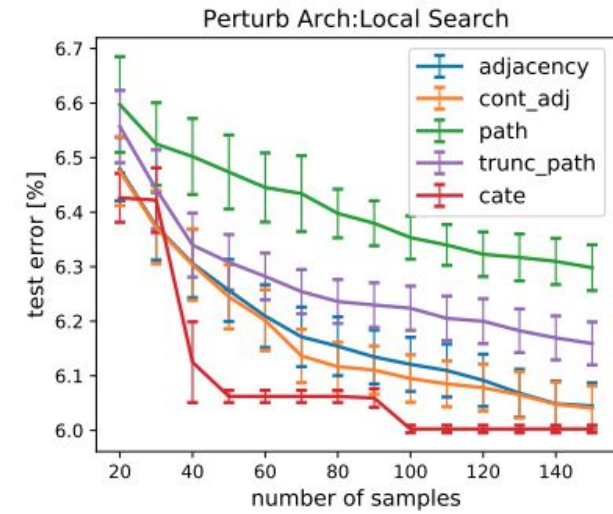
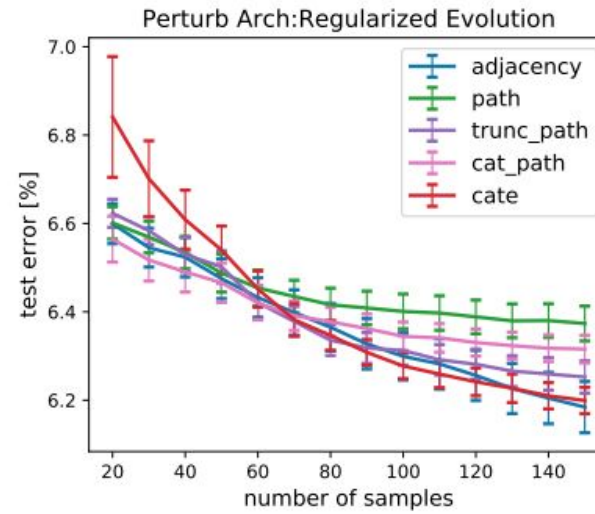
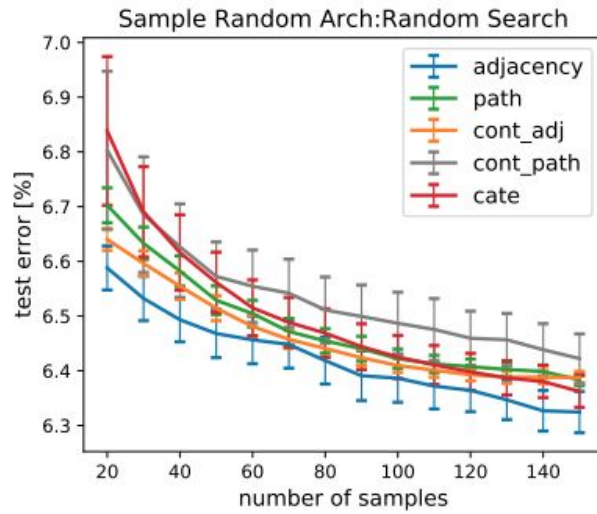
## 12 different encodings:

- One-hot/Categorical/Continuous adjacency matrix encoding (3)
- One-hot/Categorical/Continuous path encoding (3)
- The truncated counterparts (3)
- D-VAE
- arch2vec
- CATE

## 3 NAS subroutine:

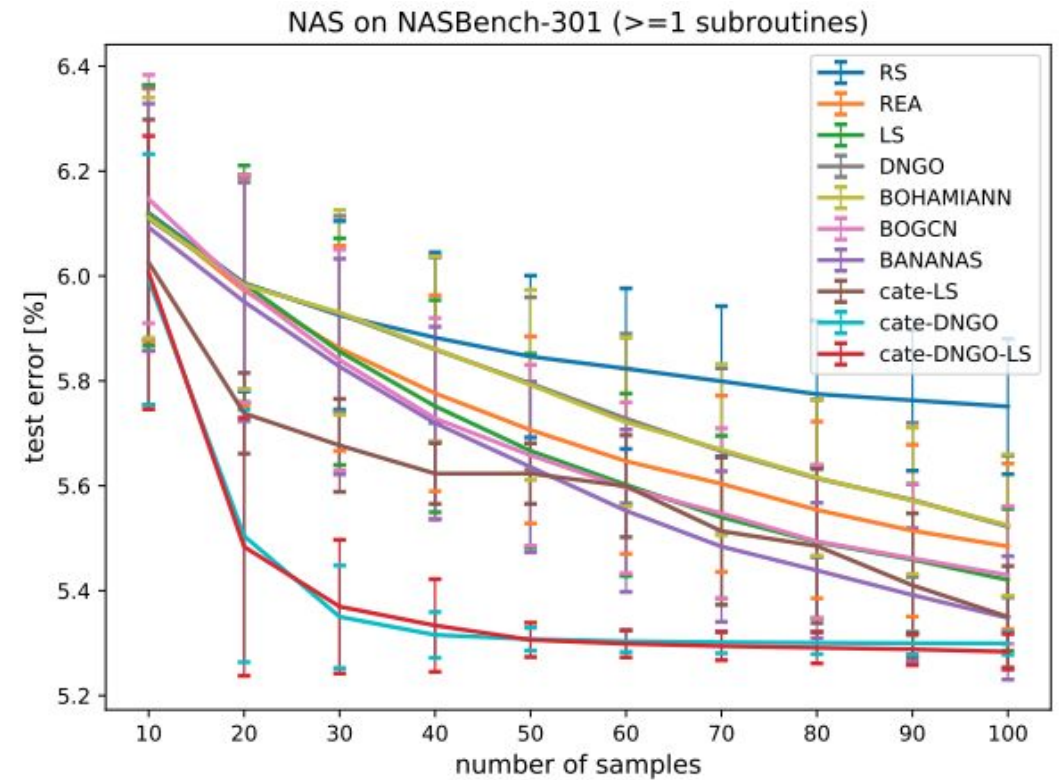
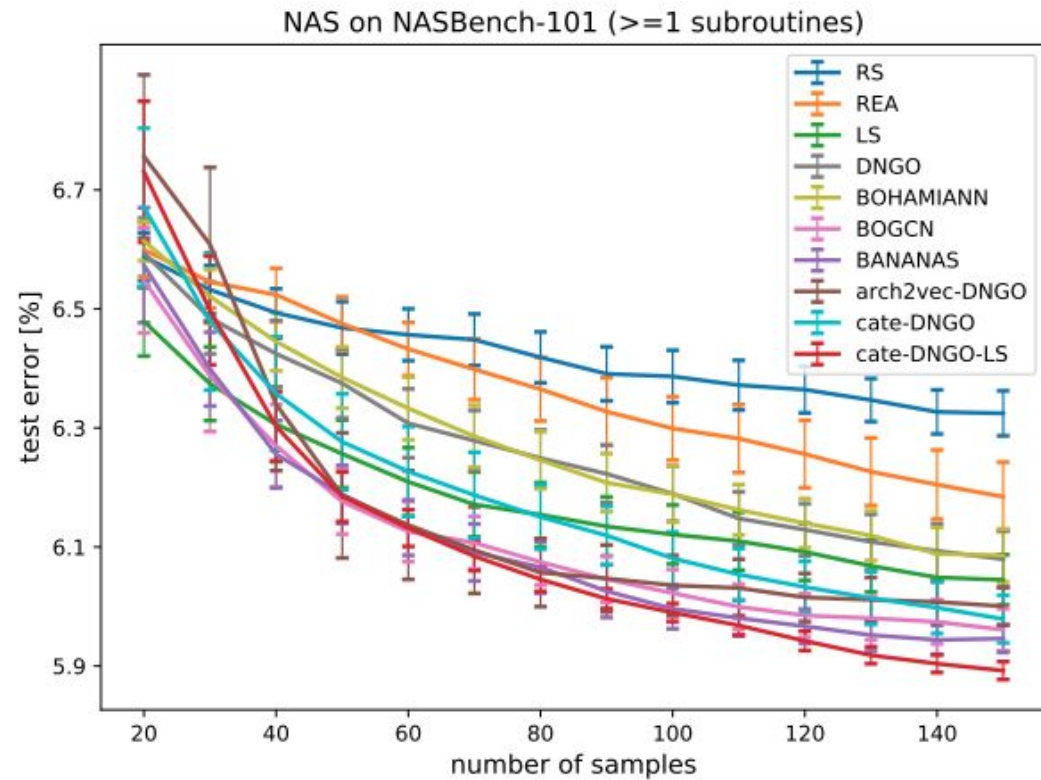
- Sample random architecture: random search
- Perturb architecture: regularized evolution, local search
- Train predictor: neural predictor, BO with GP, BO with DNGO

# Comparison between CATE and other encoding schemes





# Comparison between CATE and other NAS methods



## Evaluation on DARTS without surrogate models

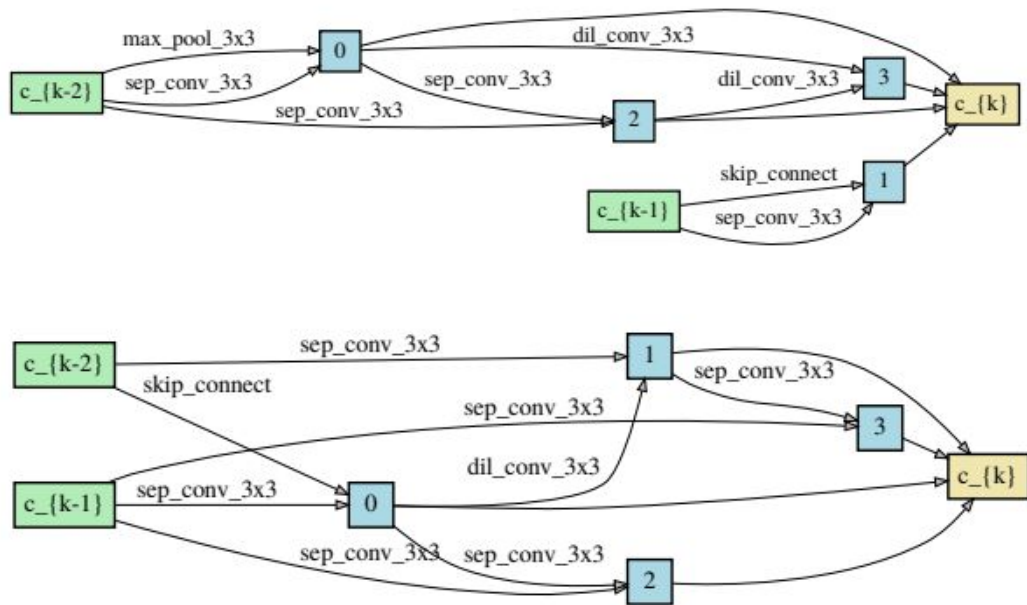


Figure 4. Top: Best found cell from CATE-DNGO-LS given the budget of 100 samples. Bottom: Best found cell from CATE-DNGO-LS given the budget of 300 samples.

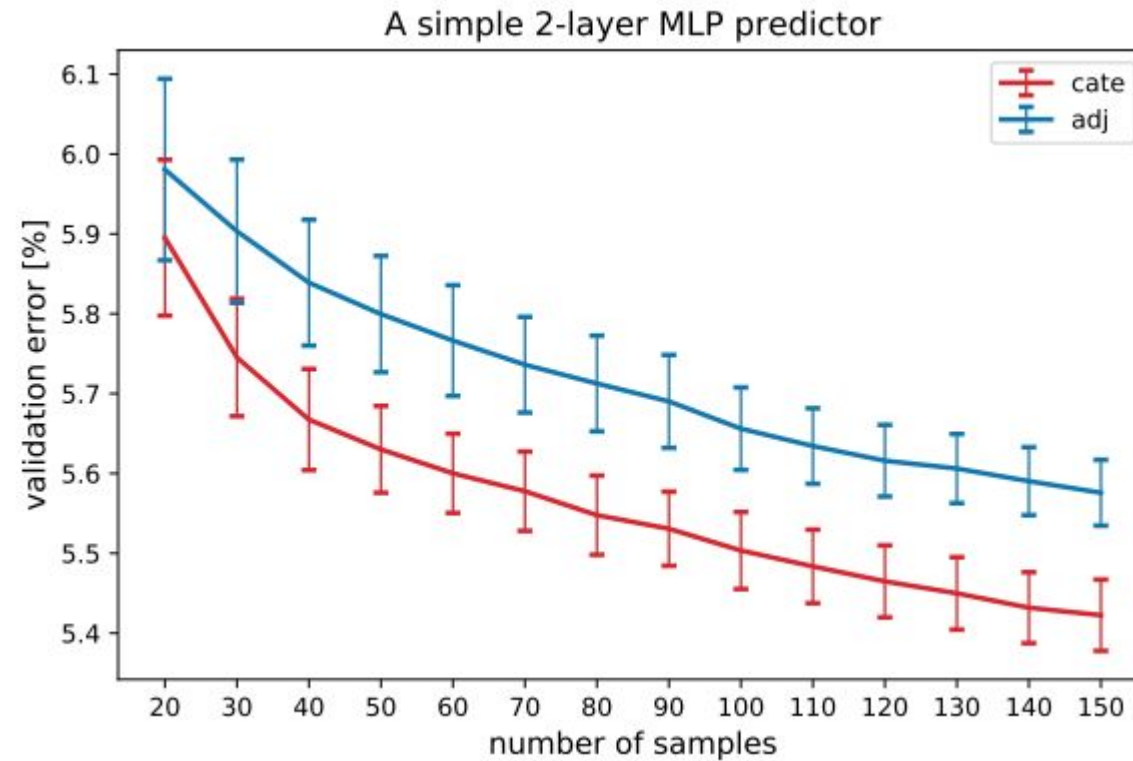
NAS Methods	Avg. Test Error (%)	Params (M)	Search Cost (GPU days)
RS (Li & Talwalkar, 2019)	$3.29 \pm 0.15$	<b>3.2</b>	4
DARTS (Liu et al., 2019a)	$2.76 \pm 0.09$	3.3	4
BANANAS (White et al., 2021)	$2.67 \pm 0.07$	3.6	11.8
arch2vec-BO (Yan et al., 2020)	$2.56 \pm 0.05$	3.6	9.2
CATE-DNGO-LS (small budget)	$2.55 \pm 0.08$	3.5	<b>3.3</b>
CATE-DNGO-LS (large budget)	<b><math>2.46 \pm 0.05</math></b>	4.1	10.3

Table 2. NAS results in DARTS search space using CIFAR-10.

NAS Methods	Params (M)	Mult-Adds (M)	Top-1 Test Error (%)
SNAS (Xie et al., 2019b)	4.3	522	27.3
DARTS (Liu et al., 2019a)	4.7	574	26.7
BayesNAS (Zhou et al., 2019)	<b>4.0</b>	<b>440</b>	26.5
arch2vec-BO (Yan et al., 2020)	5.2	580	25.5
BANANAS (ours)	5.1	576	26.3
CATE-DNGO-LS (small budget)	5.0	556	26.1
CATE-DNGO-LS (large budget)	5.8	642	<b>25.0</b>

Table 3. Transfer learning results on ImageNet.

# Evaluation on Outside Search Space



# Ablation Study

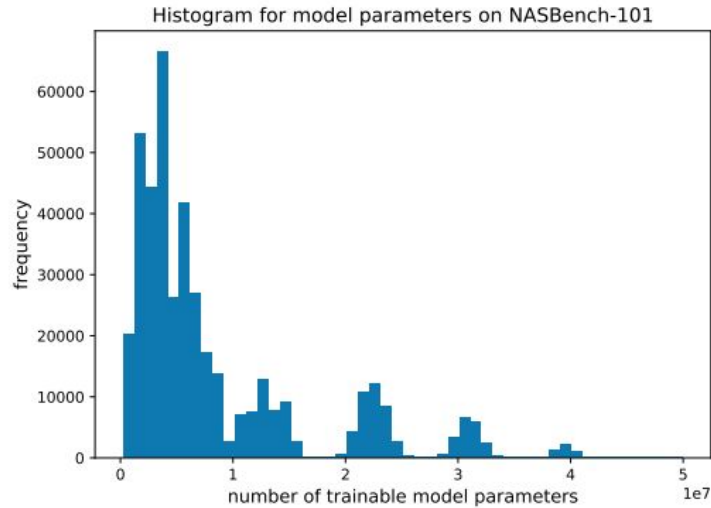


Figure 6. Histogram of model parameters on NAS-Bench-101.

$\delta \backslash K$	1	2	4	8
$1 \times 10^6$	6.02	5.95	5.99	5.95
$2 \times 10^6$	6.02	<b>5.94</b>	6.04	5.96
$4 \times 10^6$	<b>5.94</b>	6.03	6.05	5.99
$8 \times 10^6$	6.05	6.04	6.11	6.04

Table 4. Effects of  $\delta$  and  $K$  on architecture pair sampling.

$d_{ff} \backslash L_c$	6	12	24
64	6.07	5.99	5.95
128	6.01	<b>5.94</b>	5.95
256	5.97	<b>5.94</b>	<b>5.94</b>

Table 5. Effects of  $L_c$  and  $d_{ff}$  on pretraining CATE.

Mask type	NAS-Bench-101	NAS-Bench-301
Direct	6.03	5.35
Indirect	<b>5.94</b>	<b>5.30</b>

Table 6. Direct/Indirect dependency mask selection.

## Conclusion

- **A non-contrastive, pairwise pre-training method to learn computation-aware encodings with cross-attention Transformers**
- **Competitive under all encoding-dependent NAS subroutines in both small and large search spaces**
- **Superior generalization ability beyond the search space on which it was trained**

For more detailed information and code, please refer to our paper:

<https://arxiv.org/abs/2102.07108>

**Thank You**