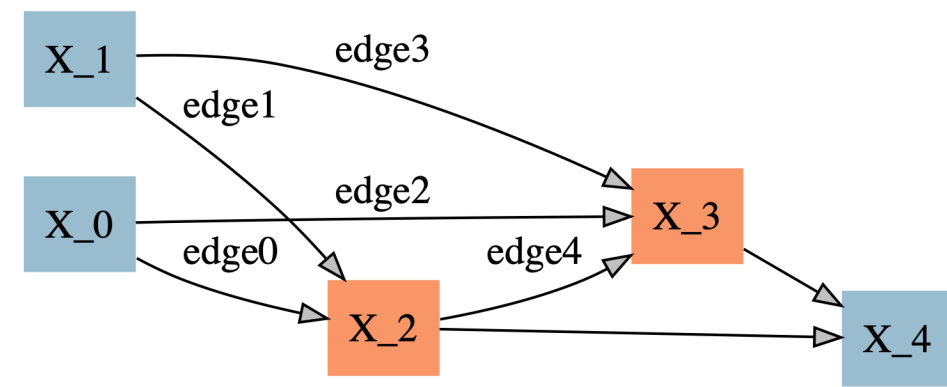


## Abstract

Motivated by three searching patterns of differentiable NAS, we

- Propose a unified view on searching algorithms of existing frameworks, transferring the global optimization to local cost minimization.
- Conduct empirical and theoretical analyses, revealing implicit biases in the cost's assignment mechanism and evolution dynamics that cause the phenomena.
- These biases indicate strong discrimination towards certain topologies.

## Search Space and Operation Selection

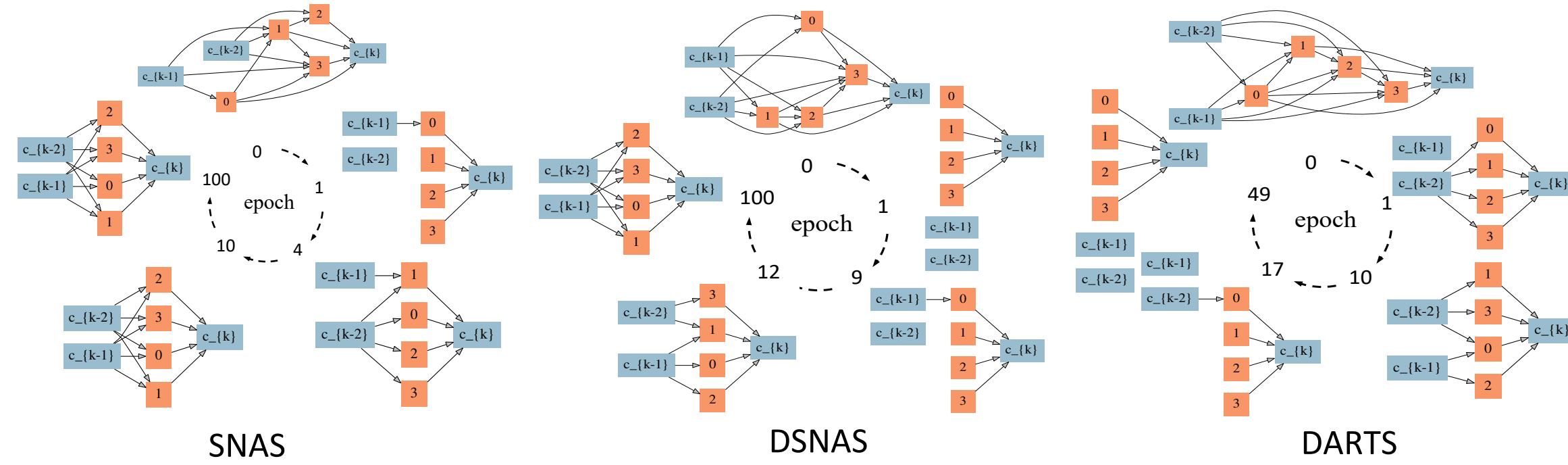


- Nodes  $X_i$  represent feature maps
- input nodes and output nodes (blue), intermediate nodes (orange)
- Operation candidates: None, skip connect, max\_pool\_3x3, avg\_pool\_3x3, sep\_conv\_3x3, dil\_conv\_3x3, dil\_conv\_5x5, sep\_conv\_5x5

- Each edge  $(i, j)$  multiply a one-hot variable  $\mathbf{Z}_{i,j}$  with operation candidates  $\mathbf{O}_{i,j}$ :

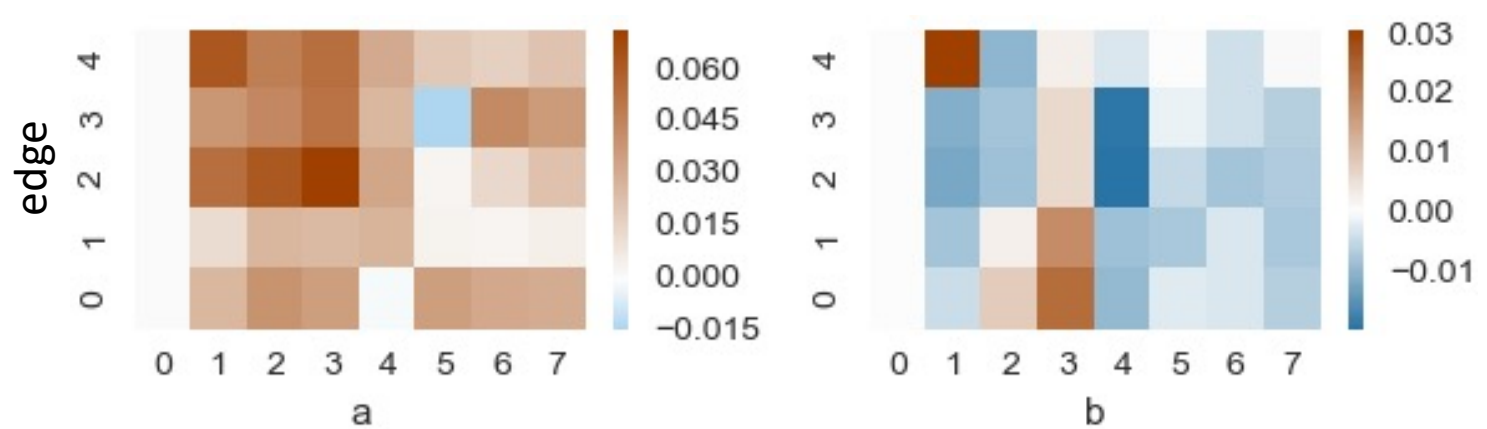
$$\tilde{\mathbf{O}}_{i,j}(\cdot) = \mathbf{Z}_{i,j}^T \mathbf{O}_{i,j}(\cdot)$$

## Evolution Pattern



- (P1)** Growing tendency: all edges are dropped in the beginning, some of them gradually recover
- (P2)** Width preference: intermediate edges barely recover from the initial drop
- (P3)** Catastrophic failure: in bi-level setting, no edge can recover from the initial drop

## Probing into the Evolution with Local Cost Mean



- the cost of None is 0
- the network topology is determined by signs of cost of other operations

**Hypothesis:** Cost of operations except *None* are positive *near* initialization and negative near convergence. Cell topology thus exhibits a tendency of growing.

## A close Look at Cost Assignment

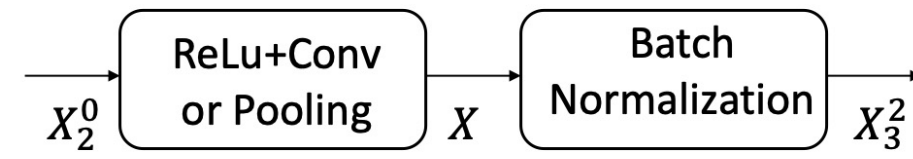
- Type1: edge 0, 1. E.g., *edge 0* involves two paths of back-propagation, i.e. 1<sup>st</sup> path (4-2-0) and 2<sup>nd</sup> path (4-3-2-0):

$$C(Z_{0,2}^s) = \frac{\partial L_{\theta}}{\partial \mathbf{X}_4^2} \mathbf{X}_2^0 + \frac{\partial L_{\theta}}{\partial \mathbf{X}_3^2} \frac{\partial \mathbf{X}_3^2}{\partial \mathbf{X}_2^0} \mathbf{X}_2^0$$

- Type2: edge 2, 3, 4. E.g., *edge 2* only involves one paths of back-propagation, i.e. path (4-3-0):

$$C(Z_{0,3}^s) = \frac{\partial L_{\theta}}{\partial \mathbf{X}_4^3} \frac{\partial \mathbf{X}_4^3}{\partial \mathbf{X}_3^0} \mathbf{X}_3^0 = \frac{\partial L_{\theta}}{\partial \mathbf{X}_3^0} \mathbf{X}_3^0$$

**Theorem 1.** A path does not distribute cost from the output edge after passing one intermediate edge.



Let  $X \in \mathbb{R}^{B \times C_{out} \times W_{out} \times H_{out}}$  denote the Conv output on edge 4, we expand  $C(Z_{0,2}^s)$  at path (4-3-2-0):

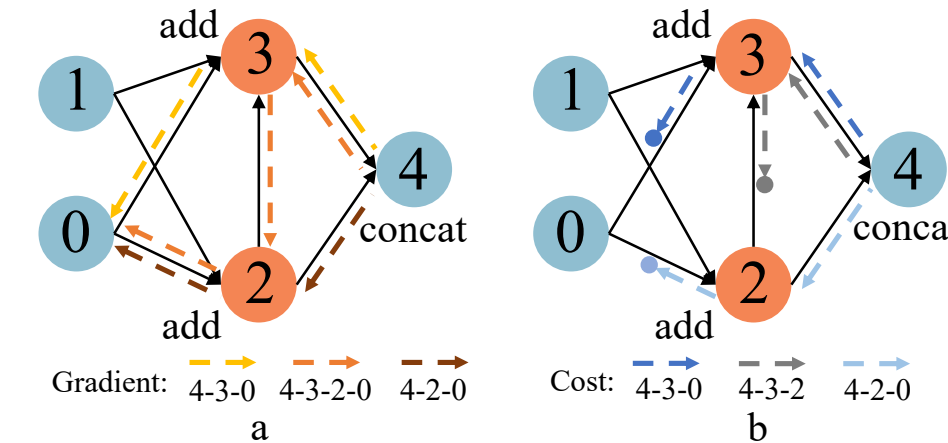
$$\frac{\partial L_{\theta}}{\partial \mathbf{X}_3^2} \frac{\partial \mathbf{X}_3^2}{\partial \mathbf{X}_2^0} \mathbf{X}_2^0 = \frac{\partial L_{\theta}}{\partial \mathbf{X}_3^2} \frac{\partial \mathbf{X}_3^2}{\partial X} \frac{\partial X}{\partial \mathbf{X}_2^0} \mathbf{X}_2^0$$

Expanding to the element-wise calculation:

$$\frac{\partial L_{\theta}}{\partial \mathbf{X}_3^2} \frac{\partial \mathbf{X}_3^2}{\partial X} \frac{\partial X}{\partial \mathbf{X}_2^0} \mathbf{X}_2^0 = \frac{\partial L_{\theta}}{\partial \mathbf{X}_3^2} \frac{\partial \mathbf{X}_3^2}{\partial X} X$$

Exploiting the property of normalization:

$$\frac{\partial L_{\theta}}{\partial \mathbf{X}_3^2} \frac{\partial \mathbf{X}_3^2}{\partial X} X = 0$$



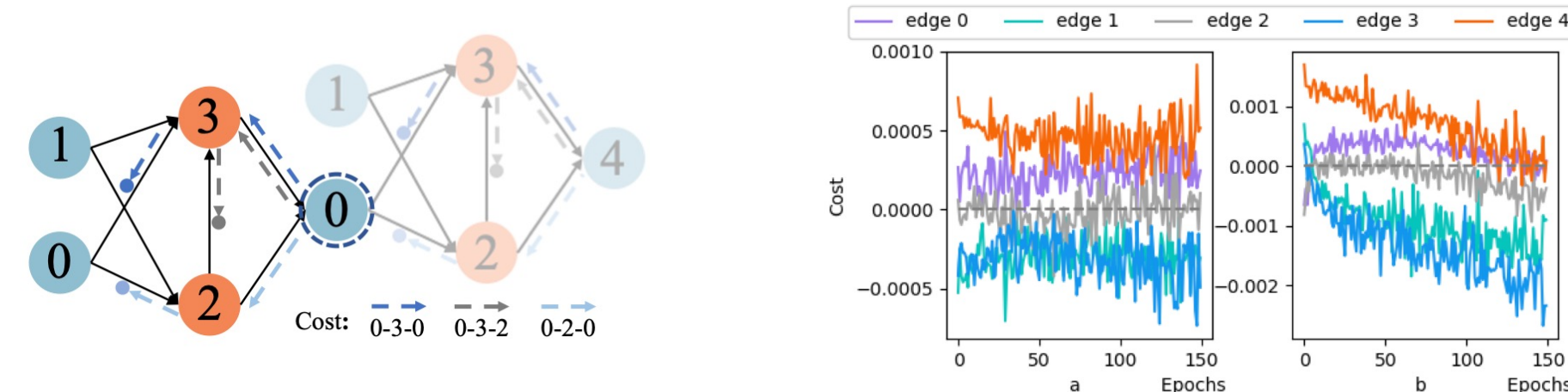
Total cost distributes to edges in the same cell as

$$\frac{\partial L_{\theta}}{\partial \mathbf{X}_{output}} \mathbf{X}_{output} = \sum_j \frac{\partial L_{\theta}}{\partial \mathbf{X}_{output}^j} (\sum_{i < j} \mathbf{X}_i^j)$$

where  $\mathbf{X}_{output}$  is the output node,  $j \in \{\text{intermediate nodes}\}$ ,  $i \in \{\text{nodes: } i < j\}$  e.g., in the minimal cell, we have

$$\frac{\partial L_{\theta}}{\partial \mathbf{X}_4} \mathbf{X}_4 = \frac{\partial L_{\theta}}{\partial \mathbf{X}_3^0} (\mathbf{X}_3^0 + \mathbf{X}_3^1 + \mathbf{X}_3^2) + \frac{\partial L_{\theta}}{\partial \mathbf{X}_2^0} (\mathbf{X}_2^0 + \mathbf{X}_2^1)$$

**Corollary.** In cells except the last, for intermediate nodes that are connected the same output node, cost of edges pointing to them sums up to 0.



Cost assignment in the last two cells. Cost of the left cell originates from node 0 of the right cell and should sum up to zero.

## Cost Dynamics

**Theorem 2.** Cost at output edges of the last cell has the form  $C_Z = L_{\theta} - H_{\theta}$ . It is **negatively related** to classification accuracy. It tends to be positive at low accuracy, negative at high accuracy.

- Negatively related:** for the last cell's output edges, cost of one batch  $M$  with sampled architecture  $\mathbf{Z}$  has an equivalent form:

$$C_Z = \sum_{b,c,d} \frac{\partial L_{\theta}}{\partial [\mathbf{X}_4]_{b,c,d}} [\mathbf{X}_4]_{b,c,d} = L_{\theta} - H_{\theta}$$

where  $L_{\theta}$  is the loss function,  $H_{\theta}$  is the entropy of network output.

- Positive at low accuracy:** Exploiting normalization and weight initialization,

$$\mathbb{E}_{\theta_0} [C_Z > 0],$$

$$\text{since } \mathbb{E}_{y_1, y_2 \sim \mathcal{N}(0, \cdot)} [y_1 \exp(y_1 + y_2)] > 0.$$

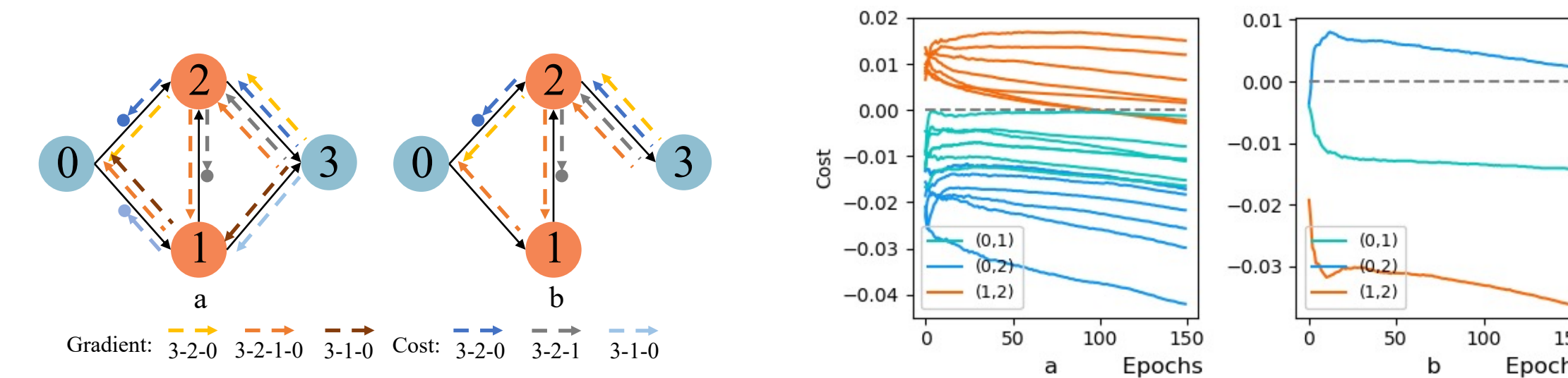
- Negative at high accuracy:** With operation parameters updated to convergence, the probability of  $b$ -th image being classified to the correct label  $n_b$  increases towards 1. Since  $Y_{bn_b} = \max\{Y_{bn}\}$ , we have

$$C_Z \propto \sum_n [-Y_{bn_b} + \frac{\exp(Y_{bn_b})}{\sum_q \exp(Y_{bq})} Y_{bn_b}] \leq \sum_n [-Y_{bn_b} + \frac{\exp(Y_{bn_b})}{\sum_q \exp(Y_{bq})} Y_{bn_b}] = 0$$

**Remark:** If the trends of cost at all edges are consistent with Thm.2, eventually the growing tendency (**P1**) occurs.

## Distinction of Intermediate Edges

By the cost assignment mechanism, all edges are born equally. Can not explain the distinction of intermediate edges. Conduct analysis on simplified minimal cell:



Cost at intermediate edge(1,2) is higher than cost of the same operation at edge(0, 2).

- $\theta_{0,1}$  and  $\theta_{0,2}$  are always trained as None is not sampled.
- $\theta_{0,1}$  is updated with gradients from two paths (3-2-1-0) and (3-1-0).
  - None* is sampled on edge(1, 2),  $\theta_{0,1}$  can be updated.
  - None* is sampled on edge (0, 1),  $\theta_{1,2}$  cannot be updated.
- Even if *None* is not included in edge (0, 1), there are more model instances on path (3-2-1-0) than path (3-2-0) and (3-1-0).

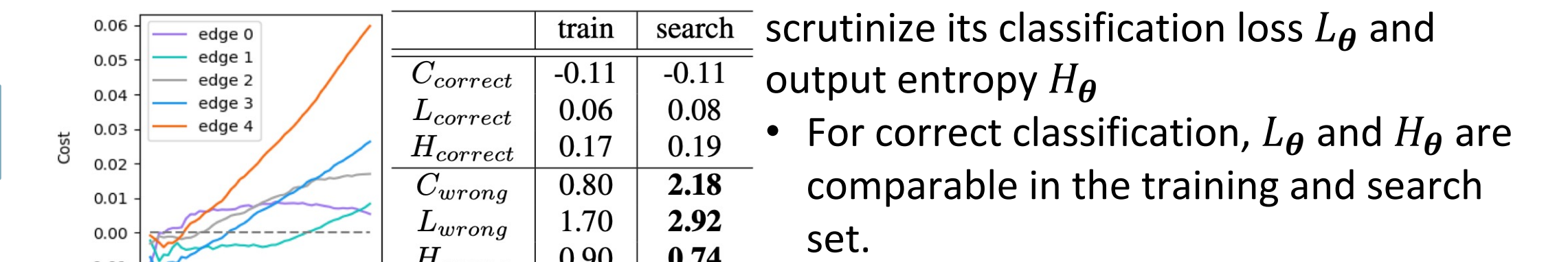
Design controlled experiment by deleting the output edge(1, 3) and fixing operation at edge(0, 1)

- path (3-2-1-0) is deeper than (3-2-0), but the cost on edge(0, 2) becomes positive
- the distinction of intermediate edges is likely due to **unequal** training

**Remark:** intermediate edge tends to have higher cost than input edges (if one with positive cost, it must be the intermediate edge.)

## Effect of Bi-level Optimization

Cost minimization formulation satisfies to single-level and bi-level ones. **Exception:** every edge drops and almost none of them finally recovers in DARTS's bi-level version



- But for data classified incorrectly, the classification loss  $L_{\theta}$  is much larger in the search set. That is, data in the search set are classified poorly.
- $H_{\theta}$  in the search set is much lower than its counterpart in the training set.

**Remark:** explains catastrophic failure (**P3**) in bi-level DARTS

## Conclusion

- The cost assignment for architecture  $\alpha$  differs from that in gradient back-prop for parameters  $\theta$ . Exaggerating discrepancy leads to distinction of intermediate edges.
- During training, cost decreases from **positive** and eventually turns **negative**, promoting the tendency of topology growth. If this cost-decreasing process is hindered, as in bi-level optimization, there will be a catastrophic failure.

Reference:

Liu et al. DARTS: Differentiable architecture search. ICLR 2019.  
Xie et al. Snas: stochastic neural architecture search. ICLR 2019.  
Hu\* and Xie\* et al. Dsnas: Direct neural architecture search without parameter retraining. CVPR 2020.

