

Automatic Relation-aware Graph Network Proliferation

– Supplementary Material –

Shaofei Cai^{1,2}, Liang Li^{1*}, Xinzhe Han^{1,2}, Jiebo Luo³, Zheng-Jun Zha⁴, Qingming Huang^{1,2,5}

¹Key Lab of Intell. Info. Process., Inst. of Comput. Tech., CAS, Beijing, China

²University of Chinese Academy of Sciences, Beijing, China, ³University of Rochester

⁴University of Science and Technology of China, China, ⁵Peng Cheng Laboratory, Shenzhen, China

{shaofei.cai, xinzhe.han}@vip1.ict.ac.cn, liang.li@ict.ac.cn,

jluo@cs.rochester.edu, zhazj@ustc.edu.cn, qmhuang@ucas.ac.cn

A. Search Space Details

In this section, we first detail how node-learning operations and relation-learning operations are formulated. Then we introduce the design of task-based layers for different graph learning tasks. $\mathbf{V} \in \mathbb{R}^{n \times d_V}$, $\mathbf{E} \in \mathbb{R}^{m \times d_E}$ are node features and relation features that are fed forward to node-learning and relation-learning operations, respectively. n and m are the number of nodes and edges of the input graph data.

A.1. Node-learning Operations

As discussed in the body content, a node-learning operation $\mathcal{O}_{\mathcal{V}}^{(i,j)}$ computes the transformed information from feature vertex \mathbf{V}_i to \mathbf{V}_j , which is denoted as $\mathbf{V}_{i \rightarrow j} = \mathcal{O}_{\mathcal{V}}^{(i,j)}(\mathbf{V}_i, \mathbf{E}_i, f_{i,j})$. Specifically, given node t on the graph, its transformed node feature $\mathbf{V}_{i \rightarrow j}^t$ is formulated as

$$\mathbf{V}_{i \rightarrow j}^t = f_{i,j}(\{\gamma_{s,t} \odot \mathbf{V}_i^s + \beta_{s,t} | s \in \mathcal{N}(t)\}), \quad (1)$$

$$\gamma_{s,t}, \beta_{s,t} = g(\mathbf{E}_i^{s,t}; \theta),$$

where, $\mathcal{N}(t)$ denotes the set of neighbors of target node t on the graph, $f_{i,j}$ is an aggregating function. $g(\cdot)$ is a two-layer multilayer perceptron to compute the affine transformation with the learnable parameters $\theta = [\mathbf{W}^1 \ \mathbf{W}^2 \ \mathbf{W}^k \ \mathbf{W}^b]$, which is formulated as

$$[\gamma_{s,t} \ \beta_{s,t}] = \sigma(\sigma(\mathbf{E}_i^{s,t} \mathbf{W}^1) \mathbf{W}^2) [\mathbf{W}^k \ \mathbf{W}^b], \quad (2)$$

where $\sigma(\cdot)$ is the rectified linear unit. The difference between different node-learning operations lies in the choice of message aggregating functions. We design 8 candidate aggregating function options, i.e., V_MEAN , V_SUM , V_MAX , V_STD , V_GEM2 , V_GEM3 , $zero$ and $skip-connect$ for capturing different types of information. For clarity, we

denote $\mathbf{M}_{s,t}$ as the modulated incoming message from node t to node s , where $\mathbf{M}_{s,t} = \gamma_{s,t} \odot \mathbf{V}_i^s + \beta_{s,t}$. The mean aggregation of neighboring messages is denoted as

$$\mu_t(\mathbf{M}) = \frac{1}{|\mathcal{N}(t)|} \sum_{s \in \mathcal{N}(t)} \mathbf{M}_{s,t}. \quad (3)$$

V_MEAN is the average of neighboring messages, which captures the mean statistics of neighboring messages [4], written as

$$\mathbf{V}_{i \rightarrow j}^t = \mu_t(\mathbf{M}). \quad (4)$$

V_SUM is the sum of neighboring messages, which captures local structural information [9], written as

$$\mathbf{V}_{i \rightarrow j}^t = |\mathcal{N}(t)| \times \mu_t(\mathbf{M}). \quad (5)$$

V_MAX is the max of neighboring messages, which captures the representative information [9], written as

$$\mathbf{V}_{i \rightarrow j}^t = \max_{s \in \mathcal{N}(t)} \mathbf{M}_{s,t}. \quad (6)$$

V_STD is the standard deviation of input feature set, which captures the stability of neighboring messages [2], i.e.,

$$\mathbf{V}_{i \rightarrow j}^t = \sqrt{ReLU(\mu_t(\mathbf{M}^2) - \mu_t(\mathbf{M})) + \epsilon}, \quad (7)$$

where $ReLU$ is the rectified linear unit used to avoid negative values caused by numerical errors and ϵ is a small positive number to ensure the output is differentiable [2].

Moreover, we also introduce Generalized Mean Pooling (GeM) for aggregating messages, which can focus on learning to propagate the prominent message [1], written as

$$GeM_t(\mathbf{M}, \alpha) = \sqrt[\alpha]{ReLU(\mu_t(\mathbf{M}^\alpha)) + \epsilon}, \quad (8)$$

where α is the hyper parameter. Here, we adopt two widely used parameters to construct the node-learning operations, i.e., V_GEM2 and V_GEM3 .

*Corresponding author.

V_GEM2:

$$\mathbf{V}_{i \rightarrow j}^t = \text{Gem}_t(\mathbf{M}, 2). \quad (9)$$

V_GEM3:

$$\mathbf{V}_{i \rightarrow j}^t = \text{Gem}_t(\mathbf{M}, 3). \quad (10)$$

skip-connect is to enhance the central node information and mitigate the gradient vanishing, written as

$$\mathbf{V}_{i \rightarrow j}^t = \mathbf{V}_i^t. \quad (11)$$

zero operation is included in the search space to indicate a lack of connection. Links that are important should have a low weight in the **zero** operation [6]. It is formulated as

$$\mathbf{V}_{i \rightarrow j}^t = \mathbf{0}. \times \mathbf{V}_i^t. \quad (12)$$

A.2. Relation-mining Operations

A relation-mining operation $\mathcal{O}_{\mathcal{E}}^{(i,j)}$ computes the transformed relational information $\mathbf{E}_{i \rightarrow j} = \mathcal{O}_{\mathcal{E}}^{(i,j)}(\mathbf{V}_i, \mathbf{E}_i, h_{i,j})$. $h_{i,j}$ is a relation-mining network, such as *subtraction*, *hardward product*, *gauss kernel*, etc. Specifically, given a specific edge (s, t) on the graph, $\mathbf{E}_{i \rightarrow j}^{s,t}$ is computed using feature-wise linear modulation:

$$\begin{aligned} \mathbf{E}_{i \rightarrow j}^{s,t} &= \gamma_{s,t} \odot \mathbf{E}_i^{s,t} + \beta_{s,t}, \\ \gamma_{s,t}, \beta_{s,t} &= h_{i,j}(\mathbf{V}_i^s, \mathbf{V}_i^t; \theta), \end{aligned} \quad (13)$$

$\gamma_{s,t}, \beta_{s,t}$ is the affine transformation learned by $h_{i,j}$ with the learnable parameters $\theta = [\mathbf{W}^1 \ \mathbf{W}^2 \ \mathbf{W}^k \ \mathbf{W}^b]$ and relation function $h^*(\cdot, \cdot)$

$$[\gamma_{s,t} \ \beta_{s,t}] = \sigma(\sigma(h^*(\mathbf{V}_i^s, \mathbf{V}_i^t) \mathbf{W}^1) \mathbf{W}^2) [\mathbf{W}^k \ \mathbf{W}^b], \quad (14)$$

where $\sigma(\cdot)$ is the rectified linear unit. The difference between different relation-mining operations lies in the choice of relation functions $h^*(\cdot, \cdot)$. We design 8 candidate relation functions, i.e., *E_SUB*, *E_GAUSS*, *E_HAD*, *E_MAX*, *E_SUM*, *E_MEAN*, *skip-connect*, and *zero* for capturing different types of relational information.

E_SUB captures the relative change between two nodes. Its relation function is computed as

$$h^*(\mathbf{V}_i^s, \mathbf{V}_i^t) = \mathbf{V}_i^s - \mathbf{V}_i^t \quad (15)$$

E_GAUSS measures the distance between the central node and its neighboring nodes:

$$h^*(\mathbf{V}_i^s, \mathbf{V}_i^t) = \exp\left(\frac{|\mathbf{V}_i^s - \mathbf{V}_i^t|^2}{2\sigma}\right) \quad (16)$$

E_HAD emphasizes on learning the commonalities between the central node and its neighbors:

$$h^*(\mathbf{V}_i^s, \mathbf{V}_i^t) = \mathbf{V}_i^s \odot \mathbf{V}_i^t \quad (17)$$

E_SUM:

$$h^*(\mathbf{V}_i^s, \mathbf{V}_i^t) = \mathbf{V}_i^s + \mathbf{V}_i^t \quad (18)$$

E_MAX:

$$h^*(\mathbf{V}_i^s, \mathbf{V}_i^t) = \max(\mathbf{V}_i^s, \mathbf{V}_i^t) \quad (19)$$

E_MEAN:

$$h^*(\mathbf{V}_i^s, \mathbf{V}_i^t) = \frac{\mathbf{V}_i^s + \mathbf{V}_i^t}{2} \quad (20)$$

skip-connect operation is to enhance the original relational information and mitigate the gradient vanishing, written as

$$[\gamma_{s,t} \ \beta_{s,t}] = [\mathbf{1} \ \mathbf{0}] \quad (21)$$

zero operation is included in the search space to indicate a lack of connection. Links that are important should have a low weight in the **zero** operation [6]. It is formulated as

$$[\gamma_{s,t} \ \beta_{s,t}] = [\mathbf{0} \ \mathbf{0}] \quad (22)$$

A.3. Task-based Layer

We design the final network layers depending on the specific task. Suppose there are 8 node feature vertices $\{\mathbf{V}_1, \dots, \mathbf{V}_8\}$ and 8 relation feature vertices $\{\mathbf{E}_1, \dots, \mathbf{E}_8\}$ in our GNN architecture, we compute the global node features \mathbf{V}_g and global edge features \mathbf{E}_g using the following formula

$$\begin{aligned} \mathbf{V}_g &= \sigma(\text{BN}([\mathbf{V}_1 \parallel \dots \parallel \mathbf{V}_8] \mathbf{W}_V)), \\ \mathbf{E}_g &= \sigma(\text{BN}([\mathbf{E}_1 \parallel \dots \parallel \mathbf{E}_8] \mathbf{W}_E)), \end{aligned} \quad (23)$$

where $\mathbf{W}_V \in \mathbb{R}^{d_V \times d_V}$, $\mathbf{W}_E \in \mathbb{R}^{d_E \times d_E}$ are the learnable parameters, BN denotes batch normalization operation, $\sigma(\cdot)$ is the rectified linear unit, $[\cdot \parallel \cdot]$ denotes the feature concatenation operation. The global graph representation \mathbf{G}_g is computed using mean-pooling readout operation over global node features and global edge features, i.e.,

$$\mathbf{G}_g = \left[\frac{1}{|\mathbf{V}_g|} \sum_{i \in \mathbf{V}_g} \mathbf{V}_g^i \parallel \frac{1}{|\mathbf{E}_g|} \sum_{(s,t) \in \mathbf{E}_g} \mathbf{E}_g^{s,t} \right], \quad (24)$$

where $|\mathbf{V}_g|$ is the number of nodes, $|\mathbf{E}_g|$ is the number of edges. Notably, this is different from the traditional GNN architecture whose global graph representation is only constructed on the readout of node features. Since our method can learn both node features and relation features, they are all leveraged to construct the global graph representation. This helps the graph representation embeds more useful relational information.

Node-level task layer. For node classification task, the prediction of node i is done as follows

$$\mathbf{y}^i = \mathbf{V}_g^i \mathbf{C}_V, \quad (25)$$

where $\mathbf{C}_V \in \mathbb{R}^{d_V \times n_V}$ is the node classifier, n_V is the number of node classes.

Edge-level task layer. For edge classification task, our method naturally makes prediction based on deep relation features E_g , formally written as

$$\mathbf{y}^{s,t} = E_g^{s,t} C_E, \quad (26)$$

where $C_E \in \mathbb{R}^{d_E \times n_E}$ is the edge classifier, n_E is the number of edge classes. This is better than traditional GNN works [3, 5, 8, 9] that concatenate the entity features as edge features, since the independent edge features (associated with the relational information) are more discriminative for edge-level tasks.

Graph-level task layer. For graph classification and regression tasks, we make the prediction based on the global graph representation, *i.e.*,

$$\mathbf{y} = G_g C_G, \quad (27)$$

where $C_G \in \mathbb{R}^{(d_v+d_E) \times n_G}$, n_G is the number of graph classes. If it is graph regression task, then $n_G = 1$.

B. Network Differentiation Details

An architecture is represented as a directed acyclic graph (DAG) with $\{\mathbb{V}, \mathbb{L}\}$, where $\mathbb{V} = \{\mathbf{X}_i\}$ is the set of feature vertices and $\mathbb{L} = \{e(\mathbf{X}_i, \mathbf{X}_j, O)\}$ is the set of directed links. Each directed link $e(\mathbf{X}_i, \mathbf{X}_j, O)$ can transform the features from \mathbf{X}_i to \mathbf{X}_j using an operation O , where O is either a specific operation \mathbf{o} or a mixture operation $\bar{\mathbf{o}}$. The mixture operation $\bar{\mathbf{o}}^{(i,j)}$ is parameterized by architectural parameters $\alpha^{(i,j)}$ as a softmax mixture over all the possible operations within the operation space \mathcal{O} , *i.e.*,

$$\bar{\mathbf{o}}^{(i,j)}(\mathbf{X}_i) = \sum_{\mathbf{o} \in \mathcal{O}} \frac{\exp(\alpha_{\mathbf{o}}^{(i,j)})}{\sum_{\mathbf{o}' \in \mathcal{O}} \exp(\alpha_{\mathbf{o}'}^{(i,j)})} \mathbf{o}(\mathbf{X}_i). \quad (28)$$

For each operation $\mathbf{o}^{(i,j)}$, it is associated with the network weights $\mathbf{w}^{(i,j)}$.

The network differentiation aims to differentiate the local supernets into several specific subnets. Specifically, after network division and before network differentiation, the current architecture contains two kinds of links $e(\mathbf{X}_i, \mathbf{X}_j, \mathbf{o})$ and $e(\mathbf{X}_i, \mathbf{X}_j, \bar{\mathbf{o}})$. After network differentiation, there is only one kind of links, *i.e.*, $e(\mathbf{X}_i, \mathbf{X}_j, \mathbf{o})$, where a valid architecture is obtained. This procedure can be implemented by some differentiable architecture search strategies (DARTS [7], SGAS [6], *etc.*). Taking DARTS as an example, the learning procedure of the architectural parameters involves a bi-level optimization problem, *i.e.*,

$$\min_{\mathcal{A}} \mathcal{L}_{val}(\mathcal{W}^*(\mathcal{A}), \mathcal{A}), \quad (29)$$

$$s.t. \mathcal{W}^*(\mathcal{A}) = \underset{\mathcal{W}}{\operatorname{argmin}} \mathcal{L}_{train}(\mathcal{W}, \mathcal{A}), \quad (30)$$

where \mathcal{L}_{train} and \mathcal{L}_{val} are the training and validation loss, respectively. \mathcal{W} is the set of network weights $\{\mathbf{w}^{(i,j)}\}$

and \mathcal{A} is the set of the architectural parameters $\{\alpha^{(i,j)}\}$. DARTS [7] proposes to solve the bi-level problem by a first/second order approximation. At the end of the search, the final architecture is derived by selecting the operation with the highest weight for each mixture operation, *i.e.*,

$$\bar{\mathbf{o}}^{(i,j)} \leftarrow \underset{\mathbf{o} \in \mathcal{O}}{\operatorname{argmax}} \alpha_{\mathbf{o}}^{(i,j)}. \quad (31)$$

C. Visualizing Hierarchical Features

In this section, we provide more examples of the learned relation and node features on ModelNet40. The visualization results are reported in Figure 1.

D. Searched Architectures

We illustrate our searched GNN architectures using the proposed *network proliferation search paradigm* in Figure 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20.

References

- [1] Maxim Berman, Hervé Jégou, Andrea Vedaldi, Iasonas Kokkinos, and Matthijs Douze. Multigrain: a unified image embedding for classes and instances. *ArXiv*, abs/1902.05509, 2019. 1
- [2] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Lio', and Petar Velickovic. Principal neighbourhood aggregation for graph nets. *NIPS*, 2020. 1
- [3] Vijay Prakash Dwivedi, Chaitanya K Joshi, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Benchmarking graph neural networks. *arXiv preprint arXiv:2003.00982*, 2020. 3
- [4] William L. Hamilton, Zhitaoying, and Jure Leskovec. Inductive representation learning on large graphs. In *NIPS*, 2017. 1
- [5] Thomas Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *ArXiv*, abs/1609.02907, 2017. 3
- [6] Guohao Li, Guocheng Qian, Itzel C Delgadillo, Matthias Muller, Ali Thabet, and Bernard Ghanem. Sgas: Sequential greedy architecture search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1620–1630, 2020. 2, 3
- [7] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *ArXiv*, abs/1806.09055, 2019. 3
- [8] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio', and Yoshua Bengio. Graph attention networks. *ArXiv*, abs/1710.10903, 2018. 3
- [9] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *ArXiv*, abs/1810.00826, 2019. 1, 3

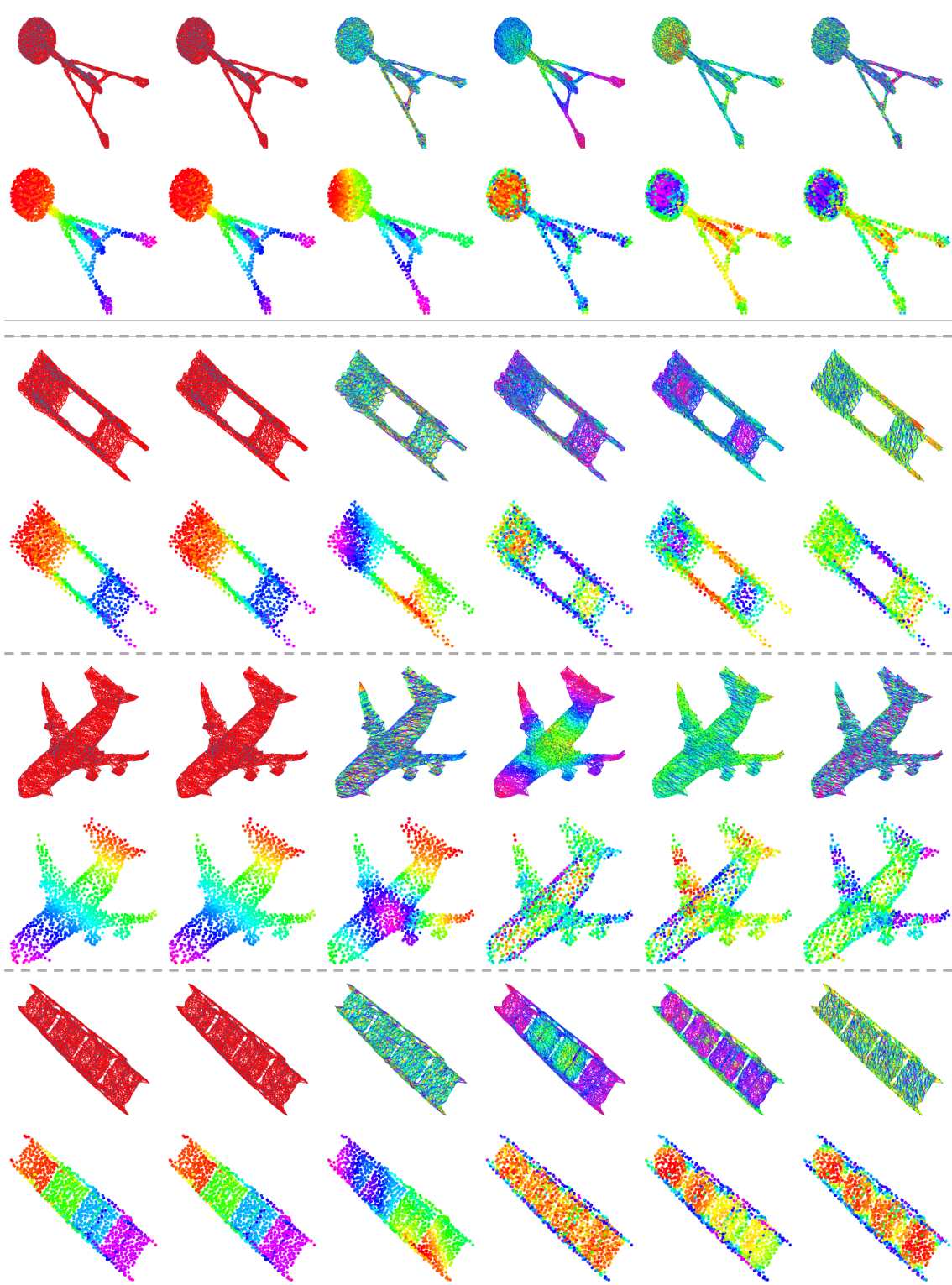


Figure 1. Visualization of the learned hierarchical relation and node features for 3D point cloud recognition. Relation features with different edge color distributions have different message-passing preferences. Node features with different node color distributions represent different clustering effects.

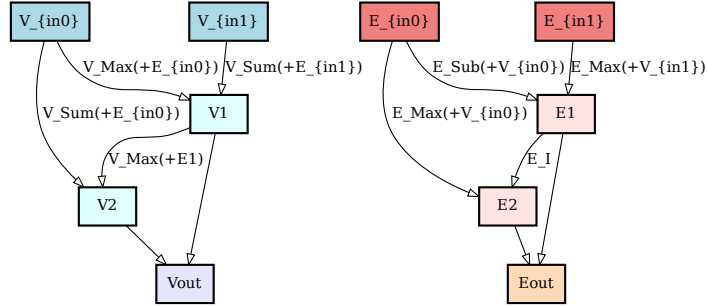


Figure 2. Illustration of the searched architecture with the size of 2 on the ZINC dataset.

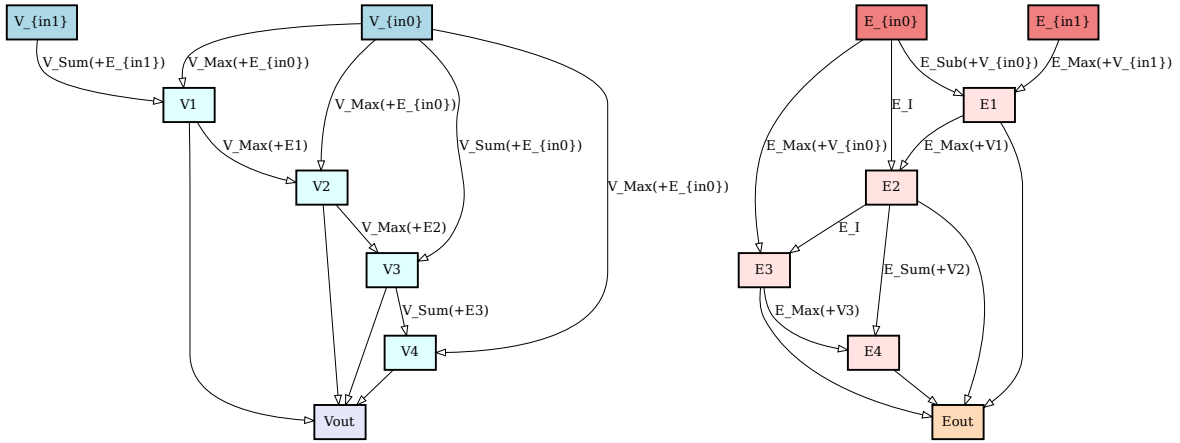


Figure 3. Illustration of the searched architecture with the size of 4 on the ZINC dataset.

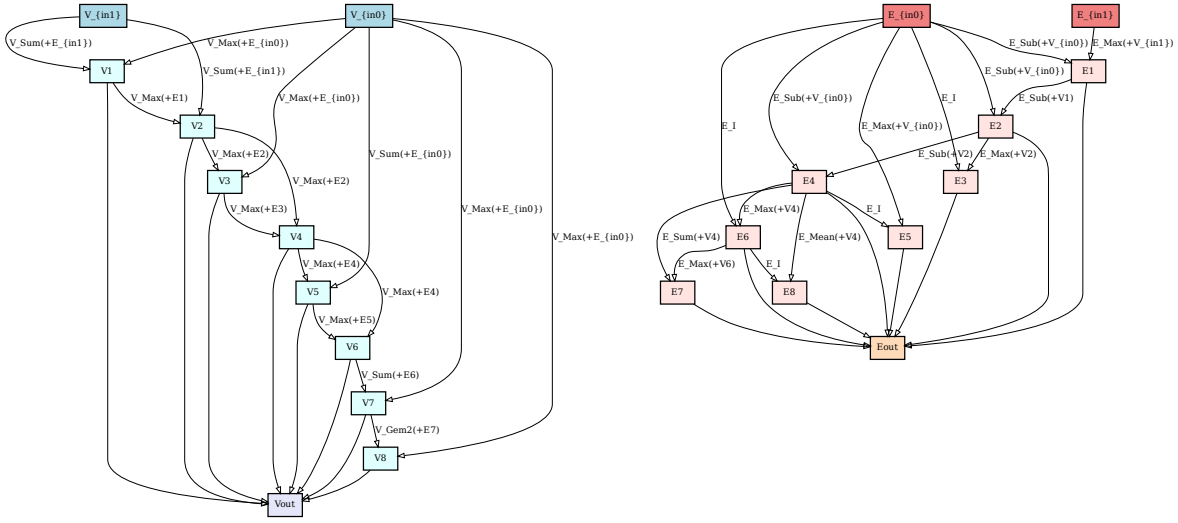


Figure 4. Illustration of the searched architecture with the size of 8 on the ZINC dataset.

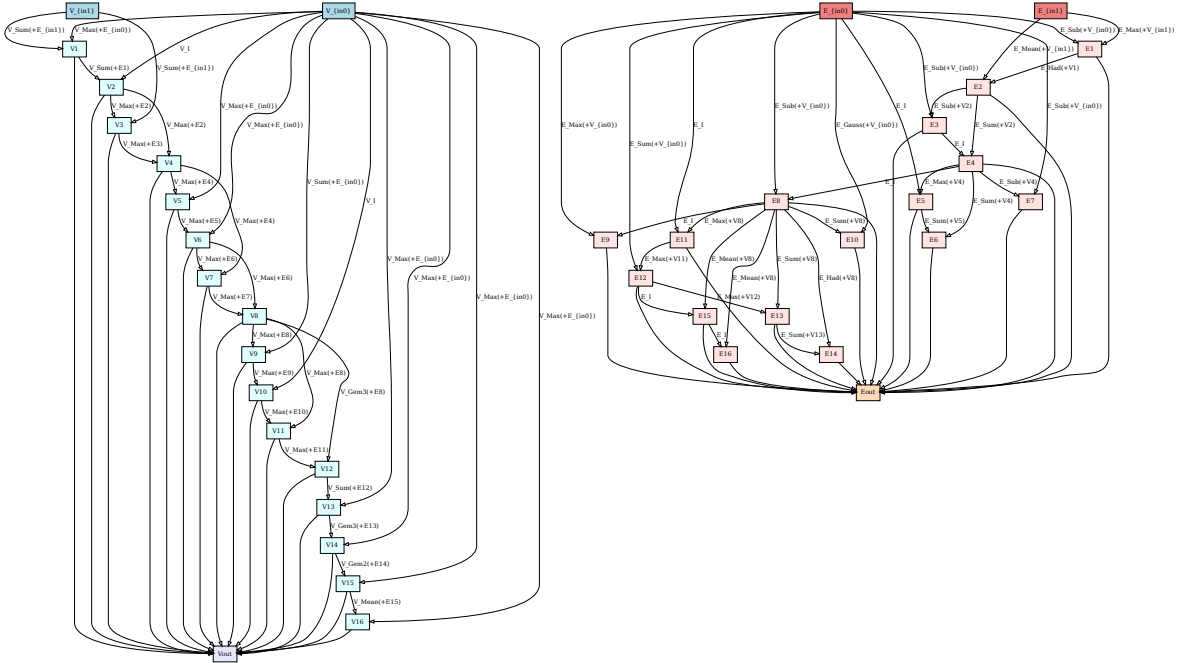


Figure 5. Illustration of the searched architecture with the size of 16 on the ZINC dataset.

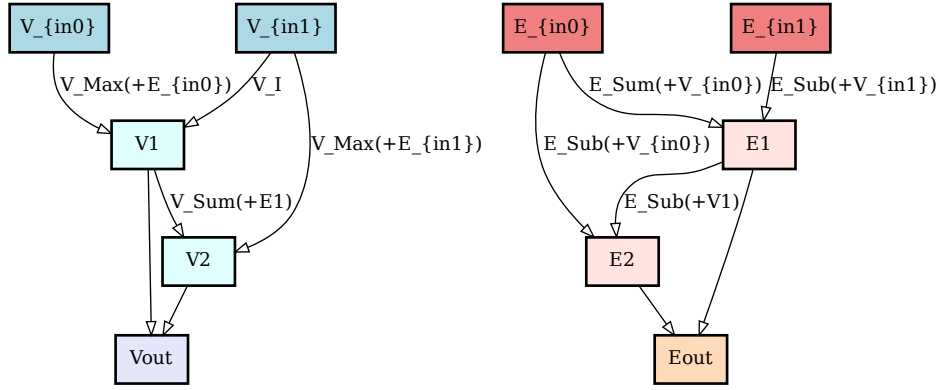


Figure 6. Illustration of the searched architecture with the size of 2 on the CLUSTER dataset.

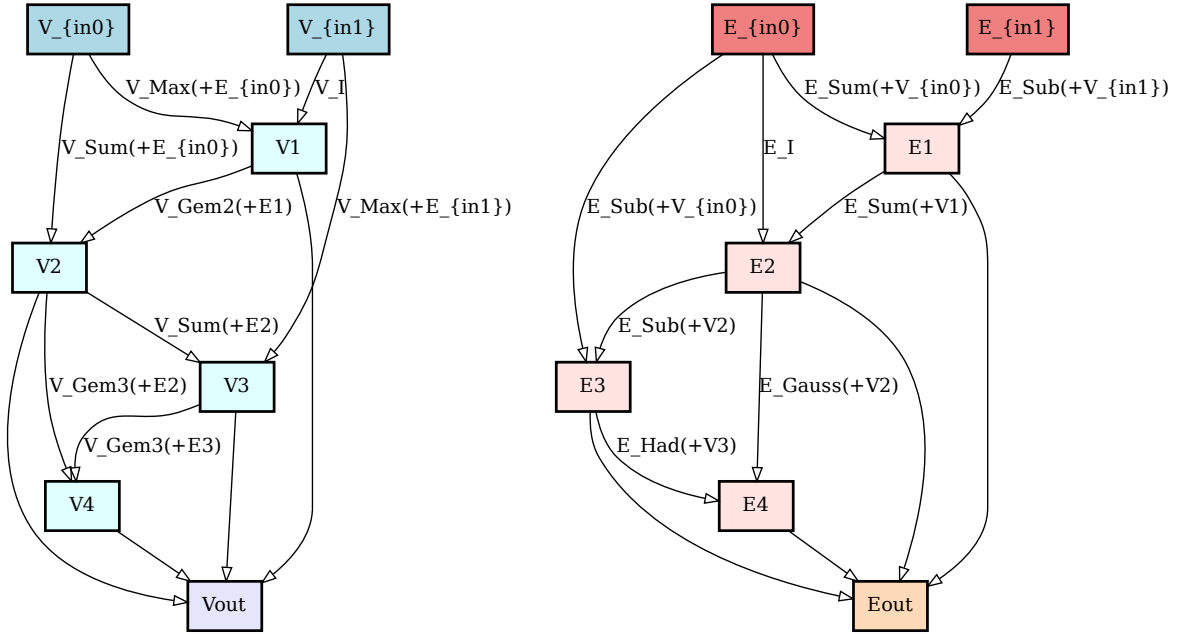


Figure 7. Illustration of the searched architecture with the size of 4 on the CLUSTER dataset.

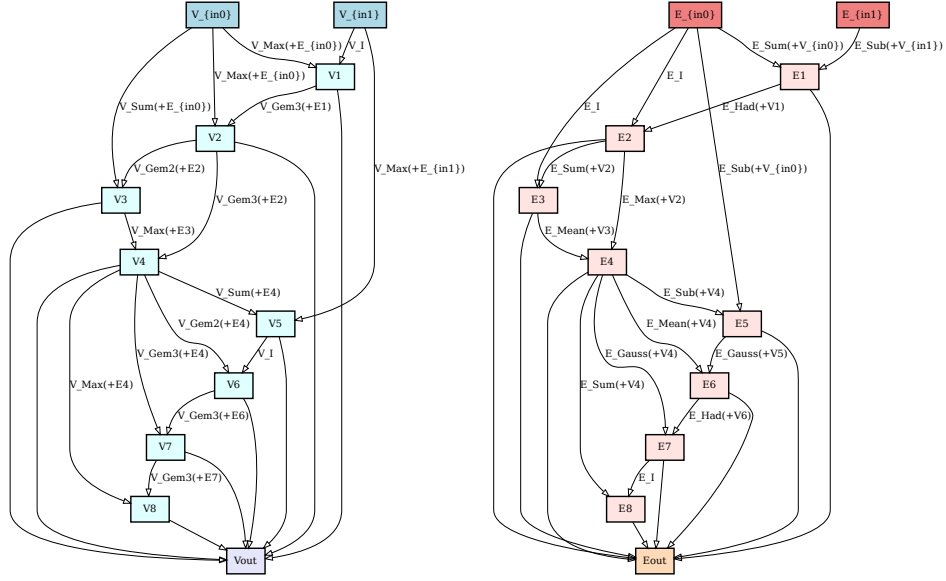


Figure 8. Illustration of the searched architecture with the size of 8 on the CLUSTER dataset.

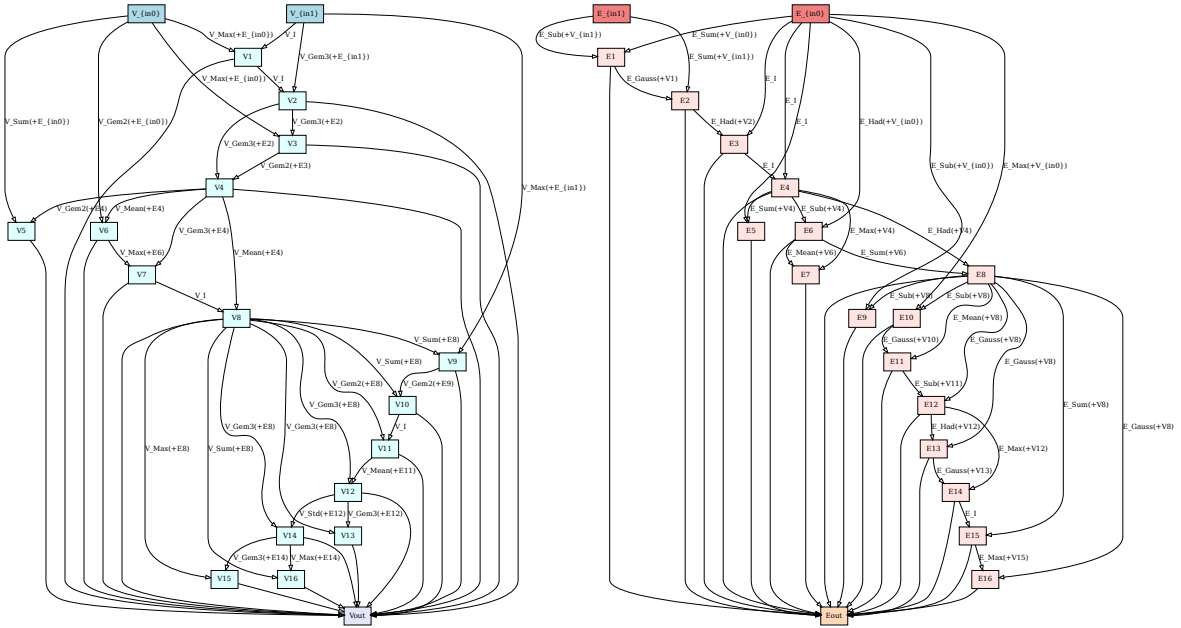


Figure 9. Illustration of the searched architecture with the size of 16 on the CLUSTER dataset.

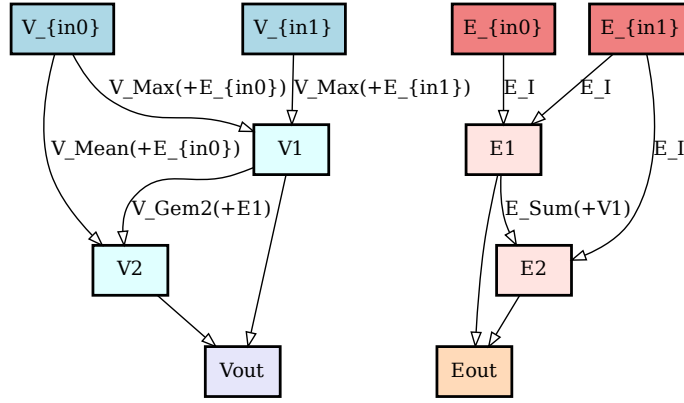


Figure 10. Illustration of the searched architecture with the size of 2 on the TSP dataset.

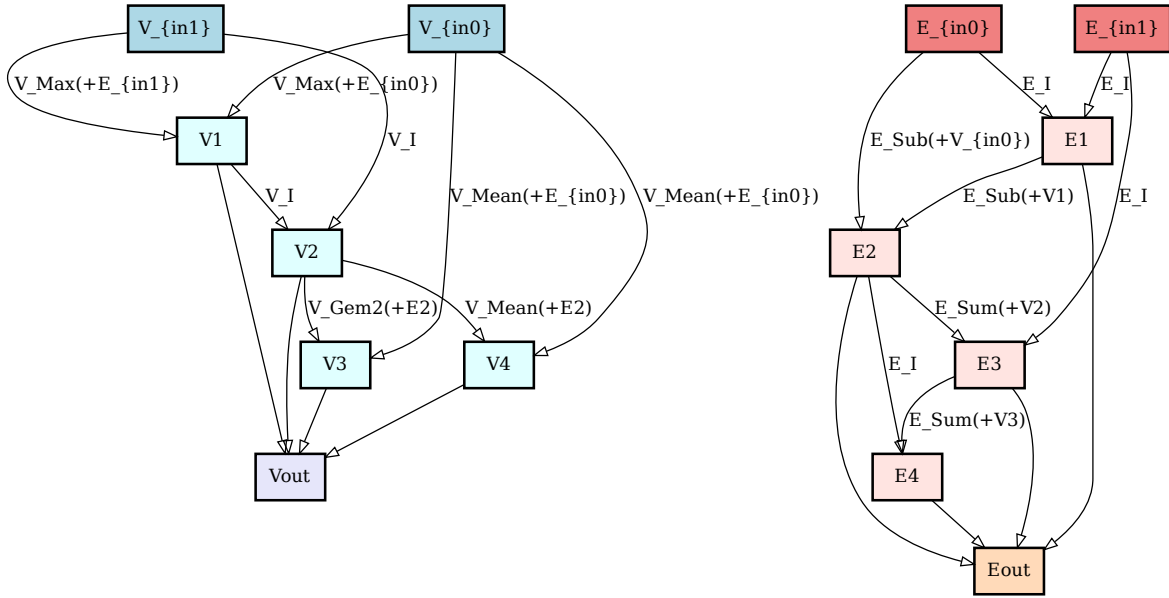


Figure 11. Illustration of the searched architecture with the size of 4 on the TSP dataset.

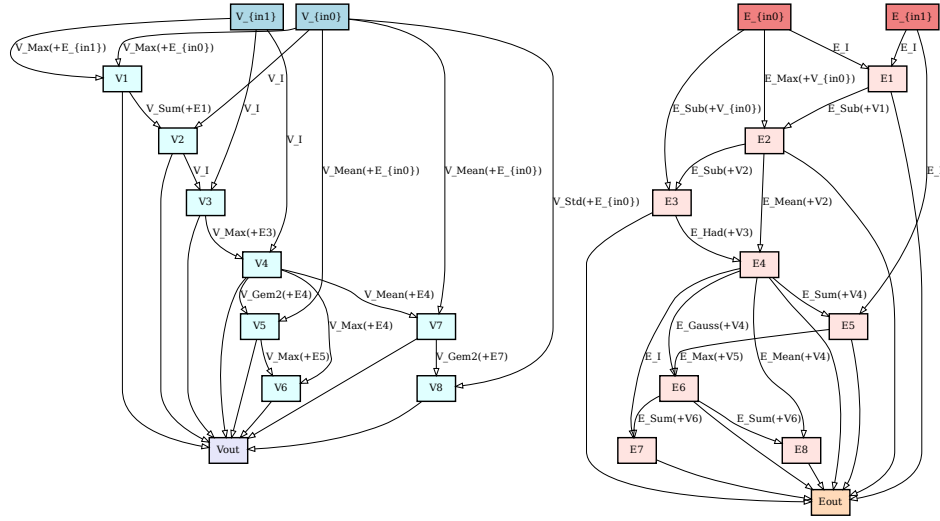


Figure 12. Illustration of the searched architecture with the size of 8 on the TSP dataset.

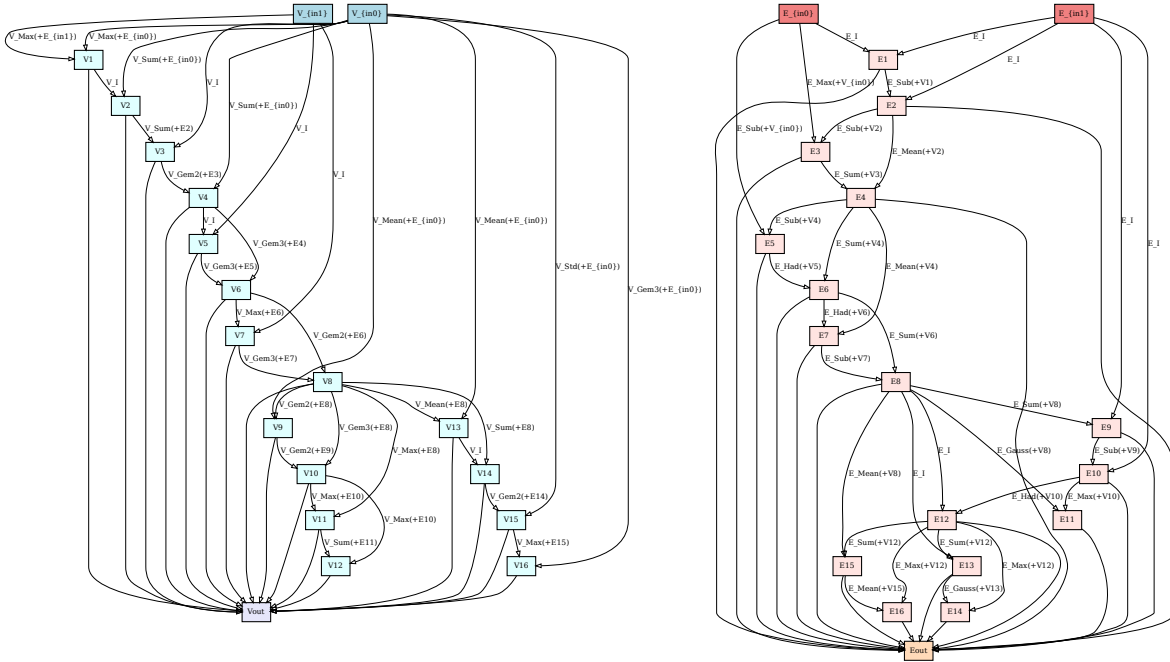


Figure 13. Illustration of the searched architecture with the size of 16 on the TSP dataset.

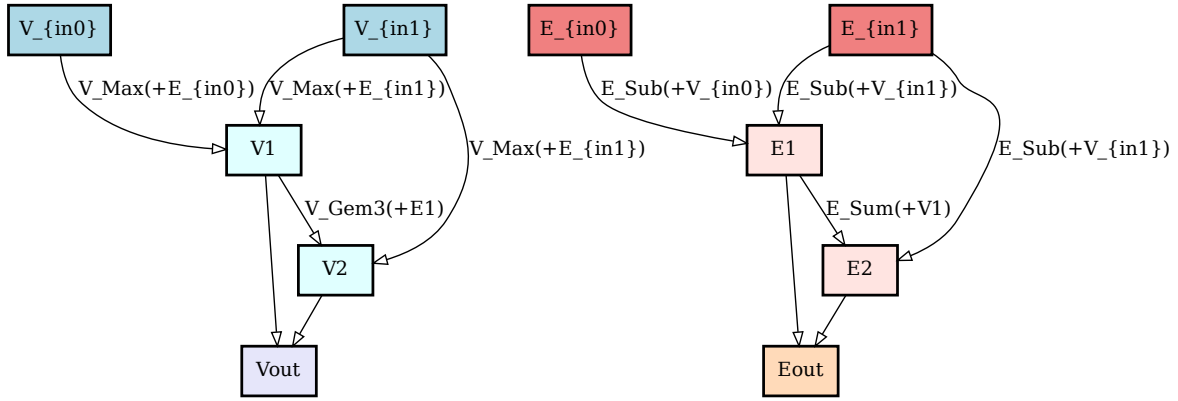


Figure 14. Illustration of the searched architecture with the size of 2 on the CIFAR10 dataset.

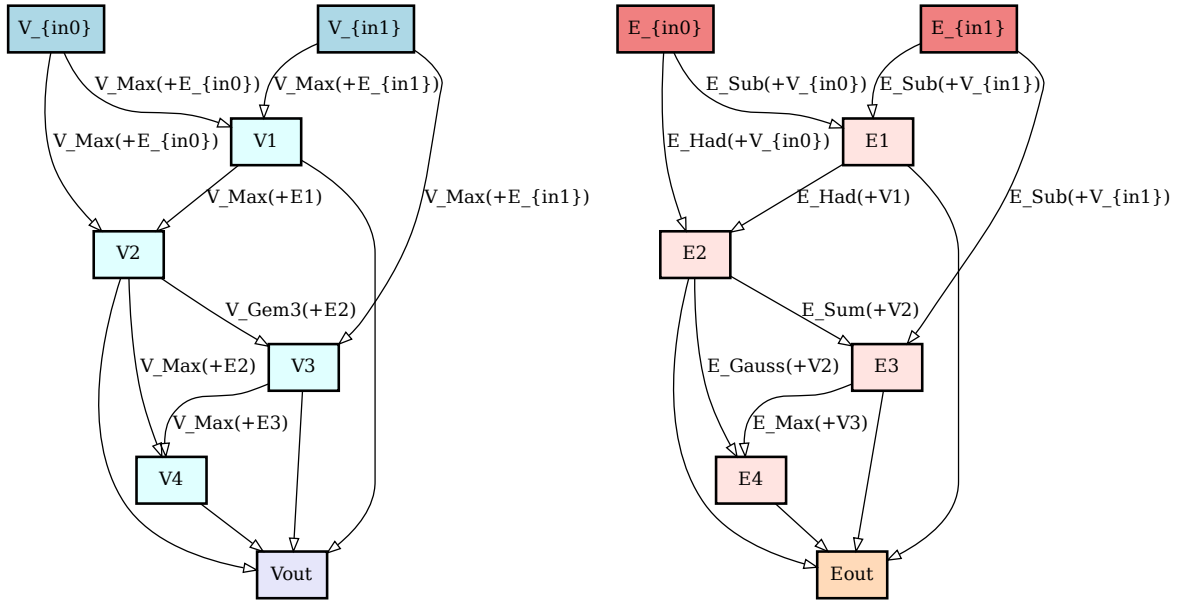


Figure 15. Illustration of the searched architecture with the size of 4 on the CIFAR10 dataset.

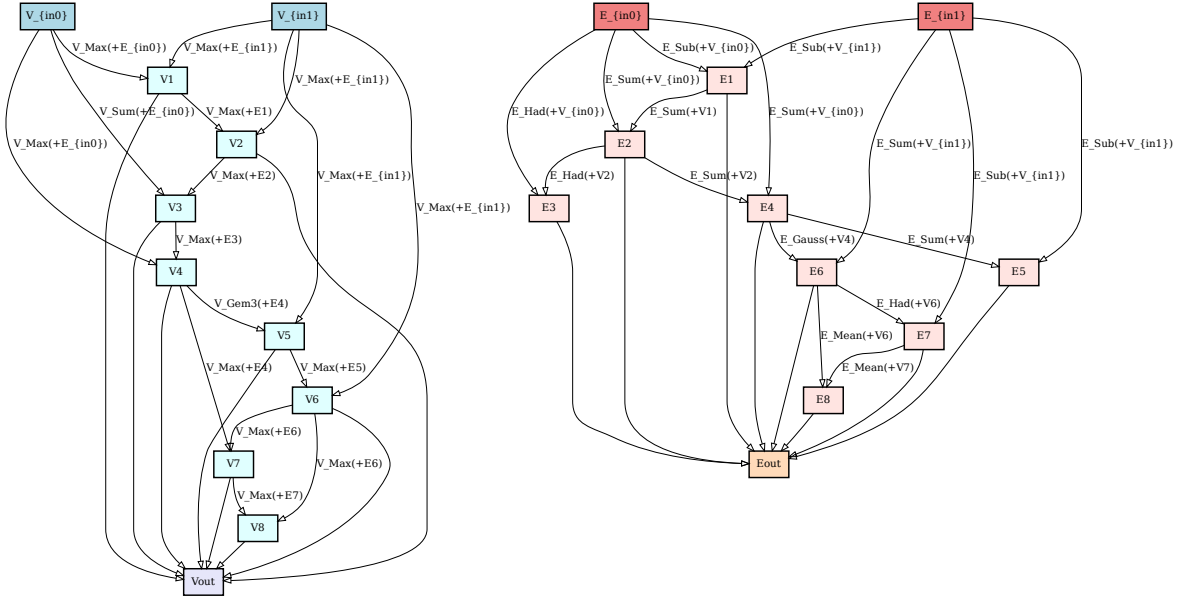


Figure 16. Illustration of the searched architecture with the size of 8 on the CIFAR10 dataset.

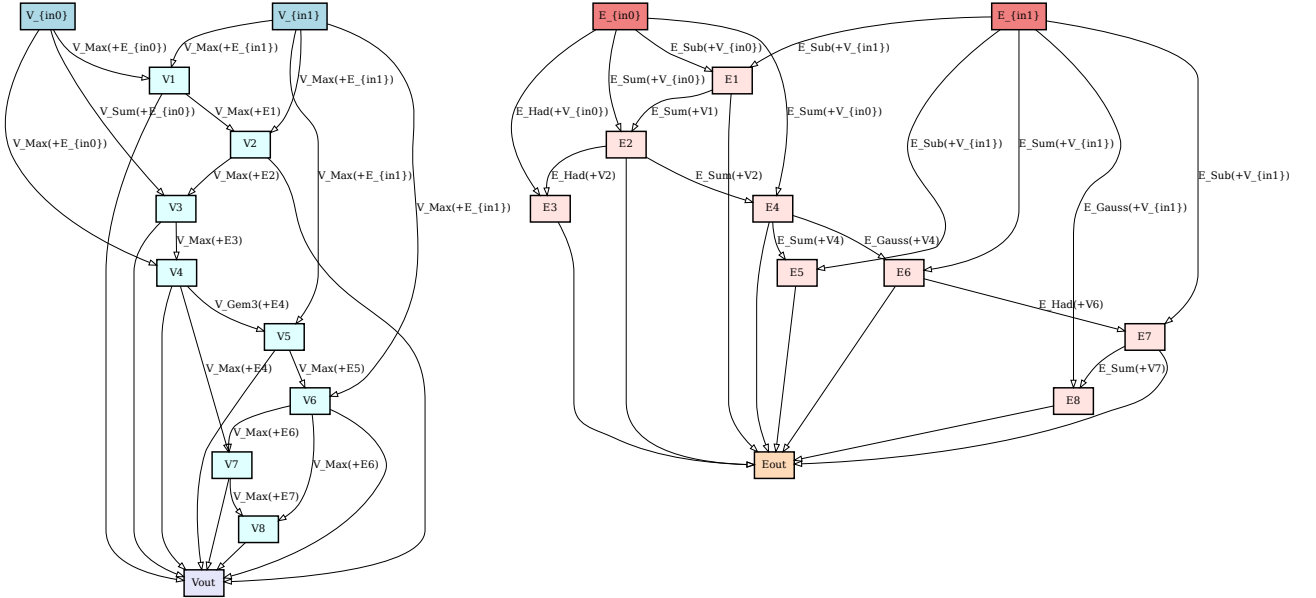


Figure 17. Illustration of the searched architecture with the size of 16 on the CIFAR10 dataset.

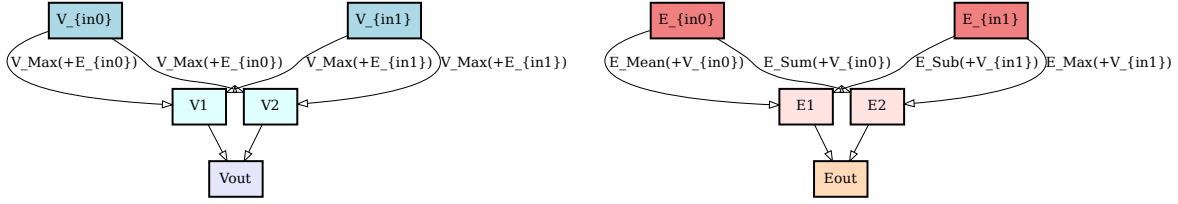


Figure 18. Illustration of the searched architecture with the size of 2 on the ModelNet10 dataset.

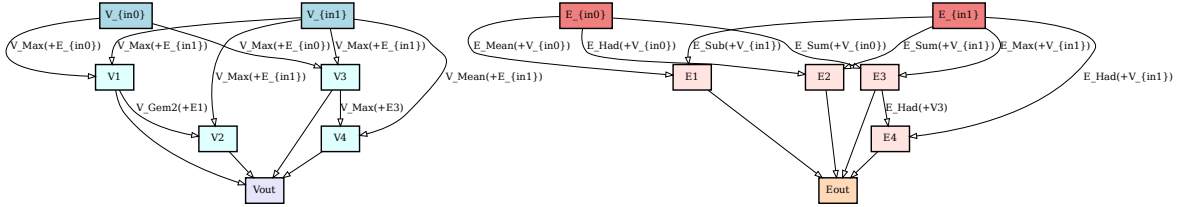


Figure 19. Illustration of the searched architecture with the size of 4 on the ModelNet10 dataset.

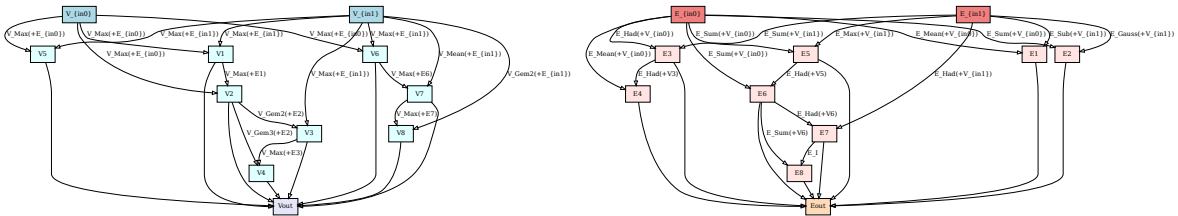


Figure 20. Illustration of the searched architecture with the size of 8 on the ModelNet10 dataset.