# Point2Node: Correlation Learning of Dynamic-Node for Point Cloud Feature Modeling

**Wenkai Han,**[1] **Chenglu Wen,**[1*] **Cheng Wang,**[1] **Xin Li,**[2] **Qing Li**[1]

[1]School of Informatics, Xiamen University
422 Siming South Road, Xiamen 361005, China
[2]School of Electrical Engineering and Computer Science, Louisiana State University
Baton Rouge, LA 70803, USA
hlxwk0525@gmail.com, clwen@xmu.edu.cn, cwang@xmu.edu.cn, xinli@cct.lsu.edu, hello.qingli@gmail.com

## Abstract

Fully exploring correlation among points in point clouds is essential for their feature modeling. This paper presents a novel end-to-end graph model, named Point2Node, to represent a given point cloud. Point2Node can dynamically explore correlation among all graph nodes from different levels, and adaptively aggregate the learned features. Specifically, first, to fully explore the spatial correlation among points for enhanced feature description, in a high-dimensional node graph, we dynamically integrate the node's correlation with self, local, and non-local nodes. Second, to more effectively integrate learned features, we design a data-aware gate mechanism to self-adaptively aggregate features at the channel level. Extensive experiments on various point cloud benchmarks demonstrate that our method outperforms the state-of-the-art.

## Introduction

The recent advance of 3D technologies in autonomous driving, robotics, and reverse engineering has attracted greater attention in understanding and analyzing 3D data. Often in the form of point clouds, 3D data is a set of unordered 3D points, sometimes with additional features (e.g., RGB, normal) defined on each point. Due to the inherent unorderness and irregularity property of 3D data, it is difficult to extend traditional convolutional neural networks to 3D data, because a permutation in the order of 3D points should not affect their geometric structures.

Multiple approaches (Qi et al. 2017a; Qi et al. 2017b; Wang et al. 2018c; Zhao et al. 2019) have been developed recently to process raw point clouds. These approaches can be grouped into several categories according to the point correlation they exploit for feature learning: *(1) Self correlation.* A seminal work on point cloud feature learning is the PointNet (Qi et al. 2017a), which uses shared Multi-Layer Perceptions (MLPs) on each point individually, then gathers global representation with pooling functions. Its description ability is limited by not being able to incorporate contextual information from each point's neighbors. *(2)*

*Local correlation.* PointNet++ (Qi et al. 2017b) improves the PointNet by applying the PointNet recursively on partitioned pointsets and encoding local features with max pooling. PointNet++ essentially treats each point individually and hence still neglects its correlation between neighbors. Later, DGCNN (Wang et al. 2018c) employs shared MLPs to process edges linking the centroid point and its neighboring points. But because the edges only connect to the centroid point, the neighboring correlation is not fully exploited. The more recent approach, PointWeb (Zhao et al. 2019), considers a point's correlation by weighting all edges among its neighbors. However, although PointWeb explores the complete correlation from a point's neighbors, it stores redundant connections between two points of different labels. Furthermore, to our best knowledge, in all these existing approaches, most existing approaches do not consider non-local correlation among distant points that are not neighbors but could potentially exhibit long-range dependencies.

Representing point cloud with the graph structure is appropriate for correlation learning. There are some graph-based methods (Simonovsky and Komodakis 2017; Verma, Boyer, and Verbeek 2017) that use graph convolutions to analyze point cloud on local regions; however, explicit correlations of the neighboring points are not well captured by the predefined convolution operations and this limits their feature modeling capability.

To better leverage correlations of points from different levels, we formulate a novel high-dimensional node graph model, named Point2Node, to learn correlation among 3D points. Point2Node proposes a Dynamic Node Correlation module, which constructs a high-dimensional node graph to sequentially explore correlations at different levels. Specifically, it explores *self correlation* between point's different channels, *local correlation* between neighboring points, and also, *non-local correlation* among distant points having long-range dependency. Meanwhile, we dynamically update nodes after each correlation learning to enhance node characteristics. Furthermore, for aggregating more discriminate features, our Point2Node employs a data-aware gate embedding, named Adaptive Feature Aggregation, to pass or block channels of high-dimensional nodes from different correlation learning, thereby generating new nodes to balance self,

local and non-local correlation. Experimental results on various datasets show better feature representation capability. Our contributions are summarized as follows:

- We propose a Point2Node model to dynamically reasons the different-levels correlation including self correlation, local correlation, and non-local correlation. This model largely enhances nodes characteristic from different scale correlation learning, especially embedding long-range correlation.

- We introduce a data-aware gate mechanism, which self-adaptively aggregates features from different correlation learning and maintains effective components in different levels, resulting in more discriminate features.

- Our end-to-end Point2Node model achieves state-of-the-art performance on various point cloud benchmark, thus demonstrating its effectiveness and generalization ability.

## Related work

**View-based and volumetric methods.** A classic category of deep-learning based 3D representations is to use the multi-view approach (Qi et al. 2016; Su et al. 2015; Kalogerakis et al. 2017), where a point cloud is projected to a collection of images rendered from different views. By this means, point data can be recognized using conventional convolution in 2D images. However, these viewed-based methods only achieve a dominant performance in classification. Scene segmentation is nontrivial due to the loss of the discriminate geometric relationship of the 3D points during 2D projection.

Influenced by the success of traditional convolution in regular data format (e.g., 2D image), another strategy is to generalize CNN to 3D voxel structures (Maturana and Scherer 2015; Wu et al. 2015; Dai et al. 2017), then assign 3D points to regular volumetric occupancy grids. Due to the intractable computational cost, these methods are usually constrained by low spatial resolution. To overcome this limitation, Kd-Net (Klokov and Lempitsky 2017) and Oct-Net (Riegler, Ulusoy, and Geiger 2017) skip the computation on empty voxels so that they can handle higher-resolutional grids. However, during voxelization, it is inevitable to lose geometric information. To avoid geometric information loss, unlike view-based and volumetric methods, our method processes 3D point clouds directly.

**Point cloud-based methods.** The pioneering method in this category is PointNet (Qi et al. 2017a), which considers each point independently with shared MLPs for permutation invariance and aggregates the global signature with max-pooling. PointNet is limited by the fine-grained information in local regions. PointNet++ (Qi et al. 2017b) extends Point-Net by exploiting a hierarchical structure to incorporate local dependencies with Farthest Point Sampling (FPS) and a geometric grouping strategy. DGCNN (Wang et al. 2018c) proposes an edge-conv operation which encodes the edges between the center point and k nearest neighbors with shared MLPs. PointWeb (Zhao et al. 2019) enriches point features by weighting all edges in local regions with point-pair difference. Our method also exploits local edges, but we explore the adaptive connection of similar nodes in local regions, to

avoid redundant connections between two points of different labels, and this leads to better shape-awareness.

A few other methods extend the convolution operator to process unordered points. PCCN (Wang et al. 2018a) proposes parametric continuous convolutions that exploit the parameterized kernel functions on point clouds. By reordering and weighting local points with an $\chi$-transformation, PointCNN (Li et al. 2018) exploits the canonical order of points and then applies convolution. PointConv (Wu, Qi, and Li 2019) extends the dynamic filter to a new convolution operation and build deep convolutional networks for point clouds. A-CNN (Komarichev, Zhong, and Hua 2019) employs annular convolution to capture the local neighborhood geometry of each point. Unlike these methods that use convolution in point clouds, we pay our attention to correlation learning among 3D points in point clouds.

**Graph-based methods.** Because it is flexible to represent irregular point clouds with graphs, there is an increasing interest in generalizing convolutions to the graphs. In line with the domain in which convolutions operate, graph-based methods can be categorized into spectral and non-spectral methods. Spectral methods (Kipf and Welling 2016; Yi et al. 2017) perform convolution in the spectral representation, which is constructed by computing with a Laplacian matrix. The major limitation of these methods is that they require a fixed graph structure. Thus they have difficulty in adapting to point clouds that have varying graph structures. Instead, we perform correlation learning on dynamic nodes, which are not limited by varying graph structures. Alternatively, non-spectral methods (Simonovsky and Komodakis 2017; Atwood and Towsley 2016; Verma, Boyer, and Verbeek 2017) perform convolution-like operations on the spatially close nodes of a graph, thereby propagating information along the neighboring nodes. But, hand-crafted ways of encoding local contextual information limits their feature modeling capability; instead, for more fine-grained contextual information in the local region, we explore the direct correlation between local nodes.

## Method

### Overview

We illustrate our end-to-end Point2Node framework in Fig. 1. Given a set of unordered points $P = \{p_i \in R^{3+d}, i = 1, 2, ..., N\}$, where $N$ is the number of points and $3 + d$ denotes the $xyz$-dimension and additional properties (e.g., $d = 3$ for RGB or normal information), through *Feature Extraction*, we construct their high-level representation, forming high-dimensional nodes $V^{(0)} = \{v_i^{(0)} \in R^C, i = 1, 2, ..., N\}$, where each node $v_i^{(0)}$ is a learned multi-channel feature vector of point $p_i$, and $C$ is the dimension of each node. The *Feature Extraction* module employs a $\chi$-conv operator from PointCNN (Li et al. 2018) and builds a U-Net (zgn Çiçek et al. 2016) structure for feature learning. Then through a *Dynamic Node Correlation* (DNC) module, we explore point correlations from different levels, and dynamically generate high-dimensional nodes $V^{(1)}$, $V^{(2)}$ and $V^{(3)}$, respectively. Simultaneously, we build an *Adaptive*
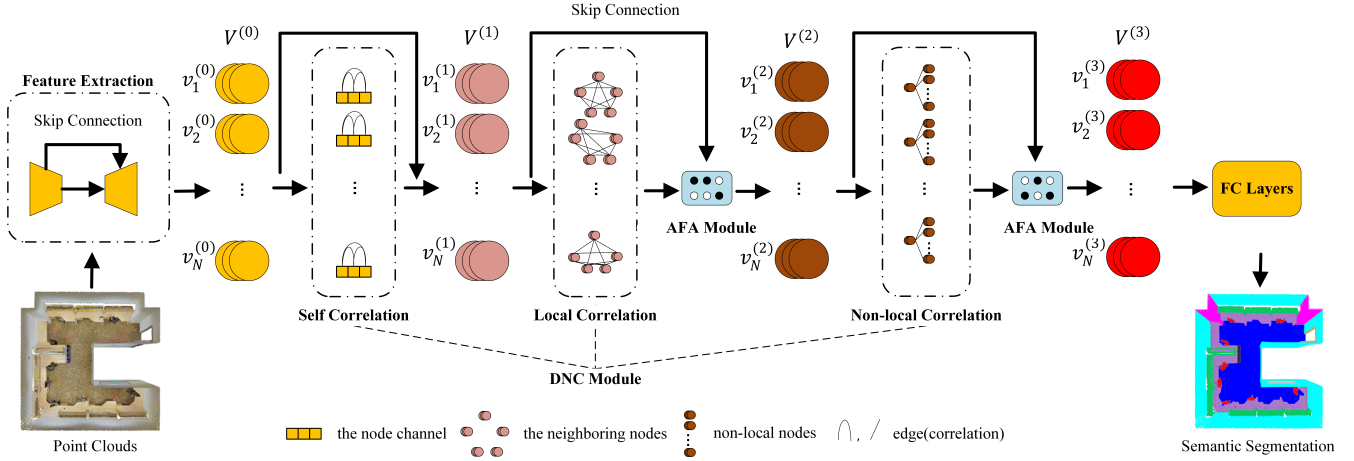
Figure 1: Our Point2Node architecture. First, Point2Node converts input points into high-level representation through Feature Extraction, then constructs high-dimension nodes $V^{(0)}$. Second, the Dynamic Node Correlation (DNC) module (which consists of self correlation, local correlation and non-local correlation) and the Adaptive Feature Aggregation (AFA) module are employed to enhance the node characteristics, resulting in different nodes (different colors). Finally, the nodes are assigned different labels through FC Layers.

*Feature Aggregation* (AFA) module to self-adaptively aggregate the features of nodes learned in different levels. Finally, the high-dimensional nodes are fed into fully-connected layers for label assignment.

## Dynamic Nodes Correlation Module

To fully explore correlation among the points from different levels, we propose a *Dynamic Nodes Correlation* module. Given the initial high-dimensional nodes $V^{(0)} = \{v_i^{(0)} \in R^C, i = 1, 2, ..., N\}$, we introduce the *cross-channel* edges and *cross-node* edges. *Cross-channel edges* are defined between two channels of the same node, and they are used to explore correlation between this node's different channels. *Cross-node edges* are defined between different nodes, and have two types. Between neighboring (spatially close) nodes, *local edges* are defined to explore local contextual information; while between non-local (spatially distant) nodes, *non-local edges* are defined to exploit latent long-range dependencies. Therefore, three types of correlation: *Self Correlation*, *Local Correlation*, and *Non-local Correlation* will be learned in this DNC module. Meanwhile, we dynamically update the high-dimensional nodes after each correlation learning. The learning of these three types of correlation is elaborated in the following sections.

## Self Correlation

The dependencies among channels of the same node are denoted by *self correlation*. Such dependencies are informative for feature representation. To model self correlation of each node, we can construct edges between different channels. However, for each channel, we do not compute its correlation weights w.r.t. other channels, but make a simplification. We use convolution to integrate these weights from different channels and encode them into a single weight.

**Channel Weights.** Inspired by traditional 1D convolution that aggregates all channels of one pixel in a 2D image, we similarly do such an encoding to integrate each channel's correlations into one weight through Multi-Layer Perception (MLP). To reduce the number of parameters, we implement the MLP with 1D separable convolutions (Chollet 2017). Consequently, for each $C$-dimensional node $v_i^{(0)} \in V^{(0)}$, we can calculate a $C$-dimensional channel weight $w_i \in R^C$

$$w_i = MLP(v_i^{(0)}). \tag{1}$$

And to make weight coefficients comparable between different channels, we perform a normalization using $softmax$,

$$\overline{w_{i,c}} = \frac{exp(w_{i,c})}{\sum_{j=1}^{C} exp(w_{i,j})}, \tag{2}$$

where $w_{i,c}$ is the $c$-th component of $w_i$, and $\overline{w_{i,c}}$ is the normalized $w_{i,c}$.

**Updating Nodes.** After correlation learning, we update nodes with residual connections:

$$v_i^{(1)} = v_i^{(0)} \oplus \alpha(\overline{w_i} \odot v_i^{(0)}), \tag{3}$$

where $\oplus$ and $\odot$ denote channel-wise summation and multiplication, respectively, and $\alpha$ denotes a learnable scale parameter. This update transforms nodes $V^{(0)}$ into new nodes $V^{(1)} = \{v_i^{(1)} \in R^C, i = 1, 2, ..., N\}$.

Note that unlike directly transforming $V^{(0)}$ to $V^{(1)}$ by convolutions, which share weights among all the nodes, our *self correlation* module learns different correlations in different nodes and can embody different nodes' individual characteristics.

## Local Correlation

To fully capture local contextual information, which is of central importance in point cloud understanding, we also

build a new *Local Correlation* sub-module. Existing approaches (Wang et al. 2018c; Zhao et al. 2019) usually construct partial or complete edges of local regions for contextual information, but these tend to store redundant connections between two points of different labels. In contrast, our sub-module adaptively connects locally similar nodes to construct a local graph, whose adjacency matrix weights correlation of relevant node pairs.

**Local Graph.** Given $V^{(1)} = \{v_i^{(1)} \in R^C, i = 1, 2, ..., N\}$, we consider each node as a centroid and construct a local graph. To increase the receptive field without redundant computational cost, we employ the dilated KNN (Li et al. 2018) to search for neighboring nodes, and form *Dilated K-Nearest Neighbor Graphs* (DKNNG) $G^{(l)} = \{G_i^{(l)} \in R^{K \times C}, i = 1, 2, ..., N\}$ in local regions. Specifically, we sample $K$ nodes at equal intervals from the top $K \times d$ neighbor nodes, where $K$ denotes the number of samples, and $d$ is the dilate rate.

**Adjacency Matrix.** For local correlation learning of node $v_i^{(1)}$, given its DKNNG $G_i^{(l)}$ where $V_i^{(l)} = \{V_i^{(l)}(k) \in R^C, k = 1, 2, ..., K\}$ is the neighboring nodes set of centroid node $v_i^{(1)}$, we construct its adjacency matrix $m_i^{(l)} \in R^{K \times K}$ for all the directed edges in $G_i^{(l)}$. Inspired by the Non-Local Neural Network (Wang et al. 2018b), which has been proven to be effective in capturing pixel correlation in 2D images, we generalize its non-local operation to our graph structure, to generate the adjacency matrix. Please refer to Appendix A for more detail about our generalization.

Specifically, we first map $V_i^{(l)}$ to low-dimension feature space with two 1-D convolutional mapping functions, $\theta^{(l)}$ and $\varphi^{(l)}$, resulting in low-dimension node representations $Q_i^{(l)} \in R^{K \times \frac{C}{r}}$ and $X_i^{(l)} \in R^{K \times \frac{C}{r}}$ respectively,

$$Q_i^{(l)} = \theta^{(l)}(V_i^{(l)}), X_i^{(l)} = \varphi^{(l)}(V_i^{(l)}), \qquad (4)$$

where $r$ is the reduction rate and is a shared super-parameter in the full paper. With $Q_i^{(l)}$ and $X_i^{(l)}$, we compute $m_i^{(l)}$ by matrix multiplication,

$$m_i^{(l)} = Q_i^{(l)} \otimes (X_i^{(l)})^T, \qquad (5)$$

where $\otimes$ denotes matrix multiplication.

Next, we prune the adjacency matrix $m_i^{(l)}$ with the $softmax$ function (Eq.6) to cut off redundant connections that link nodes of different labels, resulting in $\overline{m_i^{(l)}}$.

$$\overline{m_i^{(l)}}(x, y) = \frac{exp(m_i^{(l)}(x, y))}{\sum_{k=1}^{K} exp(m_i^{(l)}(x, k))}, \qquad (6)$$

where $\overline{m_i^{(l)}}(x, y)$ is the masked $m_i^{(l)}(x, y)$, and the final weight of the edge that links the $x$-th node and $y$-th node in the neighboring nodes set $V_i^{(l)}$.

**Update Nodes.** We simultaneously update the neighboring nodes $V_i^{(l)}$ using the correlation function $\overline{m_i^{(l)}}$ (Eq. 7) Finally, we encode learned local contextual information to the centroid node with $Max\_Pooling$ function (Eq. 8).

$$V_i^{(l)}{}_{up} = \overline{m_i^{(l)}} \otimes V_i^{(l)}. \qquad (7)$$

$$v_i^{(1)}{}_{up} = Max\_Pooling(V_i^{(l)}{}_{up}), \qquad (8)$$

where $V_i^{(l)}{}_{up}$ and $v_i^{(1)}{}_{up}$ are updated neighboring nodes and centroid node, respectively.

## Non-local Correlation

To better exploit latent long-range dependency, which is essential in capturing global characteristics of point clouds, we also propose a *Non-local Correlation* sub-module. Traditional point cloud networks usually construct global representations by stacking multi-layers local operations, but such a high-level global representation can hardly capture direct correlation between spatially distant nodes. To model such non-local correlation, our proposed module adaptively connects non-local nodes to construct a global graph, whose adjacency matrix reveals latent long-range dependencies between nodes.

**Adjacency Matrix.** Given nodes $V^{(2)} = \{v_i^{(2)} \in R^C, i = 1, 2, ..., N\}$, we construct a global graph $G^{(g)}$ and an adjacency matrix $m^{(g)} \in R^{N \times N}$. $m^{(g)}$ is constructed in a similar way to *local correlation*. We construct the low-dimensional $Q^{(g)} \in R^{N \times \frac{C}{r}}$ and $X^{(g)} \in R^{N \times \frac{C}{r}}$ from mapping functions $\theta^{(g)}$ and $\varphi^{(g)}$,

$$Q^{(g)} = \theta^{(g)}(V^{(2)}), X^{(g)} = \varphi^{(g)}(V^{(2)}), \qquad (9)$$

then compose the adjacency matrix by

$$m^{(g)} = Q^{(g)} \otimes (X^{(g)})^T. \qquad (10)$$

Unlike *local correlation*, the non-local correlation learning can explore dependence between nodes that have different features in different scales. Therefore, to make weight coefficients comparable, we perform a normalization on the adjacency matrix $\overline{m^{(g)}}$,

$$\overline{m^{(g)}}(x, y) = \frac{exp(m^{(g)}(x, y))}{\sum_{k=1}^{N} exp(m^{(g)}(x, k))}, \qquad (11)$$

where $\overline{m^{(g)}}(x, y)$ is the normalized weight of the edge linking nodes $x$ and $y$.

**Updating Nodes.** We update all the nodes $V^{(2)}$ with the correlation function $\overline{m^{(g)}}$, resulting in $V^{(2)}{}_{up}$:

$$V^{(2)}{}_{up} = \overline{m^{(g)}} \otimes V^{(2)}. \qquad (12)$$

Note that all the three aforementioned sub-modules perform node-wise operations. These operations are permutation-equivalent (See a proof in Appendix B) and hence can effectively process unordered 3D points.

## Adaptive Feature Aggregation

Traditional neural networks aggregate different features using linear aggregation such as skip connections. But linear aggregations simply mix up channels of the features and are not data-adaptive. To better reflect the characteristics of given point cloud and make the aggregation data-aware, we

propose an *Adaptive Feature Aggregation* (AFA) module [1] and use a gate mechanism to integrate node features according to their characteristics. To model the channel-wise characteristics and their relative importance, we consider two models, namely, the parameter-free to parameterized descriptions.

**Parameter-free Characteristic Modeling.** Given high-dimensional nodes $V^{(1)}$ and $V^{(1)}_{up}$ that are from different correlation learning, we first describe their characteristics using a global node average pooling, and obtain channel-wise descriptors $s_1 \in R^{1 \times C}$ and $s_2 \in R^{1 \times C}$:

$$s_{1,c} = \frac{1}{N} \sum_{i=1}^{N} v_{i,c}^{(1)}, \; s_{2,c} = \frac{1}{N} \sum_{i=1}^{N} v_{i,c}^{(1)}{}_{up}, \qquad (13)$$

where $s_{t,c}$ refers to the $c$-th channel of $s_t$, $v_{i,c}^{(1)}$ denotes the $c$-th channel of the $i$-th node $v_i^{(1)}$ in $V^{(1)}$, and $v_{i,c}^{(1)}{}_{up}$ denotes the $c$-th channel of the $i$-th node $v_i^{(1)}{}_{up}$ in $V^{(1)}{}_{up}$.

**Parameterized Characteristic Modeling.** To further explore the characteristics of the high-dimensional nodes, we introduce learnable parameters into characteristics. Specifically, we employ different MLPs to transform channel-wise descriptors $s_1$ and $s_2$ to new descriptors $z_1 \in R^{1 \times C}$ and $z_2 \in R^{1 \times C}$. To reduce computational costs, we simplify the MLPs through one dimension-reduction convolution and one dimension-increasing convolution layers:

$$z_1 = MLP_1(s_1), \; z_2 = MLP_2(s_2). \qquad (14)$$

**Gate Mechanism.** Informative descriptions of the aggregated nodes features consist of discriminate correlation representations of different levels. Linear aggregations mix up different correlation representations. To combine features in a non-linear and data-adaptive way, we construct a gate mechanism (Fig. 2) by applying a $softmax$ on the compact features $z_1$ and $z_2$, by

$$m_{1,c} = \frac{exp(z_{1,c})}{exp(z_{1,c}) + exp(z_{2,c})}, m_{2,c} = \frac{exp(z_{2,c})}{exp(z_{1,c}) + exp(z_{2,c})}, \qquad (15)$$

where $m_{1,c}$ and $m_{2,c}$ denote the masked $z_1$ and $z_2$ on the $c$-th channel, and $m_{1,c} + m_{2,c} = 1$.

**Aggregation.** Finally, we aggregate the ultimate nodes $V^{(u)} = \{V_{1:N,c}^{(u)} \in R^N, c = 1, 2, ..., C\}$ by the masked $m_1$ and $m_2$,

$$V_{1:N,c}^{(u)} = m_{1,c} \cdot V_{1:N,c}^{(1)} + m_{2,c} \cdot V_{1:N,c}^{(1)}{}_{up}, \qquad (16)$$

where $V_{1:N,c}^{(u)}$ denotes the $c$-th channel set of $V^{(u)}$.

This AFA module performs node-wise operations, which are again permutation-equivariant (A proof is given in Appendix C) and suitable for unordered point clouds.

## Experiments

We evaluate our model on multiple point cloud benchmarks, including Stanford Large-Scale 3D Indoor Space (S3DIS)

---

[1]Note that because the *self correlation* sub-module already enhances the cross-channel attention by linear aggregation, we do not place an AFA after *self correlation*, as Fig. 1 shows.
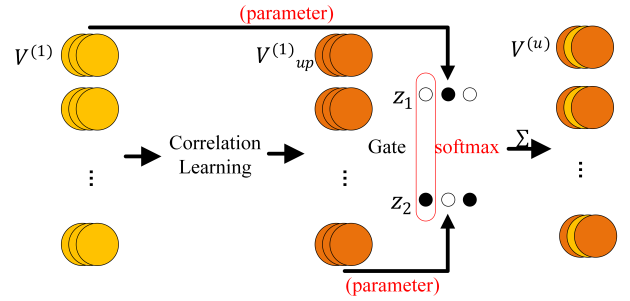


Figure 2: Illustration of Adaptive Feature Aggregation.

(Armeni et al. 2016) for semantic segmentation, ScanNet (Dai et al. 2017) for semantic voxel labeling, and ModelNet40 (Wu et al. 2015) for shape classification.

### S3DIS Semantic Segmentation

**Dataset and Metric.** The S3DIS (Armeni et al. 2016) dataset contains six sets of point cloud data from three different buildings (including 271 rooms). Each point, with xyz coordinates and RGB features, is annotated with one semantic label from 13 categories. We followed the same data preprocessing approaches given in PointCNN (Li et al. 2018).

We conducted our experiments in 6-fold cross validation (Qi et al. 2017a) and fifth fold cross validation (Tchapmi et al. 2017). Because there are no overlaps between Area 5 and the other areas, experiments on Area 5 could measure the framework's generalizability. For the evaluation metrics, we calculated mean of class-wise intersection over union (mIoU), mean of class-wise accuracy (mAcc), and overall point-wise accuracy (OA).

**Performance Comparison.** The quantitative results are given in Table 1 and Table 2. Our Point2Node performs better than the other competitive methods on both evaluation settings. In particular, despite the dataset contains objects of various sizes, our approach achieves considerable gains on both relatively small objects (e.g., chair, table, and board) and big structures (e.g., ceiling, wall, column and bookcase), which demonstrates that our model can capture discriminate features on both local and global level. For some similar geometric structures such as column vs. wall, our model also works very well, which indicates that our model is capable of distinguishing ambiguous structures due to the latent correlation they capture between the structures and other structures of the scene. Some qualitative results on different scenes are visualized in Fig. 3.

### Ablation Study

For ablation studies, we integrate each correlation learning sub-module and its corresponding feature aggregation methods, forming *self correlation* block, *local correlation* block and *non-local correlation* block, respectively. Our baseline method employs $\chi$-conv (Li et al. 2018) and builds Encoder-Decoder structure without any correlation learning block. We conducted ablation studies on Area 5 of the S3DIS.

| Methods | OA | mAcc | mIoU | ceiling | flooring | wall | beam | column | window | door | table | chair | sofa | bookcase | board | clutter |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PointNet | 78.50 | 66.20 | 47.80 | 88.00 | 88.70 | 69.30 | 42.40 | 23.10 | 47.50 | 51.60 | 54.10 | 42.00 | 9.60 | 38.20 | 29.40 | 35.20 |
| SPGraph | 85.50 | 73.00 | 62.10 | 89.90 | 95.10 | 76.40 | **62.80** | 47.10 | 55.30 | **68.40** | 69.20 | **73.50** | 45.90 | **63.20** | 8.70 | 52.90 |
| RSNet | - | 66.45 | 56.47 | 92.48 | 92.83 | 78.56 | 32.75 | 34.37 | 51.62 | 68.11 | 59.72 | 60.13 | 16.42 | 50.22 | 44.85 | 52.03 |
| DGCNN | 84.10 | - | 56.10 | - | - | - | - | - | - | - | - | - | - | - | - | - |
| PointCNN | 88.14 | 75.61 | 65.39 | **94.80** | **97.30** | 75.80 | 63.30 | 51.70 | 58.40 | 57.20 | 69.10 | 71.60 | 61.20 | 39.10 | 52.20 | 58.60 |
| PointWeb | 87.31 | 76.19 | 66.73 | 93.54 | 94.21 | 80.84 | 52.44 | 41.33 | 64.89 | 68.13 | 71.35 | 67.05 | 50.34 | 62.68 | 62.20 | 58.49 |
| Point2Node | **89.01** | **79.10** | **70.00** | 94.08 | 97.28 | **83.42** | 62.68 | **52.28** | **72.31** | 64.30 | **75.77** | 70.78 | **65.73** | 49.83 | 60.26 | **60.90** |

Table 1: Semantic segmentation results on S3DIS dataset with 6-folds cross validation.

| Methods | OA | mAcc | mIoU | ceiling | flooring | wall | beam | column | window | door | table | chair | sofa | bookcase | board | clutter |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PointNet | - | 48.98 | 41.09 | 88.80 | 97.33 | 69.80 | 0.05 | 3.92 | 46.26 | 10.76 | 58.93 | 52.61 | 5.85 | 40.28 | 26.38 | 33.22 |
| SegCloud | - | 57.35 | 48.92 | 90.06 | 96.05 | 69.86 | 0.00 | 18.37 | 38.35 | 23.12 | 70.40 | 75.89 | 40.88 | 58.42 | 12.96 | 41.60 |
| SPGraph | 86.38 | 66.50 | 58.04 | 89.35 | 96.87 | 78.12 | 0.00 | **42.81** | 48.93 | 61.58 | 84.66 | 75.41 | 69.84 | 52.60 | 2.10 | 52.22 |
| PointCNN | 85.91 | 63.86 | 57.26 | 92.31 | 98.24 | 79.41 | 0.00 | 17.60 | 22.77 | 62.09 | 74.39 | 80.59 | 31.67 | 66.67 | 62.05 | 56.74 |
| PCCN | - | 67.01 | 58.27 | 92.26 | 96.20 | 75.89 | **0.27** | 5.98 | **69.49** | **63.45** | 66.87 | 65.63 | 47.28 | 68.91 | 59.10 | 46.22 |
| PointWeb | 86.97 | 66.64 | 60.28 | 91.95 | **98.48** | 79.39 | 0.00 | 21.11 | 59.72 | 34.81 | 76.33 | **88.27** | 46.89 | 69.30 | 64.91 | 52.46 |
| GACNet | 87.79 | - | 62.85 | 92.28 | 98.27 | 81.90 | 0.00 | 20.35 | 59.07 | 40.85 | **85.80** | 78.54 | **70.75** | 61.70 | **74.66** | 52.82 |
| Point2Node | **88.81** | **70.02** | **62.96** | **93.88** | 98.26 | **83.30** | 0.00 | 35.65 | 55.31 | 58.78 | 79.51 | 84.67 | 44.07 | **71.13** | 58.72 | **55.17** |

Table 2: Semantic segmentation results on S3DIS dataset evaluated on Area 5.

To show the effectiveness of each proposed block in our model, we separately added each block after the Encoder-Decoder structure while keeping the rest unchanged in our baseline, forming Self only, Local only and Non-local only. Our proposed Point2Node employs all three blocks to form a more powerful framework. As shown in Table 3, each block plays an important role in our framework. Moreover, compared with the baseline, the performance improvement (6.26%) of Point2Node exceeds the simple summation (5.02%) of the separate block, demonstrating that our sequential structure further enhances the node features.

| Ablation studies | mIoU | Δ |
|---|---|---|
| Baseline | 56.70 | +0.00 |
| Self only | 57.45 | +0.75 |
| Local only | 57.91 | +1.23 |
| Non-local only | 59.74 | +3.04 |
| Parallel_1 | 59.51 | +2.81 |
| Parallel_2 | 60.57 | +3.87 |
| LA | 58.17 | +1.47 |
| Parameter-free AA | 60.57 | +3.87 |
| Point2Node | **62.96** | **+6.26** |

Table 3: Ablation studies on the S3DIS dataset evaluated on Area 5.

To study the effectiveness of sequentially dynamic nodes update strategy in our model, we compared our original sequential structure with parallel structure variants. Specifically, (1) by placing three blocks in parallel, denoted by "Parallel_1", the nodes get updated just once after these parallel blocks; (2) by placing the *local correlation* block and *non-local correlation* block in parallel, denoted by "Parallel_2", the nodes get updated after *self correlation* and after the parallel blocks. As shown in Table 3, "Parallel_2" performs better than "Parallel_1". The proposed Point2Node model, with sequential and fully updated nodes, has the best performance and indicates the effectiveness of this node update mechanism.

To show our adaptive aggregation can lead to more dis-

criminate features than linear aggregation, we replaced the adaptive aggregation of *local correlation* block and *non-local correlation* block with the linear aggregation in the architecture of Point2Node, to form a new network (LA). For a fair comparison, we also employed the parameter-free characteristics of the AFA module to construct a parameter-free adaptive aggregation, named "Parameter-free AA". The results demonstrate that our adaptive aggregation captures more discriminate features than linear aggregation does, and the parameterized adaptive aggregation (Point2Node) is more effective than parameter-free AA.

| | model size | mIoU |
|---|---|---|
| Baseline | **10.98M** | 56.70 |
| Point2Node | 11.14M(+1.46%) | **62.96(+11.04%)** |

Table 4: Comparisons results on model complexity.

As shown in Table 4, we significantly improve the performance (+11.03%) while increasing very few parameters (+1.46%). In our Point2Node, we frequently employ different MLPs, and we control the reduction rate $r$ ($r = 8$ in our experiments) to reduce parameters in adjacency matrix computing and parameterized characteristics.

| Mehtods | OA |
|---|---|
| 3DCNN (Bruna et al. 2014) | 73.0 |
| PointNet (Qi et al. 2017a) | 73.9 |
| TCDP (Tatarchenko et al. 2018) | 80.9 |
| PointNet++ (Qi et al. 2017b) | 84.5 |
| PointCNN (Li et al. 2018) | 85.1 |
| A-CNN (Komarichev, Zhong, and Hua 2019) | 85.4 |
| PointWeb (Zhao et al. 2019) | 85.9 |
| Point2Node | **86.3** |

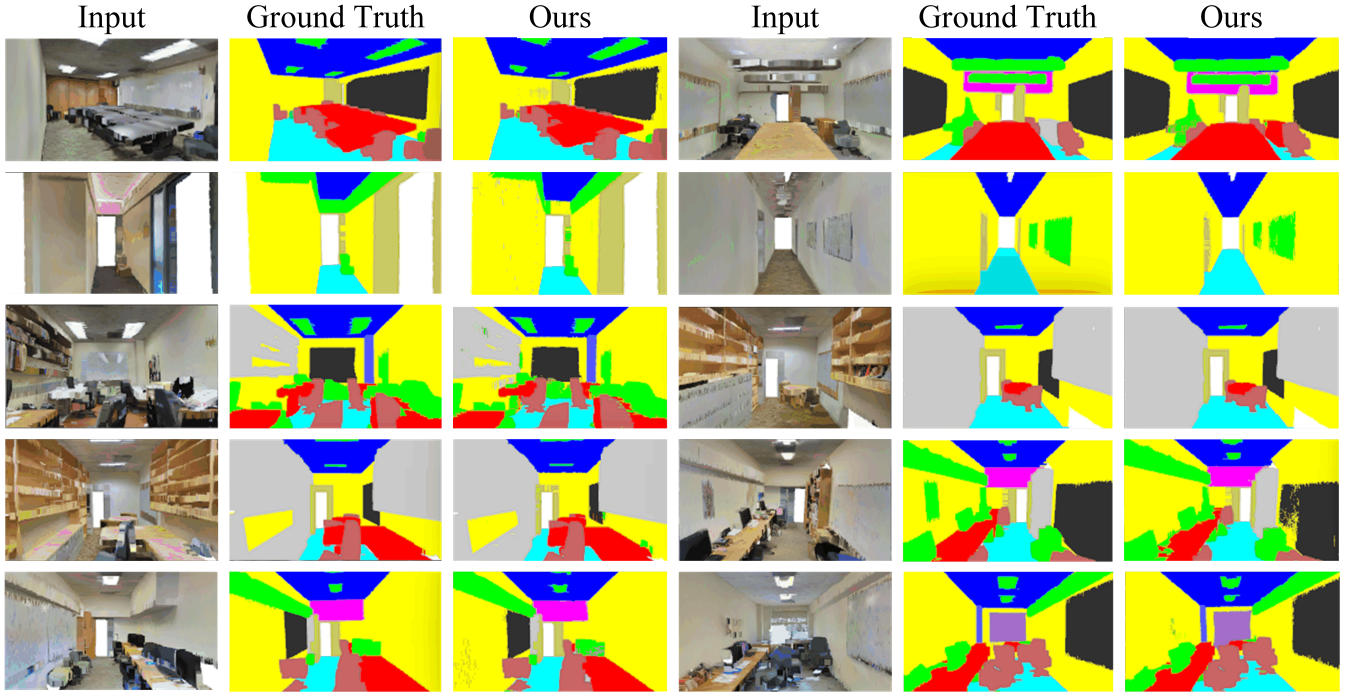Table 5: Results on ScanNet dataset.

Figure 3: Visualization of semantic segmentation results on S3DIS dataset.

## ScanNet Semantic Voxel Labeling

The ScanNet (Dai et al. 2017) dataset contains over 1,500 scanned and reconstructed indoor scenes, which consist of 21 semantic categories. The dataset provides a 1,201/312 scene split for training and testing. We followed the same data pre-processing strategies as with S3DIS, and reported the per-voxel accuracy (OA) as evaluation metrics. Table 5 shows the comparisons between our Point2Node and other competitive methods. Our method achieves the state-of-the-art performances, due to its more effective feature learning.

| Methods | input | #points | OA |
|---|---|---|---|
| PointNet (Qi et al. 2017a) | $xyz$ | $1k$ | 89.2 |
| SCN (Xie et al. 2018) | $xyz$ | $1k$ | 90.0 |
| PointNet++ (Qi et al. 2017b) | $xyz$ | $1k$ | 90.7 |
| KCNet (Shen et al. 2018) | $xyz$ | $1k$ | 91.0 |
| PointCNN (Li et al. 2018) | $xyz$ | $1k$ | 92.2 |
| DGCNN (Wang et al. 2018c) | $xyz$ | $1k$ | 92.2 |
| PCNN (Atzmon, Maron, and Lipman 2018) | $xyz$ | $1k$ | 92.3 |
| Point2Sequence (Liu et al. 2019) | $xyz$ | $1k$ | 92.6 |
| A-CNN (Komarichev, Zhong, and Hua 2019) | $xyz$ | $1k$ | 92.6 |
| Point2Node | $xyz$ | $1k$ | **93.0** |
| SO-Net (Li, Chen, and Lee 2018) | $xyz$ | $2k$ | 90.9 |
| PointNet++ (Qi et al. 2017b) | $xyz, normal$ | $5k$ | 91.9 |
| PointWeb (Zhao et al. 2019) | $xyz, normal$ | $1k$ | 92.3 |
| PointConv (Wu, Qi, and Li 2019) | $xyz, normal$ | $1k$ | 92.5 |
| SpiderCNN (Xu et al. 2018) | $xyz, normal$ | $5k$ | 92.4 |
| SO-Net (Li, Chen, and Lee 2018) | $xyz, normal$ | $5k$ | 93.4 |

Table 6: Shape classification results on ModelNet40 dataset.

## ModelNet40 Shape Classification

The ModelNet40 (Wu et al. 2015) dataset is created using 12,311 3D mesh models from 40 man-made object categories, with an official 9,843/2,468 split for training and

testing. Following PointCNN (Li et al. 2018), we uniformly sampled 1,024 points from each mesh without normal information to train our model. For classification, we removed the decoder in the Point2Node model and reported the overall accuracy (OA) as the evaluation metric.

To make a fair comparison with previous approaches that may take different information as input, we list detailed experiment setting and results in Table 6. Point2Node achieves the highest accuracy among the $1k$-$xyz$-input methods. Furthermore, our $1k$-$xyz$-input model outperforms most additional-input methods, even if performing slightly worse than the best additional-input method SO-Net (Li, Chen, and Lee 2018). Experimental results also demonstrate the effectiveness of our model in point cloud understanding.

## Conclusions

We present a novel framework, *Point2Node*, for correlation learning among 3D points. Our Point2Node dynamically enhances the nodes representation capability by exploiting correlation of different levels and adaptively aggregates more discriminate features to better understand point cloud. Ablation studies show that our performance gain comes from our new correlation learning modules. Results from a series of experiments demonstrate the effectiveness and generalization ability of our method.

*Limitations.* The correlation learning in Point2Node is based on refining the features (currently PointCNN (Li et al. 2018) is used) extracted on nodes. While it can be considered as a general feature enhancement tool, the final performance is also limited by the capability of the pre-determined features. An ideal strategy would be to integrate the feature

learning and correlation learning into a joint pipeline, so that the inter-channel and inter-node correlations could impact the feature extraction. We will explore the design of such an end-to-end network in the future.

## Acknowledgments

## References

[Armeni et al. 2016] Armeni, I.; Sener, O.; Zamir, A. R.; Jiang, H.; Brilakis, I.; Fischer, M.; and Savarese, S. 2016. 3D semantic parsing of large-Scale indoor spaces. In *CVPR*, 1534–1543.

[Atwood and Towsley 2016] Atwood, J., and Towsley, D. 2016. Diffusion-convolutional neural networks. In *NIPS*. 1993–2001.

[Atzmon, Maron, and Lipman 2018] Atzmon, M.; Maron, H.; and Lipman, Y. 2018. Point convolutional neural networks by extension operators. *ACM* 1–12.

[Bruna et al. 2014] Bruna, J.; Zaremba, W.; Szlam, A.; and YannLecun. 2014. Spectral networks and locally connected networks on graphs. In *ICLR*.

[Chollet 2017] Chollet, F. 2017. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, 1800–1807.

[Dai et al. 2017] Dai, A.; Chang, A. X.; Savva, M.; Halber, M.; Funkhouser, T.; and Niessner, M. 2017. ScanNet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 5828–5839.

[Jiancheng Yang 2019] Jiancheng Yang, Qiang Zhang, B. N. L. L. J. L. M. Z. Q. T. 2019. Modeling point clouds with self-attention and gumbel subset sampling. In *CVPR*.

[Kalogerakis et al. 2017] Kalogerakis, E.; Averkiou, M.; Maji, S.; and Chaudhuri, S. 2017. 3D shape segmentation with projective convolutional networks. In *CVPR*, 6630–6639.

[Kipf and Welling 2016] Kipf, T. N., and Welling, M. 2016. Semi-supervised classification with graph convolutional networks. *arXiv:1609.02907*.

[Klokov and Lempitsky 2017] Klokov, R., and Lempitsky, V. 2017. Escape from Cells: Deep kd-networks for the recognition of 3D point cloud models. In *ICCV*, 863–872.

[Komarichev, Zhong, and Hua 2019] Komarichev, A.; Zhong, Z.; and Hua, J. 2019. A-CNN: Annularly convolutional neural networks on point clouds. In *CVPR*.

[Li et al. 2018] Li, Y.; Rui, B.; Sun, M.; and Chen, B. 2018. PointCNN: Convolution on x-transformed points. *arXiv:1801.07791*.

[Li, Chen, and Lee 2018] Li, J.; Chen, B. M.; and Lee, G. H. 2018. SO-Net: Self-organizing network for point cloud analysis. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 93979406.

[Liu et al. 2019] Liu, X.; Han, Z.; Liu, Y. S.; and Zwicker, M. 2019. Point2Sequence: Learning the shape representation of 3D point clouds with an attention-based sequence to sequence network. *arXiv:1811.02565*.

[Maturana and Scherer 2015] Maturana, D., and Scherer, S. 2015. VoxNet: A 3D convolutional neural network for real-time object recognition. In *IROS*, 922–928.

[Qi et al. 2016] Qi, C. R.; Su, H.; NieBner, M.; Dai, A.; Yan, M.; and Guibas, L. J. 2016. Volumetric and multi-view CNNs for object classification on 3d data. In *CVPR*, 5648–5656.

[Qi et al. 2017a] Qi, C. R.; Su, H.; Mo, K.; and Guibas, L. J. 2017a. PointNet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, 652–660.

[Qi et al. 2017b] Qi, C. R.; Yi, L.; Su, H.; and Guibas, L. J. 2017b. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NIPS*. 5099–5108.

[Riegler, Ulusoy, and Geiger 2017] Riegler, G.; Ulusoy, A. O.; and Geiger, A. 2017. OctNet: Learning deep 3D representations at high resolutions. In *CVPR*, 6620–6629.

[Shen et al. 2018] Shen, Y.; Chen, F.; Yang, Y.; and Dong, T. 2018. Mining point cloud local structures by kernel correlation and graph pooling. In *CVPR*, 4548–4557.

[Simonovsky and Komodakis 2017] Simonovsky, M., and Komodakis, N. 2017. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In *CVPR*, 3693–3702.

[Su et al. 2015] Su, H.; Maji, S.; Kalogerakis, E.; and Learned-Miller, E. 2015. Multi-view convolutional neural networks for 3d shape recognition. In *ICCV*, 945–953.

[Tatarchenko et al. 2018] Tatarchenko, M.; Park, J.; Koltun, V.; and Zhou, Q. 2018. Tangent convolutions for dense prediction in 3D. In *CVPR*, 3887–3896.

[Tchapmi et al. 2017] Tchapmi, L. P.; Choy, C. B.; Armeni, I.; Gwak, J. Y.; and Savarese, S. 2017. Segcloud: Semantic segmentation of 3D point clouds. In *3DV*.

[Verma, Boyer, and Verbeek 2017] Verma, N.; Boyer, E.; and Verbeek, J. 2017. Dynamic filters in graph convolutional networks. *arXiv:1706.05206*.

[Wang et al. 2018a] Wang, S.; Suo, S.; Ma, W.-C.; Pokrovsky, A.; and Urtasun, R. 2018a. Deep parametric continuous convolutional neural networks. In *CVPR*, 2589–2597.

[Wang et al. 2018b] Wang, X.; Girshick, R.; Gupta, A.; and He, K. 2018b. Non-local neural networks. In *CVPR*, 7794–7803.

[Wang et al. 2018c] Wang, Y.; Sun, Y.; Liu, Z.; Sarma, S. E.; Bronstein, M. M.; and Solomon, J. M. 2018c. Dynamic graph CNN for learning on point clouds. *arXiv:1801.07829*.

[Wu et al. 2015] Wu, Z.; Song, S.; Khosla, A.; Yu, F.; Zhang, L.; Tang, X.; and Xiao, J. 2015. 3D ShapeNets: A deep representation for volumetric shapes. In *CVPR*, 1912–1920.

[Wu, Qi, and Li 2019] Wu, W.; Qi, Z.; and Li, F. 2019. PointConv: Deep convolutional networks on 3D point clouds. In *CVPR*.

[Xie et al. 2018] Xie, S.; Liu, S.; Chen, Z.; and Tu, Z. 2018. Attentional ShapeContextNet for point cloud recognition. In *CVPR*, 4606–4615.

[Xu et al. 2018] Xu, Y.; Fan, T.; Xu, M.; Long, Z.; and Yu, Q. 2018. Spidercnn: Deep learning on point sets with parameterized convolutional filters. In *ECCV*, 90–105.

[Yi et al. 2017] Yi, L.; Su, H.; Guo, X.; and Guibas, L. 2017. SyncSpecCNN: Synchronized spectral CNN for 3D shape segmentation. In *CVPR*, 2282–2290.

[Zhao et al. 2019] Zhao, H.; Jiang, L.; Fu, C.-W.; and Jia, J. 2019. PointWeb: Enhancing local neighborhood deatures for point cloud processing. In *CVPR*.

[zgn Çiçek et al. 2016] zgn Çiçek; Abdulkadir, A.; Lienkamp, S. S.; Brox, T.; and Ronneberger, O. 2016. 3D U-Net: Learning dense volumetric segmentation from sparse annotation. In *MICCAI*. 424–432.

# Appendix

## A. Relationship to standard Non-local operations

Standard Non-Local operations (Wang et al. 2018b) are widely employed in recent literature to capture long-range dependencies in 2D images and videos. To capture correlation among pixels in videos, these operations map input features with 3D convolution, and construct adjacency matrix with reshape operation and matrix multiplication.

In our papers, to capture the correlation among nodes, we generalize standard Non-Local operations to the graph domain, where the adjacency matrix is more naturally applied for weighting edges. We not only employ Non-local operations on non-local nodes to capture long-range dependencies, but also extend them to neighboring nodes to encode local contextual information.

Moreover, our computation is less expensive. In (Wang et al. 2018b), the dimension for video data fed to the network is $T \times H \times W \times C$, where $T$ is the time of videos, $H$ is the height of videos, $W$ is the width of videos, and $C$ is the channel of videos. The dimension of our node tensor is $N \times C$, where $N$ is the number of nodes and $C$ is the channel of nodes. Because the 3D convolutions are replaced with 1D convolutions, and the frequent "reshape" operations are no needed, our computation is much more efficient.

## B. Proof of Permutation Equivariance of Dynamic Nodes Correlation

Following (Jiancheng Yang 2019), we prove the permutation equivariance of our proposed modules, and provide the proof process as follows.

**Lemma 1.** (Permutation matrix and permutation function). $\forall A \in R^{N \times N}$, $\forall$ permutation matrix $P$ of size $N$, $\exists p : \{1, 2, ..., N\} \to \{1, 2, ..., N\}$ is a permutation function:

$$A_{ij} = (P \cdot A)_{p(i)j} = (A \cdot P^T)_{ip(j)} = (P \cdot A \cdot P^T)_{p(i)p(j)} \tag{17}$$

**Lemma 2.** Given $A \in R^{N \times N}$, $\forall$ permutation matrix $P$ of size $N$,

$$softmax(P \cdot A \cdot P^T) = P \cdot softmax(A) \cdot P^T \tag{18}$$

*Proof.*
$softmax(A)_{ij} = \frac{e^{A_{ij}}}{\sum_{n=1}^N e^{A_{ij}}}$. *Consider a permutation function p for P, using Eq. 17, we get:*

$$(P \cdot softmax(A) \cdot P^T)_{p(i)p(j)} = softmax(A)_{ij}$$
$$= \frac{e^{A_{ij}}}{\sum_{n=1}^N e^{A_{ij}}}$$
$$= \frac{e^{(PAP^T)_{p(i)p(j)}}}{\sum_{n=1}^N e^{(PAP^T)_{p(i)p(j)}}}$$
$$= softmax(P \cdot A \cdot P^T)_{p(i)p(j)}$$

*which implies,*

$$P \cdot softmax(A) \cdot P^T = softmax(P \cdot A \cdot P^T)$$

*similarly,*

$$P \cdot softmax(A) = softmax(P \cdot A) \tag{19}$$

$$softmax(A) \cdot P^T = softmax(A \cdot P^T) \tag{20}$$

**Proposition 1.** *Self Correlation (SC) operation is permutation-equivariant, i.e., given input $X \in R^{N \times C}$, $\forall$ permutation matrix $P$ of size $N$,*

$$SC(P \cdot X) = P \cdot SC(X)$$

*Proof.*

$$SC(X) = X + \alpha(softmax(MLP(X)) \odot X)$$

*where MLP is a node-wise function that shares weights among all the nodes,*

$$MLP(P \cdot X) = P \cdot MLP(X)$$

*using Eq. 19, we get:*

$$SC(P \cdot X) = P \cdot X + \alpha(softmax(MLP(P \cdot X)) \odot (P \cdot X))$$
$$= P \cdot X + \alpha(softmax(P \cdot MLP(X)) \odot (P \cdot X))$$
$$= P \cdot X + \alpha(P \cdot softmax(MLP(X)) \odot (P \cdot X))$$
$$= P \cdot X + P \cdot \alpha(softmax(MLP(X)) \odot X)$$
$$= P \cdot (X + \alpha(softmax(MLP(X)) \odot X))$$
$$= P \cdot SC(X)$$

**Proposition 2.** *Local Correlation (LC) operation is permutation-equivariant, i.e., given input $X \in R^{N \times C}$, $\forall$ permutation matrix $P$ of size $N$,*

$$LC(P \cdot X) = P \cdot LC(X)$$

*Proof.*

$$LC(X) = LNMP(softmax(\theta^{(l)}(DKNN(X)) \cdot (\varphi^{(l)}(DKNN(X)))^T) \cdot DKNN(X))$$

*where Dilated K Nearest Neighbors (DKNN) operation is a node-wise operation,*

$$DKNN(P \cdot X) = P \cdot DKNN(X)$$

$\theta^{(l)}$ *and* $\varphi^{(l)}$ *are both node-wise functions that share weights among all the nodes,*

$$\theta^{(l)}(P \cdot X) = P \cdot \theta^{(l)}(X)$$
$$\varphi^{(l)}(P \cdot X) = P \cdot \varphi^{(l)}(X)$$

*Local Nodes Max-Pooling (LNMP) operation is invariant to local nodes' order, and permutation-equivariant to input order,*

$$LNMP(P \cdot X) = P \cdot LNMP(X)$$

*using Eq. 18, we get:*

$$LC(P \cdot X) = LNMP(softmax(\theta^{(l)}(DKNN(P \cdot X)) \cdot$$
$$(\varphi^{(l)}(DKNN(P \cdot X)))^T) \cdot DKNN(P \cdot X))$$
$$= LNMP(softmax(\theta^{(l)}(P \cdot DKNN(X)) \cdot$$
$$(\varphi^{(l)}(P \cdot DKNN(X)))^T) \cdot P \cdot DKNN(X))$$
$$= LNMP(softmax(P \cdot \theta^{(l)}(DKNN(X)) \cdot$$
$$(\varphi^{(l)}(DKNN(X)))^T \cdot P^T) \cdot P \cdot DKNN(X))$$
$$= LNMP(P \cdot softmax(\theta^{(l)}(DKNN(X)) \cdot$$
$$(\varphi^{(l)}(DKNN(X)))^T) \cdot P^T \cdot P \cdot DKNN(X))$$
$$= LNMP(P \cdot softmax(\theta^{(l)}(DKNN(X)) \cdot$$
$$(\varphi^{(l)}(DKNN(X)))^T) \cdot DKNN(X))$$
$$= P \cdot LNMP(softmax(\theta^{(l)}(DKNN(X)) \cdot$$
$$(\varphi^{(l)}(DKNN(X)))^T) \cdot DKNN(X))$$
$$= P \cdot LC(X)$$

**Proposition 3.** *Non-local Correlation (NLC) operation is permutation-equivariant, i.e., given input* $X \in R^{N \times C}$, $\forall$ *permutation matrix P of size N,*

$$NLC(P \cdot X) = P \cdot NLC(X)$$

*Proof.*

$$NLC(X) = softmax(\theta^{(g)}(X) \cdot (\varphi^{(g)}(X))^T) \cdot X$$

*where* $\theta^{(g)}$ *and* $\varphi^{(g)}$ *are both node-wise functions that share weights among all the nodes,*

$$\theta^{(g)}(P \cdot X) = P \cdot \theta^{(g)}(X)$$
$$\varphi^{(g)}(P \cdot X) = P \cdot \varphi^{(g)}(X)$$

*using Eq. 18, we get:*

$$NLC(P \cdot X) = softmax(\theta^{(g)}(P \cdot X) \cdot (\varphi^{(g)}(P \cdot X))^T) \cdot$$
$$(P \cdot X)$$
$$= softmax(P \cdot \theta^{(g)}(X) \cdot (\varphi^{(g)}(X))^T \cdot P^T) \cdot$$
$$(P \cdot X)$$
$$= P \cdot softmax(\theta^{(g)}(X) \cdot (\varphi^{(g)}(X))^T) \cdot P^T \cdot$$
$$(P \cdot X)$$
$$= P \cdot softmax(\theta^{(g)}(X) \cdot (\varphi^{(g)}(X))^T) \cdot X$$
$$= P \cdot NLC(X)$$

## C. Proof of Permutation Equivariance of Adaptive Feature Aggregation

**Proposition 1.** *Adaptive Feature Aggregation (AFA) operation is permutation-equivariant, i.e., given input* $X \in R^{N \times C}$ *and* $Y \in R^{N \times C}$, $\forall$ *permutation matrix P of size N,*

$$AFA(P \cdot X, P \cdot Y) = P \cdot AFA(X, Y)$$

*Proof.*

$$AFA(X, Y) = Gate(GNP(X, Y))_X \odot X$$
$$+ Gate(GNP(X, Y))_Y \odot Y$$

*where Global Node Pooling (GNP) operation is invariant to input order,*

$$GNP(P \cdot X, P \cdot Y) = GNP(X, Y)$$

*Gate operation is shared among all the points,*

$$Gate(A)_{P \cdot X} = Gate(A)_X$$

*in this way,*

$$AFA(P \cdot X, P \cdot Y)$$
$$= Gate(GNP(P \cdot X, P \cdot Y))_{P \cdot X} \odot (P \cdot X)$$
$$+ Gate(GNP(P \cdot X, P \cdot Y))_{P \cdot Y} \odot (P \cdot Y)$$
$$= Gate(GNP(X, Y))_{P \cdot X} \odot (P \cdot X)$$
$$+ Gate(GNP(X, Y))_{P \cdot Y} \odot (P \cdot Y)$$
$$= Gate(GNP(X, Y))_X \odot (P \cdot X)$$
$$+ Gate(GNP(X, Y))_Y \odot (P \cdot Y)$$
$$= P \cdot (Gate(GNP(X, Y))_X \odot X$$
$$+ Gate(GNP(X, Y))_Y \odot Y)$$
$$= P \cdot AFA(X, Y)$$