
GSIP: GREEN SEMANTIC SEGMENTATION OF LARGE-SCALE INDOOR POINT CLOUDS

A PREPRINT

Min Zhang
Media Communications Lab
University of Southern California
Los Angeles, CA, USA
zhan980@usc.edu

Pranav Kadam
Media Communications Lab
University of Southern California
Los Angeles, CA, USA
pranavka@usc.edu

Shan Liu*
Tencent Media Lab
Tencent America
Palo Alto, CA, USA
shanl@tencent.com

C.-C. Jay Kuo
Media Communications Lab
University of Southern California
Los Angeles, CA, USA
cckuo@sipi.usc.edu

November 15, 2022

ABSTRACT

An efficient solution to semantic segmentation of large-scale indoor scene point clouds is proposed in this work. It is named GSIP (Green Segmentation of Indoor Point clouds) and its performance is evaluated on a representative large-scale benchmark — the Stanford 3D Indoor Segmentation (S3DIS) dataset. GSIP has two novel components: 1) a room-style data pre-processing method that selects a proper subset of points for further processing, and 2) a new feature extractor which is extended from PointHop. For the former, sampled points of each room form an input unit. For the latter, the weaknesses of PointHop’s feature extraction when extending it to large-scale point clouds are identified and fixed with a simpler processing pipeline. As compared with PointNet, which is a pioneering deep-learning-based solution, GSIP is green since it has significantly lower computational complexity and a much smaller model size. Furthermore, experiments show that GSIP outperforms PointNet in segmentation performance for the S3DIS dataset.

Keywords Point clouds · unsupervised learning · large scale · semantic segmentation

1 Introduction

Given a point cloud set, the goal of semantic segmentation is to label every point as one of the semantic categories. Semantic segmentation of large-scale point clouds find a wide range of real world applications such as autonomous driving in an out-door environment and robotic navigation in an in- or out-door environment. As compared with the point cloud classification problem that often targets at small-scale objects, a high-performance point cloud semantic segmentation method demands a good understanding of the complex global structure as well as the local neighborhood of each point. Meanwhile, efficiency measured by computational complexity and memory complexity is important for practical real-time systems.

State-of-the-art point cloud classification and segmentation methods are based on deep learning. Raw point clouds captured by the LiDAR sensors are irregular and unordered. They cannot be directly processed by deep learning networks designed for 2D images. This problem was addressed by the pioneering work on the PointNet [1]. PointNet and its follow-ups [2–5] achieved impressive performance in small-scale point cloud classification and segmentation

*This work was supported by Tencent Media Lab.

tasks. A representative point cloud classification dataset is ModelNet40 [6], where each object has 1,024 or 2,048 points.

Although the above-mentioned methods work well for small-scale point clouds, they can't be generalized to handle large-scale point cloud directly due to the memory and time constraints. Large-scale point cloud semantic segmentation methods are often evaluated on the S3DIS dataset [7] for the indoor scene and the Semantic3D [8] or SemanticKITTI [9] datasets for the outdoor scenes. These datasets have millions of points, covering up to 200×200 meters in 3D real-world space. In order to feed such a large amount of data to deep learning networks, it is essential to conduct pre-processing steps on the data such as block partitioning [1]. Recently, efforts have been made in [10–13] to tackle with large-scale point clouds directly. For example, RandLA-Net [13] abandons block partitioning and enlarges an input unit from 4,096 points to 40,960 points.

Green point cloud learning methods are proposed recently, e.g., [14–18]. They are called “green” since they have smaller model sizes and can be trained by CPU with much less time and complexity. The main saving comes from the one-pass feedforward feature learning mechanism, where no backpropagation neither supervision label is needed. The unsupervised feature learning is built upon statistical analysis of points in a local neighborhood of a point cloud set. Along this direction, we attempt to develop a green semantic segmentation solution for large-scale indoor point clouds in this work, named GSIP (Green Segmentation of Indoor Point clouds).

GSIP has two novel components: 1) a room-style data pre-processing method and 2) an enhanced PointHop feature extractor. For the former, we compare existing data pre-process techniques and identify the most suitable data preparation method. In deep-learning methods, which are implemented on GPUs, data pre-processing is used to ensure that input units contain same number of points so that the throughput is maximized by exploiting GPU's built-in parallel processing capability. For the S3DIS dataset, a unit has 4,096 points in PointNet which adopts the block-style pre-processing while a unit has 40,960 points in RandLA-Net which adopts the view-style pre-processing. Since feature extraction for green point cloud learning runs on a CPU, we are able to relax the constraint that each input unit can have different number of points. We propose a new room-style pre-processing method for GSIP and show its advantages. For the latter, we point out some weaknesses of PointHop's feature extraction when extending it to large-scale point clouds and fix them with a simpler processing pipeline.

The proposed GSIP method can be stated below. A raw point cloud set is voxel-downsampled with a fixed grid size. All down-sampled points in one room form one unit whose point number ranges from 10K to 200K. Then, feature of points in each unit are fed into an unsupervised feature extractor to obtain point-wise features at various hops. Finally, features are fed to the XGBoost classifier [19] for point-wise classification, i.e., semantic segmentation. We conduct experiments on the S3DIS dataset [6] and show that, besides a smaller model size and lower computational complexity, GSIP outperforms PointNet in segmentation performance. Our work has the following three main contributions:

1. We compare existing data pre-processing methods for point cloud semantic segmentation and identify the most suitable data preparation method.
2. We propose a lightweight and efficient unsupervised feature extraction method in GSIP.
3. We demonstrate the performance, efficiency and model size gains of GSIP over PointNet on the S3DIS dataset.

2 Related Work

2.1 Traditional Methods

Traditional features for 3D point clouds are extracted using a hand-crafted solution. It does not involve any training data but relies on local geometric properties of points. FPFH [20] and SHOT [21] are exemplary feature descriptors. Semantic segmentation [22, 23] is formulated as a point-wise classification problem. Each point is described by one or multiple feature descriptors, and its extracted features are concatenated to form a feature vector. Then, the feature vector is fed into a classifier such as the Support Vector Machine (SVM) [24] and the Random Forest (RF) [25]. The classification stage does need representative training data for classifier's training. In the inference stage, the classifier will assign a class label to each point of the target point cloud set.

2.2 Deep Learning Methods

One pioneering work in deep learning methods is PointNet [1]. It adopts multi-layer perceptrons (MLPs) and max pooling to learn point-wise features. Yet, PointNet fails to capture the local structure of a point. Its follow-ups [2–5] attempt to capture the local information of a point to learn a richer context. For example, PointNet++ [2] applies PointNet to the local region of each point, and then aggregates local features in a hierarchical architecture. DGCNN [5]

Table 1: Comparison of traditional, deep learning (DL) and green learning (GL) methods.

Feature	Traditional	DL	GL
Data eagerness	low	high	middle
Supervision	low	high	middle
Model size	small	large	middle
Time complexity	low	high	middle
Interpretation	easy	hard	easy
Performance	poor	good	good

exploits another idea to learn better local features. It uses point features to build a dynamical graph and updates neighbor regions at every layer. A dynamic graph can capture better semantic meaning than a fixed graph. PointCNN [4] learns an χ -transformation from an input point cloud to get weights of neighbors of each point and permute points based on a latent and potentially canonical order. To make local features invariant to permutations, PointSIFT [3] discards the max pooling idea but designs an orientation-encoding unit to encode the eight orientations, which is inspired by the well-known 2D SIFT descriptor [26]. RandLA-Net [13] learns important local features through the attention mechanism. A shared MLP followed by the softmax operation is used to compute attention scores. Then, attention-weighted local features are summed together.

2.3 Green Learning Methods

Green point cloud learning has been introduced in a sequence of recent publications [14–18]. Its theoretical foundation is successive subspace learning (SSL) [27–33], which applies to 2D images as well as 3D point clouds. SSL helps reduce the model size and computation complexity, and offers mathematical transparency. As compared with traditional and deep learning methods. We compare traditional, deep learning and green learning methods qualitatively in several aspects in Table 1. As shown in the table, green point cloud learning methods have several advantages:

1. data-driven but not data-eager (i.e. robust with less data),
2. no supervision needed for feature extraction,
3. good generalization ability, where extracted feature can be used for multi-tasking such as object classification, part segmentation and registration,
4. mathematically interpretable,
5. smaller model sizes and lower computation complexity.

Apart from point cloud processing, the green learning technology has been successfully applied to other tasks such as face gender classification [34], deepfake image detection [35], and image anomaly detection [36].

2.4 Large-scale Point Cloud Learning

Point-based deep learning methods as reviewed in Sec. 2.2 are limited to small-scale point clouds. For large-scale point cloud learning, SPG [10] builds super graphs composed by super points in a pre-processing step. Then, it learns semantics for each super point rather than a point. FCPN [11] and PCT [12] combine voxelization and point-based networks to process large-scale point clouds. However, graph construction and voxelization are computationally heavy so that these solutions are not suitable for mobile or edge devices. Recently, RandLA-Net [13] revisits point-based deep learning methods. It replaces the farthest point sampling (FPS) method with random sampling (RS) to save time and memory cost. In this way, the number of points that can be processed one time is increased by 10 times. However, RS may discard key features, so it is not as accurate as FPS. To address this problem, a new local feature aggregation module is adopted to capture complex local structures.

3 Proposed GSIP Method

An overview of the proposed GSIP method is given in Fig. 1. A raw point cloud set is first voxel-downsampled, where each downsampled point has 9D attributes (i.e., XYZ, RGB, and normalized XYZ). Each sampled point is augmented with 12D additional attributes. They include: surface normals (3D), geometric features (6D, planarity, linearity, surface

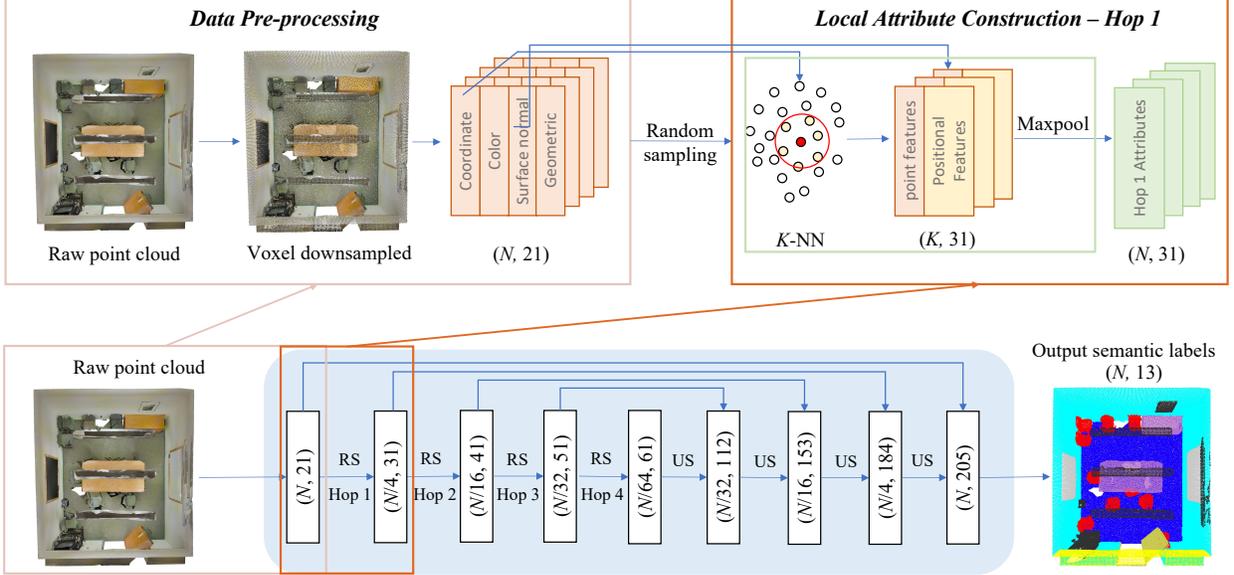


Figure 1: An overview of the proposed GSIP method, where the upper left block shows the data pre-processing, the upper right block shows the local attribute construction process and the lower block shows the encoder-decoder architecture for large-scale point cloud semantic segmentation.

variation, etc.) [22] and HSV color space (3D, converted from RGB). Surface normals and geometric features are commonly used in traditional point cloud processing. As compared with the XYZ coordinates, they describe local geometric patterns and can be easily computed. Here, we include them as additional features to the XYZ coordinates. For points that have the same geometric pattern but belong to different objects (e.g., wall and blackboard), the RGB plus HSV color spaces work better than RGB alone. Consequently, each point has 21D input features. The point-wise features are not powerful enough since they are localized in the space. To obtain more powerful features, we need to group points based on different neighborhood sizes (e.g., points in a small neighborhood, a mid-size neighborhood and a large neighborhood). By borrowing the terminology from graph theory, we call these neighborhoods hops, where hop 1 denotes the smallest neighborhood size. To carry out the semantic segmentation task, we need to extract features at various hops, which is achieved by an unsupervised feature extractor. The feature extractor adopts the encoder-decoder architecture. It has four encoding hops followed by four decoding hops. Finally, a classifier is trained and used to classify each point into a semantic category based on its associated hop features.

3.1 Data Pre-processing

Data pre-processing is used to prepare the input data in proper form so that they can be fed into the learning pipeline for effective processing. Although the generic principle holds for any large-scale point clouds, the implementation detail depends on the application context. Data pre-processing techniques targeting at large-scale indoor scene point clouds are discussed here. Specifically, we use the S3DIS dataset [7], which is a subset of the Stanford 2D-3D-Semantics dataset [37], as an example. S3DIS is one of the benchmark datasets for point cloud semantic segmentation. It contains point clouds scanned from 6 indoor areas with 271 rooms. There are 13 object categories in total, including ceiling, floor, wall, door, etc. Each point has 9 dimensions, i.e., XYZ, RGB and normalized XYZ. Data pre-processing techniques can be categorized into the following three styles.

Block Style. Block partitioning was proposed by PointNet and adopted by its follow-ups. The 2D horizontal plane of a room is split into 1×1 meter blocks while its 1D vertical direction is kept as it is to form an input unit as shown in Fig. 2(a). Each unit contains 4,096 points, which is randomly sampled from its raw block. In the inference stage, 4,096 points of each unit are classified and, then, their predicted labels are propagated to their neighbors. Typically, the k-fold strategy is adopted for train and test. For example, if area 6 is selected as the test area, the remaining 5 areas are used for training. Under this setting, the block-style data pre-processing has 20,291 training units and 3,294 testing units, where each unit has 4,096 points.

View Style. View-style data pre-processing was adopted by RandLA-Net, but not explained in the paper. We get its details from the codes and describe its process below. It first partitions the 3D space of a room into voxel grids and

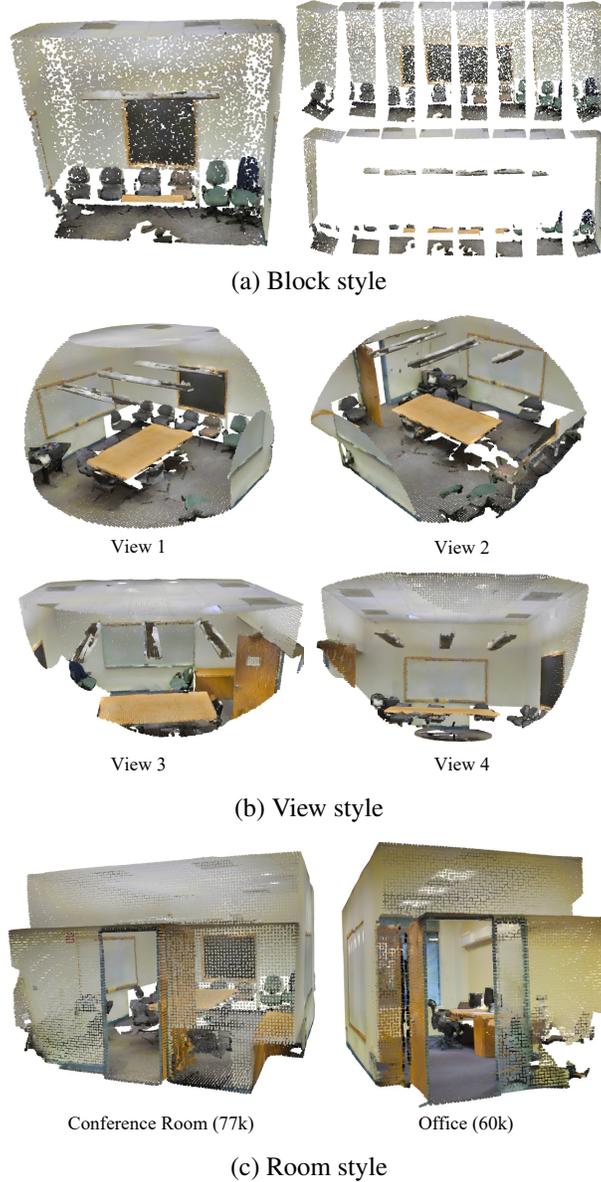


Figure 2: Comparison of three data pre-processing methods.

randomly selects one point per voxel. For instance, the first conference room of area 1 has about one million points. One can obtain 70k points by setting the voxel grid size to 0.04, around 7.7% points are kept in this example. This is a commonly used point downsampling procedure. Then, it iteratively selects a fixed number of points generating input units to facilitate GPU parallel processing. For initialization, it randomly selects a room and a point in the room as the reference point. Then, it finds the 40,960 nearest points around the reference point by K -NN algorithm to form the first unit. This process is repeated to get a sequence of input units until it reaches the target unit number. For the S3DIS dataset, the target training and test unit numbers for each fold are set to 3K and 2K, respectively. To reduce the overlapping of different units, it assigns a possibility to every point randomly at the beginning and updates the possibilities of the selected points in each round as inversely proportional to their distances to the reference point. Thus, unselected points will be more likely to be chosen as the next reference point. Four examples are shown in Fig. 2(b) to visualize points inside the same unit. We see that they offer certain views to a room.

There are however several problems for the above two data pre-processing methods. First, they do not have a global view of the whole room, resulting in inaccurate nearest neighbor search in unit boundaries. Second, the view-style method

generates 3K units for 200 rooms (i.e., 15 units per room on the average) in the training. There are redundant points in the intersection of views of the same room. The total number of training points increases from 80M ($\approx 20,291 \times 4,096$) of the block style to 120M ($\approx 3,000 \times 40,960$) of the view style. We may ask whether it is essential to have so many training points.

Room Style. It is desired to increase input scale and reduce the total number of training points while keeping the same segmentation performance. To achieve this goal, we propose a room-style pre-processing method and adopt a flexible feature extraction pipeline, which will be discussed in Sec. 3.2. The most distinctive aspect of the room style is that each unit can have a different number of points. Thus, we include all points in a room in one unit and call it room-style pre-processing. This is a possible solution since our point cloud feature extractor can be implemented in the CPU effectively. By relaxing the constraint that all units have the same number of points demanded by the GPU implementation, the data pre-processing problem can be greatly simplified. By following the first step of RandLA-Net’s pre-processing, we downsample raw points by voxel downsampling method with a fixed grid size. Afterwards, the number of points for each room ranges from 10K to 200K. By taking rooms in areas 1-5 as training and rooms in area 6 as testing as an example, we have 224 rooms (or 224 units) for training and 48 rooms (or 48 units) for testing. The training and testing sets contain 15M and 2M points, respectively. The total number of training points is much smaller than those of the block-style and the view-style methods while the input scale is much larger.

3.2 Feature Extractor and Classifier

Points in a room-style input unit are fed into a feature extractor to obtain point-wise features as shown in Fig. 1. The upper right block shows the local attribute construction process and the lower block shows the encoder-decoder architecture for large-scale point cloud semantic segmentation. It is developed upon our previous work PointHop. The reason for developing new feature learner is due to several shortcomings when extending PointHop from small-scale to large-scale point cloud learning. It is detailed below.

PointHop and PointHop++ Feature Extractor. PointHop [14] and PointHop++ [15] are unsupervised feature extractors proposed for small-scale point cloud classification. They have been successfully applied to joint point cloud classification and part segmentation [16] and point cloud registration [17, 18]. PointHop constructs attributes from local to global by stacking several hops, covering small-, mid- and long-range neighborhoods. In each hop, the local neighborhood of each point are divided by eight octants using the the 3D coordinates of local points. Then, point attributes in each octant are aggregated and concatenated to form a local descriptor, which keeps more information than naive max pooling. Due to the fast dimension growing, the Saab transform [29], which is a variant of Principal Component Analysis (PCA) [38], is used for dimension reduction. Between two consecutive units, FPS is used to downsample the point cloud to increase the speed as well receptive field of each point. PointHop++ [15] is an extension of PointHop. It has a lower model size and training complexity by leveraging the observation that spectral features are uncorrelated after PCA. Thus, we can conduct PCA to each spectral feature independently, which is called the channel-wise Saab transform.

There are three shortcomings for the pipeline used in PointHop and PointHop++ when it is ported to large-scale point cloud data. First, the computational efficiency is limited by FPS between two consecutive hops. Second, the memory cost of eight-octant partitioning and grouping is high. Third, the covariance matrix computation in the Saab transform is computationally intensive with a higher memory cost. To address them, we propose several changes.

Proposed Feature Extractor. As shown in Fig. 1 (enclosed by the orange block), the new processing module contains random sampling (RS), K -NN, relative point positional encoding [13], max pooling, and point feature standardization. First, we use RS rather than FPS between hops for faster computation of large-scale point clouds. Second, to reduce the memory cost of the eight-octant partitioning and grouping, we adopt max pooling. Since the feature dimension remains the same with max pooling prevents, no dimension reduction via the Saab transform is needed, which further helps save memory and time cost. However, max pooling may drop important information occasionally. To address it, we add relative point positional encoding before max pooling to ensure that point features are aware of their relative spatial positions. Specifically, for point p_i and its K neighbors $\{p_i^1, \dots, p_i^k, \dots, p_i^K\}$, the relative point position r_i of each neighbor p_i^k is encoded as

$$r_i = p_i \oplus p_i^k \oplus (p_i - p_i^k) \oplus \|p_i - p_i^k\|, \quad (1)$$

where \oplus denotes concatenation, and $\|\cdot\|$ is the Euclidean distance. We will show that the new pipeline is much more economic than PointHop, PointHop++ and deep-learning methods in terms of memory consumption and computational complexity in Sec. 4.

Classifier. The S3DIS dataset is highly imbalanced in point labels. Among the 13 object categories, the ceiling, floor and wall three classes take around 75% of the data. We adopt the XGBoost classifier [19] to help reduce the influence of imbalanced data. Other classifiers such as Linear Least Square Regression, Random Forest and Support Vector Machine

(SVM) demand higher computational speed and memory cost in the large-scale point cloud segmentation problem. Overall, XGBoost can handle the large-scale point cloud data with good performance, fast speed and lower memory requirement.

4 Experiments

Experimental Setup. We adopt the following setting to evaluate the semantic segmentation performance for the S3DIS dataset. The grid size is 0.04 in voxel-based downsampling. The feature extractor has 4 hops. We set $K = 64$ in K -NN search. Some methods may choose smaller K to save computational complexity, i.e., 16, 20 and 32, but our method can afford a larger K . Because the input is a large-scale indoor scene, a larger K will lead to a larger receptive field which helps learn the structure of the scene. Thus, we choose 64 here. The subsampling ratios between two consecutive hops are 0.25, 0.25, 0.5, 0.5, respectively. Three nearest neighbors’ features are used to interpolate in the upsampling module. For example, to upsample from $N/64$ points to $N/32$ points (see Fig.1), we first find the three nearest points in the $N/64$ point set for each point in the $N/32$ point set. Then, we average the features of the three points. In this way, we get the point features for $N/32$ points. The output features are truncated to 32-bit before fed into the XGBoost classifier. The standard 6-fold cross validation is used in the experiment, where one of the six areas is selected as test area in each fold. By following [1], we consider two the evaluation metrics: the mean Intersection-over-Union (mIoU) and the Overall Accuracy (OA) of the total 13 classes.

Comparison of Data Pre-processing Methods. The statistics of the S3DIS dataset using three data pre-processed methods are compared in Table 2. The proposed room-style method has more points in each input unit, i.e., input scale, ranging from 10K to 200K. We also compare the data size when the training data are collected from areas 1-5 and the test data are collected from area 6. The total number of training points of the room-style method is 18.75% of the block-style method and 12.5% of the view-style method. Points inside each unit of the three methods are visualized in Fig. 2, which includes a conference room and an office from area 1. The block-style method loses the structure of a complete room. As to the view-style method, the four views of a conference room overlap a lot with each other, producing significant redundancy. The units obtained by the room-style method look more natural. They offer a view of the whole room while keeping a small data size.

Table 2: Comparison of data statistics of three pre-processing methods.

Method		Block	View	Room
Input Scale ($\times 10^3$)		4	40	10-200
Total Data Size ($\times 10^6$)	Train	80	120	15
	Test	10	80	2

Semantic Segmentation Performance. We compare the semantic segmentation performance on the S3DIS dataset of PointNet and the proposed GSIP in Table 3, where the better results are shown in bold. It is common to use area 5 as the test. Thus, we show performance for area 5 as well as the 6-fold. As shown in the table, GSIP outperforms PointNet in mIoU by 2.7% and 0.9% for area 5 and 6-fold, respectively. The cross validation results of all 6 areas of GSIP are shown in Table 4. We see that area 6 is the easiest one (64.5% mIoU and 86.5% OA) while area 2 is the hardest one (31% and 68.8%). The mIoU and OA over 13 classes averaged by the 6-fold are 48.5% and 79.8%, respectively. Visualization of GSIP’s segmentation results and the ground truth of two room in area 6 is given in Fig. 3. The ceiling is removed to show inner objects clearly.

It is worthwhile to point out that some categories have extremely low mIoU (close to 0%) in more than 2 areas. For example, sofa got 1.7%, 6.4%, 4.4%, and 3.4% mIoU in areas 2, 3, 4 and 5. This is attributed to data imbalance. The data imbalance problem is commonly seen in large-scale segmentation tasks. For example, in a regular room, it is highly possible that more points are from the wall, ceiling and floor, while less points from chairs, desks, and other small objects. The beam, column, sofa, and board are even less common. Without seeing enough samples during training, the decrease in prediction performance is valid.

Comparison of Model and Time Complexities. We compare the model size and training time complexity of GSIP and PointNet in Table 5. PointNet takes 22 hours to train in a single GeForce GTX TITAN X GPU. GSIP takes around 40 minutes for feature extraction with a Intel(R)Xeon(R) CPU and 20 minutes to train the XGBoost classifier with 4 GeForce GTX TITAN X GPUs. The total training time is around 1 hour for each fold, which is much faster than PointNet. If we use a single GPU for XGBoost training, the overall training time is still significantly less than PointNet. To justify our claim, we calculate the time complexity of algorithm theoretically in terms of big O . The proposed method is composed of data pre-processing, feature extractor and classifier. For each sample with N points,

Table 3: Comparison of semantic segmentation performance (%) for S3DIS.

	Method	OA	mIoU	ceiling	floor	wall	beam	column	window	door	table	chair	sofa	bookcase	board	clutter
Area 5	PointNet	-	41.1	88.0	97.3	69.8	0.05	3.90	42.3	10.8	58.9	52.6	5.90	40.3	26.4	33.2
	GSIP	79.9	43.8	89.2	97.0	72.2	0.10	18.4	37.3	22.5	64.3	59.5	3.40	47.2	22.9	35.7
6-fold	PointNet	78.5	47.6	88.0	88.7	69.3	42.4	23.1	47.5	51.6	54.1	42.0	9.60	38.2	29.4	35.2
	GSIP	79.8	48.5	91.8	89.8	73.0	26.3	24.0	44.6	55.8	55.5	51.1	10.2	43.8	21.8	43.2

Table 4: Class-wise semantic segmentation performance (%) of GSIP for S3DIS.

	1	2	3	4	5	6	mean
ceiling	91.9	88.9	95.2	89.9	89.2	95.6	91.8
floor	93.7	58.5	97.7	95.2	97.0	96.8	89.8
wall	71.5	70.9	73.3	72.5	72.2	77.5	73.0
beam	21.7	5.90	64.9	0.20	0.10	65.2	26.3
column	38.7	12.8	20.9	15.9	18.4	37.0	24.0
window	77.5	18.3	39.7	15.1	37.3	79.5	44.6
door	71.5	46.0	69.0	51.7	22.5	73.8	55.8
table	61.7	23.3	62.5	49.8	64.3	71.3	55.5
chair	49.4	21.8	60.4	49.3	59.5	66.3	51.1
sofa	20.4	1.70	6.40	4.40	3.40	24.7	10.2
bookcase	41.2	22.8	58.1	34.4	47.2	58.9	43.8
board	29.4	5.00	22.8	14.2	22.9	36.7	21.8
clutter	46.6	27.6	51.3	42.4	35.7	55.7	43.2
mIoU	55.0	31.0	55.6	41.2	43.8	64.5	48.5
OA	81.1	68.8	83.9	78.8	79.9	86.5	79.8

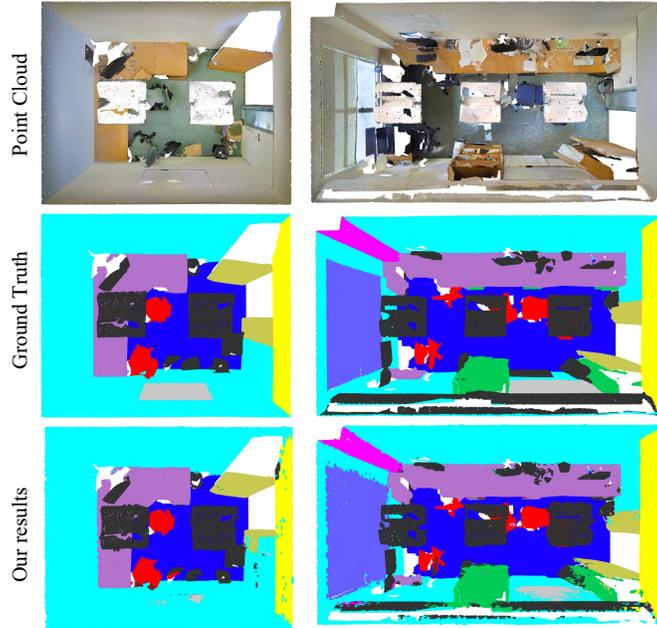


Figure 3: Qualitative results of the proposed GSIP method.

1. data pre-processing: voxel downsampling takes $O(N \log N)$; geometric feature calculation takes $O(NK^2)$, K is the number of nearest points;
2. feature extraction: random sampling takes $O(1)$; K -NN search takes $O(KND)$, D is feature dimension; relative point position encoding takes $O(KND)$; max pooling takes $O(N)$. To sum up, the feature extractor costs $O(KND)$;
3. XGBoost classifier: the XGBoost classifier’s complexity can be found in their paper, which is not our focus.

For the algorithms involved, brute force time complexities are used for calculation. There are better implementations that can optimize the complexity, which are not considered here.

As to the model size, GSIP has 24K parameters while PointNet has 900K parameters. PointNet is an end-to-end supervised method, which uses some fully connected layers at the end of their pipeline to predict point labels. We do not break the pipeline into the feature extractor and the classifier for comparison. The parameters mainly come from the XGBoost classifier, which has 128 trees and the maximum depth of a tree is 6. For each tree, there are 2 parameters per intermediate node and 1 parameter per leaf node. The feature extractor only has 2 parameter in each hop, mean and standard deviation, in point feature standardization.

Table 5: Comparison of time and model complexities.

Method	Device	Training time	Parameter No.
PointNet	GPU	22 hours	900169
GSIP	Feature	CPU	40 minutes
	XGBoost	GPU	20 minutes

Other Comparisons. As discussed in Sec. 3.2, the new feature extractor is more effective than the PointHop feature extractor in terms of computational and memory efficiency. It is worthwhile to compare the segmentation performance of the two to see whether there is any performance degradation. We compare the effectiveness of the new feature extractor and the PointHop feature extractor under the same GSIP framework for the S3DIS dataset in Table 6. Their performance is comparable as shown in the table. Actually, the new one achieves slightly better performance. Furthermore, we compare the quantization effect of extracted features before feeding them into the XGBoost classifier in Table 7. We see little performance degradation for features to be quantized to 16 or 32 bits. Thus, we can save computation and memory for classifier training and testing by using 16-bit features.

Table 6: Performance comparison of two feature extractors (%).

Method	mIoU	OA
GSIP	48.5	79.8
PointHop	47.9	79.1

Table 7: Impact of quantized point-wise features (%).

Quantization	mIoU	OA
32-bit	64.5	86.5
16-bit	64.9	86.5

5 Conclusion and Future Work

A green semantic segmentation method for large-scale indoor point clouds, called GSIP, was proposed in this work. It contains two novel ingredients: a new room-style method for data pre-processing and a new point cloud feature extractor which is extended from PointHop with lower memory and computational costs while preserving the segmentation performance. We evaluated the performance of GSIP against PointNet with the indoor S3DIS dataset and showed that GSIP outperforms PointNet in terms of performance accuracy, model sizes and computational complexity. As to future extension, it is desired to generalize the GSIP method from the large-scale indoor point clouds to the large-scale outdoor point clouds. The latter has many real world applications. Furthermore, the data imbalance problem should be carefully examined so as to boost the segmentation and/or classification performance.

References

- [1] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017.
- [2] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413*, 2017.
- [3] Mingyang Jiang, Yiran Wu, Tianqi Zhao, Zelin Zhao, and Cewu Lu. Pointsift: A sift-like network module for 3d point cloud semantic segmentation. *arXiv preprint arXiv:1807.00652*, 2018.
- [4] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. *Advances in neural information processing systems*, 31:820–830, 2018.
- [5] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019.
- [6] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015.
- [7] Iro Armeni, Ozan Sener, Amir R. Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition*, 2016.
- [8] Timo Hackel, Nikolay Savinov, Lubor Ladicky, Jan D Wegner, Konrad Schindler, and Marc Pollefeys. Semantic3d. net: A new large-scale point cloud classification benchmark. *arXiv preprint arXiv:1704.03847*, 2017.
- [9] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Jurgen Gall. Semantickitti: A dataset for semantic scene understanding of lidar sequences. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9297–9307, 2019.
- [10] Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4558–4567, 2018.
- [11] Dario Rehg, Johanna Wald, Jurgen Sturm, Nassir Navab, and Federico Tombari. Fully-convolutional point networks for large-scale point clouds. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 596–611, 2018.
- [12] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *arXiv preprint arXiv:2012.09688*, 2020.
- [13] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. Randla-net: Efficient semantic segmentation of large-scale point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11108–11117, 2020.
- [14] Min Zhang, Haoxuan You, Pranav Kadam, Shan Liu, and C-C Jay Kuo. Pointhop: An explainable machine learning method for point cloud classification. *IEEE Transactions on Multimedia*, 22(7):1744–1755, 2020.
- [15] Min Zhang, Yifan Wang, Pranav Kadam, Shan Liu, and C-C Jay Kuo. Pointhop++: A lightweight learning model on point sets for 3d classification. In *2020 IEEE International Conference on Image Processing (ICIP)*, pages 3319–3323. IEEE, 2020.
- [16] Min Zhang, Pranav Kadam, Shan Liu, and C-C Jay Kuo. Unsupervised feedforward feature (uff) learning for point cloud classification and segmentation. In *2020 IEEE International Conference on Visual Communications and Image Processing (VCIP)*, pages 144–147. IEEE, 2020.
- [17] Pranav Kadam, Min Zhang, Shan Liu, and C-C Jay Kuo. Unsupervised point cloud registration via salient points analysis (spa). In *2020 IEEE International Conference on Visual Communications and Image Processing (VCIP)*, pages 5–8. IEEE, 2020.
- [18] Pranav Kadam, Min Zhang, Shan Liu, and C-C Jay Kuo. R-pointhop: A green, accurate and unsupervised point cloud registration method. *arXiv preprint arXiv:2103.08129*, 2021.
- [19] Tianqi Chen, Tong He, Michael Benesty, Vadim Khotilovich, Yuan Tang, Hyunsu Cho, et al. Xgboost: extreme gradient boosting. *R package version 0.4-2*, 1(4):1–4, 2015.
- [20] Radu Bogdan Rusu, Nico Blodow, and Michael Beetz. Fast point feature histograms (fpfh) for 3d registration. In *2009 IEEE international conference on robotics and automation*, pages 3212–3217. IEEE, 2009.

- [21] Federico Tombari, Samuele Salti, and Luigi Di Stefano. Unique signatures of histograms for local surface description. In *European conference on computer vision*, pages 356–369. Springer, 2010.
- [22] Timo Hackel, Jan D Wegner, and Konrad Schindler. Fast semantic segmentation of 3d point clouds with strongly varying density. *ISPRS annals of the photogrammetry, remote sensing and spatial information sciences*, 3:177–184, 2016.
- [23] Loic Landrieu, Hugo Raguét, Bruno Vallet, Clément Mallet, and Martin Weinmann. A structured regularization framework for spatially smoothing semantic labelings of 3d point clouds. *ISPRS Journal of Photogrammetry and Remote Sensing*, 132:102–118, 2017.
- [24] Clément Mallet, Frédéric Bretar, Michel Roux, Uwe Soergel, and Christian Heipke. Relevance assessment of full-waveform lidar data for urban area classification. *ISPRS journal of photogrammetry and remote sensing*, 66(6):S71–S84, 2011.
- [25] Nesrine Chehata, Li Guo, and Clément Mallet. Airborne lidar feature selection for urban classification using random forests. In *Laserscanning*, 2009.
- [26] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.
- [27] C-C Jay Kuo. Understanding convolutional neural networks with a mathematical model. *Journal of Visual Communication and Image Representation*, 41:406–413, 2016.
- [28] C-C Jay Kuo and Yueru Chen. On data-driven saak transform. *Journal of Visual Communication and Image Representation*, 50:237–246, 2018.
- [29] C-C Jay Kuo, Min Zhang, Siyang Li, Jiali Duan, and Yueru Chen. Interpretable convolutional neural networks via feedforward design. *Journal of Visual Communication and Image Representation*, 60:346–359, 2019.
- [30] Yueru Chen, Zhuwei Xu, Shanshan Cai, Yujian Lang, and C-C Jay Kuo. A saak transform approach to efficient, scalable and robust handwritten digits recognition. In *2018 Picture Coding Symposium (PCS)*, pages 174–178. IEEE, 2018.
- [31] Yueru Chen and C-C Jay Kuo. Pixelhop: A successive subspace learning (ssl) method for object recognition. *Journal of Visual Communication and Image Representation*, 70:102749, 2020.
- [32] Yueru Chen, Mozhdeh Rouhsedaghat, Suyu You, Raghuveer Rao, and C-C Jay Kuo. Pixelhop++: A small successive-subspace-learning-based (ssl-based) model for image classification. In *2020 IEEE International Conference on Image Processing (ICIP)*, pages 3294–3298. IEEE, 2020.
- [33] Yijing Yang, Vasileios Magoulianitis, and C-C Jay Kuo. E-pixelhop: An enhanced pixelhop method for object classification. *arXiv preprint arXiv:2107.02966*, 2021.
- [34] Mozhdeh Rouhsedaghat, Yifan Wang, Xiou Ge, Shuowen Hu, Suyu You, and C-C Jay Kuo. Facehop: A light-weight low-resolution face gender classification method. *arXiv preprint arXiv:2007.09510*, 2020.
- [35] Hong-Shuo Chen, Mozhdeh Rouhsedaghat, Hamza Ghani, Shuowen Hu, Suyu You, and C-C Jay Kuo. Defakehop: A light-weight high-performance deepfake detector. In *2021 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2021.
- [36] Kaitai Zhang, Bin Wang, Wei Wang, Fahad Sohrab, Moncef Gabbouj, and C-C Jay Kuo. Anomalyhop: An ssl-based image anomaly localization method. *arXiv preprint arXiv:2105.03797*, 2021.
- [37] Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017.
- [38] Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52, 1987.