# Sparse Single Sweep LiDAR Point Cloud Segmentation via Learning Contextual Shape Priors from Scene Completion

**Xu Yan**[1,2†], **Jiantao Gao**[2,4 †], **Jie Li**[1,3], **Ruimao Zhang**[1,2], **Zhen Li**[1,2 *],
**Rui Huang**[1,3], **Shuguang Cui**[1,2]

[1] The Chinese University of Hong Kong (Shenzhen), [2] Shenzhen Research Institute of Big Data,
[3] Shenzhen Institute of Artificial Intelligence and Robotics for Society, [4] Shanghai University
{*xuyan1@link., lizhen@*}*cuhk.edu.cn*

## Abstract

LiDAR point cloud analysis is a core task for 3D computer vision, especially for autonomous driving. However, due to the severe sparsity and noise interference in the single sweep LiDAR point cloud, the accurate semantic segmentation is nontrivial to achieve. In this paper, we propose a novel sparse LiDAR point cloud semantic segmentation framework assisted by learned contextual shape priors. In practice, an initial semantic segmentation (SS) of a single sweep point cloud can be achieved by any appealing network and then flows into the semantic scene completion (SSC) module as the input. By merging multiple frames in the LiDAR sequence as supervision, the optimized SSC module has learned the contextual shape priors from sequential LiDAR data, completing the sparse single sweep point cloud to the dense one. Thus, it inherently improves SS optimization through fully end-to-end training. Besides, a Point-Voxel Interaction (PVI) module is proposed to further enhance the knowledge fusion between SS and SSC tasks, i.e., promoting the interaction of incomplete local geometry of point cloud and complete voxel-wise global structure. Furthermore, the auxiliary SSC and PVI modules can be discarded during inference without extra burden for SS. Extensive experiments confirm that our JS3C-Net achieves superior performance on both SemanticKITTI and SemanticPOSS benchmarks, i.e., 4% and 3% improvement correspondingly.

## Introduction

LiDAR point clouds, compared with data from other sensors, such as cameras and radars in autonomous driving perception, have advantages of both accurate distance measurements and fine semantic descriptions. Semantic segmentation of LiDAR point clouds is usually conducted by assigning a semantic class label to each point. It is traditionally viewed as a typical task in the computer vision community. In autonomous driving, accurate and effective point cloud semantic segmentation undoubtedly plays a critical role.

Previous studies (Thomas et al. 2019; Wu, Qi, and Fuxin 2019) about point cloud semantic segmentation mainly focused on the complete or dense point cloud scenarios, which are post-processed by merging multiple collected LiDAR or RGB-D sequences (e.g., ScanNet (Dai et al. 2017),
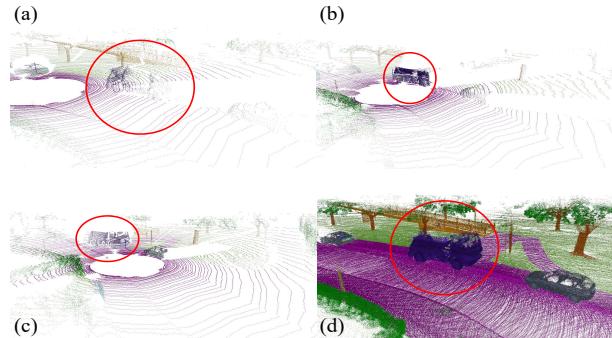
---

Figure 1: **Learning shape priors from multiple frames.** For the sparse per-sweep point cloud shown in (a), it is nontrivial for current methods to recognize the truck from partial components. However, if we introduce the auxiliary information from adjacent frames (b) and (c), it is much easier to segment the complete truck in (d).

S3DIS (Armeni et al. 2016) and Semantic3D (Hackel et al. 2017)). However, raw per-sweep LiDAR point clouds, as the original input of autonomous driving, are much sparser. Their sparsity usually increases with the reflection distance, which often leads to extremely shapes missing and uneven point sampling for various categories. Therefore, despite the promising performance on complete data (e.g., 80% mIOU on Semantic3D), the semantic segmentation of sparse single sweep LiDAR point cloud still remains a big challenge, which extremely limits its accuracy in real applications.

In this paper, we try to break through the barrier of semantic segmentation on sparse single sweep LiDAR point clouds. One plausible way to solve this problem is to fully utilize the sequential nature of LiDAR data. Taking the scenario shown in Fig. 1(a) as an example, for a per-sweep point cloud with extremely sparse points of the truck, it seems impossible for previous methods to conduct accurate segmentation. Nevertheless, such segmentation would be possible, if we introduce the richer shape information from the other two frames, i.e., Fig. 1(b) and Fig. 1(c), to reconstruct a shape-complete truck as shown in Fig. 1(d). For this purpose, some previous works utilized historical adjacent frames to supplement the local details missing from the point clouds. For instance, SpSequenceNet (Shi et al.

2020) and MeteorNet (Liu, Yan, and Bohg 2019) use the point cloud of the current frame to query the nearest neighbors from the previous frames, following which a feature aggregation is conducted to fuse the adjacent-frame information. PointRNN (Fan and Yang 2019) applies Recurrent Neural Networks (RNNs) to select available features from previous scenes. However, all of the above methods become unavailable in most real scenarios since the following reasons: (1) These methods exclusively use historical frames of the current scene in LiDAR sequence. Thus, they cannot introduce priors for newly incoming objects in this scene, i.e., they cannot utilize future frames. (2) Their proposed feature aggregation methods (i.e., through kNN or RNN) inevitably increase the computational burden, which makes it less effective and unsuitable for self-driving task.

To solve the above issues, we propose an enhanced **J**oint single sweep LiDAR point cloud **S**emantic **S**egmentation by exploiting learned shape prior form **S**cene **C**ompletion network, i.e., **JS3C-Net**. Specifically, by merging dozens of consecutive frames in a LiDAR sequence, a large complete point cloud is achieved as ground truth for the Semantic Scene Completion (SSC) task without extra annotation. The optimized SSC by using these annotations could capture the compelling shape priors, making the incomplete input complete to the acceptable shape with semantic labels (Song et al. 2017). Therefore, the completed shape priors can inherently benefit the semantic segmentation (SS) due to the fully end-to-end training strategy. Furthermore, a well-designed Point-Voxel Interaction (PVI) module is further proposed for implicit mutual knowledge fusion between the SS and SSC tasks. Concretely, the point-wise segmentation and voxel-wise completion are leveraged to maintain the coarse global structure and fine-grained local geometry through PVI module. More importantly, we design our SSC and PVI modules to be **disposable**. To achieve this, JS3C-Net combines the SS and SSC in a cascaded manner, which means that it would not influence the information flow for SS while discarding the SSC and PVI modules in inference stage. Thus, it can prevent bringing the extra computing burden from generating complete high-resolution dense volumes.

Our main contributions are: 1) To the best of our knowledge, the proposed **JS3C-Net** is the first to achieve the enhanced sparse single sweep LiDAR semantic segmentation via auxiliary scene completion. 2) For better trade-off between performance and effectiveness, our auxiliary components are designed in cascaded and disposable manners, and a novel point-voxel interaction (PVI) module is proposed for better feature interaction and fusion between the two tasks. 3) Our method shows superior results in both SS and SSC on two benchmarks, i.e., SemanticKITTI (Behley et al. 2019) and SemanticPOSS (Pan et al. 2020), by a large margin.

## Related Work

**Point Cloud Semantic Segmentation.** Unlike 2D images with regular grids, point clouds are often sparse and disordered. Thus, point clouds processing is a challenging task. There are three main strategies to approach this problem: **projection-based**, **voxel-based** and **point-based**. (1)

Projection-based methods map point clouds onto 2D pixels, so that traditional CNN can play a normal role. Previous works projected all points scanned by the rotating LiDAR onto 2D images by plane projection (Lawin et al. 2017; Boulch, Le Saux, and Audebert 2017; Tatarchenko et al. 2018) or spherical projection (Wu et al. 2018, 2019). (2) Considering the sparsity of point clouds and memory consumption, it is not very effective to directly voxelize point clouds and then use 3D convolution for feature learning. Various subsequent improved methods have been proposed, e.g., efficient spatial sparse convolution (Choy, Gwak, and Savarese 2019; Graham, Engelcke, and van der Maaten 2018) and octree based convolutional neural networks (Wang et al. 2017; Riegler, Osman Ulusoy, and Geiger 2017). Also (Tang et al. 2020) use NAS to obtain a more efficient feature representation. (3) Point-based methods directly process raw point clouds (Qi et al. 2017a,b). Usually, most methods use sampling strategies to select sub-points from the original point clouds, and then use local grouping with feature aggregation function for local feature learning of each sub-point. Among these methods, graph-based learning (Wang et al. 2019a; Landrieu and Simonovsky 2018; Landrieu and Boussaha 2019; Wang et al. 2019c) and convolution-like operations (Thomas et al. 2019; Wu, Qi, and Fuxin 2019; Hu et al. 2020) are widely used. However, previous methods often suffer from local information missing in LiDAR scenes due to insufficient priors and bias in data collection.

**Semantic Scene Completion.** Semantic scene completion (SSC) aims to produce a complete 3D voxel representation from an incomplete input. Concretely, Song *et al.*(Song et al. 2017) firstly use single-view depth as input to construct an end-to-end model SSCNet, which can predict the results of scene completion and semantic labeling simultaneously. Spatially sparse group convolution is used in Zhang *et al.* (Zhang et al. 2018) for fast 3D dense prediction. Meanwhile, coarse-to-fine strategies (e.g., LSTM based model) are used in (Dai et al. 2018; Han et al. 2017) to recover missing parts of 3D shapes. More recently, some works introduce color information (Garbade et al. 2019a) and use more powerful two-stream feature extractor (Li et al. 2019) or feature fusion (Liu et al. 2020) to enhance the performance. However, SSC is rarely studied in large-scale LiDAR scenarios, and the serious geometric details missing and real-time requirements make it difficult.

**Multi-task Learning on Segmentation.** Multi-task learning aims to improve the learning efficiency and prediction accuracy for each task through knowledge transfer, which is widely used in 2D images segmentation (Kendall, Gal, and Cipolla 2018). Wang *et al.* (Wang et al. 2019b), Pham *et al.* (Pham et al. 2019) and Wei *et al.* (Wei et al. 2020) innovatively combine semantic and instance segmentation with specific-designed fusion modules to improve the performance. OccuSeg (Han et al. 2020) proposes a 3D voxel projection-based segmentation network with voxel occupancy size regression and owns advantages of robustness in prediction. However, the shape priors brought by completion tasks are often neglected in previous works, while a proper
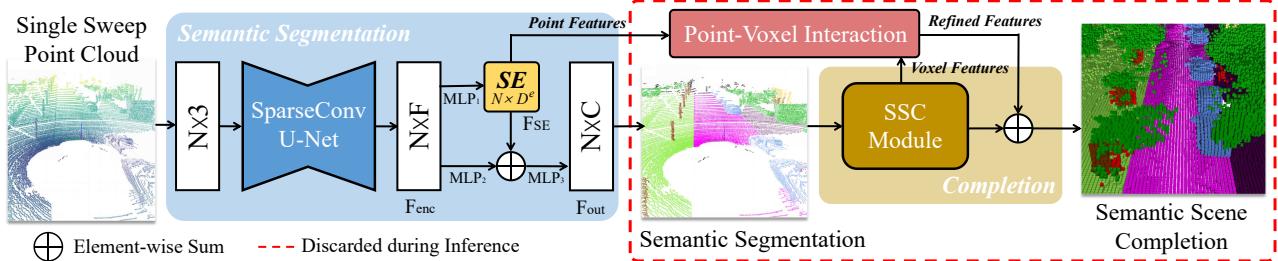
Figure 2: **Overall pipeline of JS3C-Net**. Given a sparse incomplete single sweep point cloud, it firstly uses a sparse convolution U-Net to conduct point feature encoding $F_{enc}$. Based on the initial encoding, $MLP_1$ is used to generate shape embedding (SE) $F_{SE}$, which flows into $MLP_3$ together with initial encoding transferred through $MLP_2$ to generate $F_{out}$ for point cloud semantic segmentation. Afterwards, the incomplete fine-grained point features from SE and complete voxel features from semantic scene completion (SSC) module flows into the Point-Voxel Interaction (PVI) module to achieve the refined features, which finally outputs the completion voxels with supervision. Note that the SSC and PVI modules can be discarded during inference.

use could improve the performance of segmentation.

## Method

### Overview

The pipeline of **JS3C-Net**[1] is illustrated in Fig. 2. In practice, we firstly use the general appealing point cloud segmentation network to obtain initial point semantic segmentation and a shape embedding (SE) for each incomplete single-frame point cloud. Then SSC module takes results of segmentation network as input and generates the completed voxel of the whole scene with dense convolution neural network. Meanwhile, a point-voxel interaction (PVI) module is proposed to conduct shape-aware knowledge transfer.

### Semantic Segmentation

In general, a point cloud has two components: the points $\mathcal{P} \in \mathbb{R}^{N \times 3}$ and their features $\mathcal{F} \in \mathbb{R}^{N \times D}$, where the points record spatial coordinates of $N$ points and $D$-dimensional features can include any point-wise information, e.g., RGB information. Here we only use point coordinates as inputs.

For semantic segmentation stage, we simply choose Submanifold SparseConv (Graham, Engelcke, and van der Maaten 2018) as our backbone. Unlike traditional voxel-based methods (Ronneberger, Fischer, and Brox 2015; Choy, Gwak, and Savarese 2019) directly transforming all points into the 3D voxel grids by averaging all input features, it only stores non-empty voxels by the Hash table, and conduct convolution operations only on these non-empty voxels with more efficient way. Afterwards, the voxel-based output from sparse convolution based U-Net (SparseConv U-Net) is transformed back to the point-wise features $\mathcal{F}_{enc} \in \mathbb{R}^{N \times F}$ by nearest-neighbor interpolation. To further introduce shape priors (see latter section) to point-wise features, we use multi-layer perceptions ($MLP_1$) to transfer their features to shape embedding (SE) $F_{SE} \in \mathbb{R}^{N \times D^e}$, which works as the input of the subsequent point-voxel interaction module. Furthermore, an element-wise addition operation after $MLP_2$ is used to fuse shape embedding with features from SparseConv U-Net. Finally, $F_{out} \in \mathbb{R}^{N \times C}$ are

---

[1]https://github.com/yanx27/JS3C-Net.

generated by $MLP_3$ and prepare for further semantic scene completion stage.

### Cascaded Semantic Scene Completion

The Semantic Scene Completion (SSC) module aims to introduce contextual shape priors from the entire LiDAR sequence. For stage of SSC, it takes the semantic probability $F_{out}$ from the SparseConv as input, and then predict the completion results.

The architecture of our SSC module is depicted in Fig. 3 (a). Taking an incomplete point cloud with per points categorical probability as input, the network firstly conducts voxelization to obtain high-resolution 3D volume, and uses one convolution layer following by a pooling layer to reduce the resolution and the complexity of computation. Then, several basic blocks using convolutions with skip-connection are exploited to learn a local geometry representation. Afterwards, the features from different scales are concatenated to aggregate information from multiple scales. For achieving original resolution of SSC output, we leverage dense upsampling (Liu et al. 2018) shown in Fig. 3(a) to avoid the interpolation inaccuracy, instead of dilation convolution based upsampling (Song et al. 2017). Finally, we obtain a voxel output with $C + 1$ channels ($C$ semantic categories label and one non-object label). This coarse completion will be fed into the PVI module for further mutual enhancements.

### Shape-aware Point-Voxel Interaction

To fully utilize implicit knowledge transfer for mutual improvements of two tasks, we innovatively propose a shape-aware Point-Voxel Interaction (PVI) module for knowledge fusion between incomplete point clouds and complete voxels from former two steps. Although the SSC module can generate voxel-wise output with complete shape, such output is relatively coarse due to the voxelization procedure, leading to local geometric details missing. The entire geometric representation in raw point cloud data can, nevertheless, provide semantic guidance during completion process despite missing parts.

The inner structure of PVI module is shown in Fig. 3 (b), which aims to conduct a coarse-to-fine process for the SSC
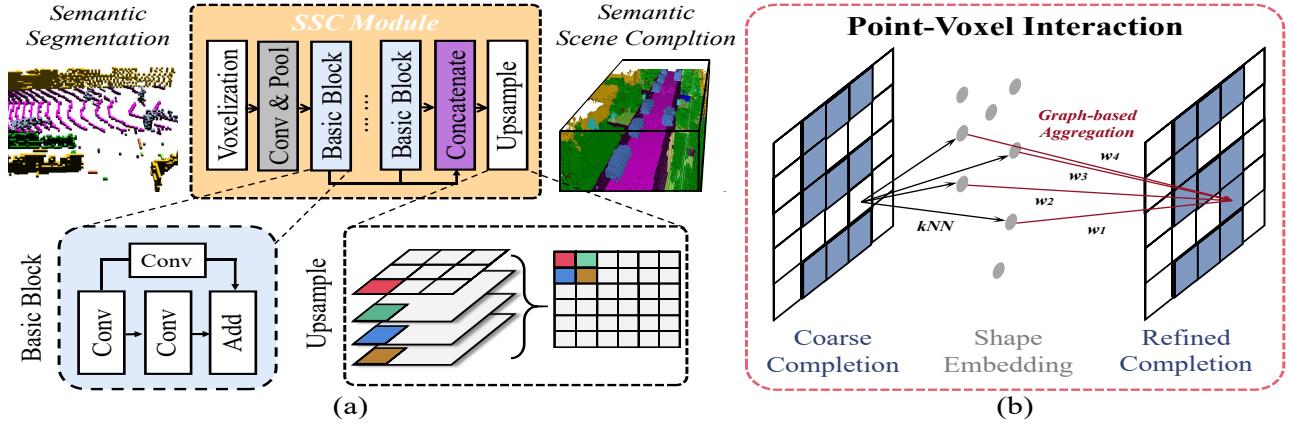
Figure 3: Part (a) shows the inner structure of SSC module, which uses semantic probability from segmentation network as inputs, generating complete volume by several convolution blocks and dense upsample. Part (b) illustrates a 2D case of PVI module, which uses the center points of the coarse global structure of number '5' to query $k$ nearest neighbors from the raw point cloud, and then applies graph-based aggregation to achieve the completed '5' through fine-grained local geometry.

prediction. To be more precise, per point shape embedding $\mathcal{F}_{SE} \in \mathbb{R}^{N \times D^e}$ and coarse completion from SSC module $\mathcal{V}$ flow into PVI module as input. Afterwards, PVI firstly selects geometric centers of all non-empty voxels from $\mathcal{V}$ as a new point cloud $\mathcal{P}^v \in \mathbb{R}^{N' \times (C+1)}$, then it uses k-nearest neighbor by Euclidean distance to query the closest points from original point cloud $\mathcal{P}$. To this end, a graph convolutional network is further employed to enhance the relation learning between $\mathcal{P}^v$ and $\mathcal{P}$ in both spatial and semantic spaces. In particular, the nodes of the graph are defined by the point positions with associated points features. For each point $p_i^v \in \mathcal{P}^v$ and its $j$-th neighboring point $p_j \in \mathcal{P}$, we adopt the convolutional operator from DGCNN (Wang et al. 2019c) to define edge-features $e_{ij}$ between two points as:

$$e_{ij} = \phi([p_i^v, f_i^v], [p_i^v, f_i^v] - [p_j, f_j]), \qquad (1)$$

where $f_i^v$ and $f_j$ are features of point $p_i^v$ and $p_j$ respectively, and $[\cdot, \cdot]$ means concatenation operation. The share-weighted non-linear fuction $\phi$ is the multi-layer perceptron (MLP) in this paper (any differentiable architecture alternative). Finally, by stack $l$ graph convolutional network (GCN) layers, we obtain final fine-grained completion.

The feature interaction process enables features of sparse point cloud to acquire the ability to predict semantics of complete voxels. Therefore, the information on complete details can positively affect the segmentation part through back propagation. Furthermore, PVI module enhances the probability of predicting whether the corresponding voxel of $p_i^v$ represents an object in the fine-grained architecture, which fully utilizes spatial and semantic relationships between each $p_j$ and $p_i^v$. Finally, this enhanced feature will be added to the original coarse completion output through a residual connection for further refinement (see the refined number '5' in Fig. 3 (b)).

### Uncertainty-weighted Multi-task Loss

In order to further balance these two tasks and avoid complicated manual attempts during the end-to-end training, we use the uncertainty weighting method proposed in (Kendall, Gal, and Cipolla 2018). It introduces acquirable parameters to automatically adjust the optimal proportion between different tasks. Specifically, the joint loss can be written as:

$$\mathcal{L}(W, \sigma_1, \sigma_2) = \frac{1}{2\sigma_1^2}\mathcal{L}_{seg}(W_{seg}) + \frac{1}{2\sigma_2^2}\mathcal{L}_{complet}(W) \\ + log\sigma_1 + log\sigma_2, \qquad (2)$$

where losses $\mathcal{L}_{seg}$ for segmentation and $\mathcal{L}_{complet}$ for completion are both weighted cross-entropy losses exploited to update the network parameters $W$. Note that gradients from segmentation outputs are only conducted on parameters of segmentation network $W_{seg} \in W$. In addition, we use trainable parameters $\sigma_1$ and $\sigma_2$ to weight the proportion of these two tasks for optimal trade-off. Their uncertainty can be deduced as two $log$ term to control their values. Finally, during the training process, these two tasks will promote each other through back-propagation of joint learning.

### Disposable Properties of Auxiliary Components

Our proposed JS3C-Net is a general joint learning framework to improve the point cloud segmentation by introducing complete shape extracted by LiDAR sequence itself. The network used in semantic segmentation stage is flexible and can be replaced by other appealing networks. Furthermore, our JS3C-Net is effective enough for real-time applications, since the auxiliary components (i.e., SSC module and PVI module) can be discarded during inference to prevent introducing any computing burden for segmentation. That is to say, the completion part are *only* exploited in the training process as the dotted line shown in Fig. 2.

## Experiments

### Dataset

To further verify the effectiveness of our method, we evaluate JS3C-Net on two benchmarks SemanticKITTI (Behley

Table 1: Semantic segmentation on the *SemanticKITTI* benchmark. Underline marks results that $\sim 10\%$ higher than baseline.

| Method | Size | mIoU | road | sidewalk | parking | other-ground | building | car | truck | bicycle | motorcycle | other-vehicle | vegetation | trunk | terrain | person | bicyclist | motorcyclist | fence | pole | traffic sign |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SqueezeSegV2 (Wu et al. 2019) | | 39.7 | 88.6 | 67.6 | 45.8 | 17.7 | 73.7 | 81.8 | 13.4 | 18.5 | 17.9 | 14.0 | 71.8 | 35.8 | 60.2 | 20.1 | 25.1 | 3.9 | 41.1 | 20.2 | 26.3 |
| DarkNet53Seg (Behley et al. 2019) | 64*2048 pixels | 49.9 | **91.8** | 74.6 | 64.8 | 27.9 | 84.1 | 86.4 | 25.5 | 24.5 | 32.7 | 22.6 | 78.3 | 50.1 | 64.0 | 36.2 | 33.6 | 4.7 | 55.0 | 38.9 | 52.2 |
| RangeNet53++ (Milioto et al. 2019) | | 52.2 | **91.8** | **75.2** | **65.0** | 27.8 | 87.4 | 91.4 | 25.7 | 25.7 | 34.4 | 23.0 | 80.5 | 55.1 | 64.6 | 38.3 | 38.8 | 4.8 | 58.6 | 47.9 | 55.9 |
| 3D-MiniNet (Alonso et al. 2020) | | 55.8 | 91.6 | 74.5 | 64.2 | 25.4 | 89.4 | 90.5 | 28.5 | 42.3 | 42.1 | 29.4 | 82.8 | 60.8 | 66.7 | 47.8 | 44.1 | 14.5 | 60.8 | 48.0 | 56.6 |
| SqueezeSegV3 (Xu et al. 2020) | | 55.9 | 91.7 | 74.8 | 63.4 | 26.4 | 89.0 | 92.5 | 29.6 | 38.7 | 36.5 | 33.0 | 82.0 | 58.7 | 65.4 | 45.6 | 46.2 | 20.1 | 59.4 | 49.6 | 58.9 |
| PointNet++ (Qi et al. 2017b) | | 20.1 | 72.0 | 41.8 | 18.7 | 5.6 | 62.3 | 53.7 | 0.9 | 1.9 | 0.2 | 0.2 | 46.5 | 13.8 | 30.0 | 0.9 | 1.0 | 0.0 | 16.9 | 6.0 | 8.9 |
| TangentConv (Tatarchenko et al. 2018) | | 40.9 | 83.9 | 63.9 | 33.4 | 15.4 | 83.4 | 90.8 | 15.2 | 2.7 | 16.5 | 12.1 | 79.5 | 49.3 | 58.1 | 23.0 | 28.4 | 8.1 | 49.0 | 35.8 | 28.5 |
| PointASNL (Yan et al. 2020) | 50K pts | 46.8 | 87.4 | 74.3 | 24.3 | 1.8 | 83.1 | 87.9 | 39.0 | 0.0 | 25.1 | 29.2 | 84.1 | 52.2 | **70.6** | 34.2 | 57.6 | 0.0 | 43.9 | 57.8 | 36.9 |
| RandLA-Net (Hu et al. 2020) | | 55.9 | 90.5 | 74.0 | 61.8 | 24.5 | 89.7 | 94.2 | 43.9 | 29.8 | 32.2 | 39.1 | 83.8 | 63.6 | 68.6 | 48.4 | 47.4 | 9.4 | 60.4 | 51.0 | 50.7 |
| KPConv (Thomas et al. 2019) | | 58.8 | 90.3 | 72.7 | 61.3 | 31.5 | 90.5 | 95.0 | 33.4 | 30.2 | 42.5 | 44.3 | 84.8 | 69.2 | 69.1 | 61.5 | 61.6 | 11.8 | 64.2 | 56.4 | 47.4 |
| PolarNet (Zhang et al. 2020) | | 54.3 | 90.8 | 74.4 | 61.7 | 21.7 | 90.0 | 93.8 | 22.9 | 40.3 | 30.1 | 28.5 | 84.0 | 65.5 | 67.8 | 43.2 | 40.2 | 5.6 | 61.3 | 51.8 | 57.5 |
| SparseConv (Baseline) | 50K pts | 61.8 | 89.9 | 72.1 | 56.5 | 29.6 | 90.5 | 94.5 | 43.5 | 51.0 | 42.4 | 31.3 | 83.9 | 67.4 | 68.3 | 60.4 | 61.3 | **41.1** | 65.6 | 57.9 | 67.7 |
| **JS3C-Net (Ours)** | | **66.0** | 88.9 | 72.1 | 61.9 | **31.9** | **92.5** | **95.8** | <u>54.3</u> | <u>59.3</u> | **52.9** | **46.0** | 84.5 | 69.8 | 67.9 | <u>69.5</u> | 65.4 | 39.9 | **70.8** | **60.7** | **68.7** |

et al. 2019) and SemanticPOSS (Pan et al. 2020). SemanticKITTI is currently the largest LiDAR sequential dataset with point-level annotations, which consists of 43552 densely annotated LiDAR scans belonging to 21 sequences. These scans are annotated with a total of 19 valid classes and each scan spans up to $160 \times 160 \times 20$ meters with more than $\sim 10^5$ points. We follow the official split for training, validation and online testing. SemanticPOSS is a newly proposed dataset with 11 similar annotated categories with SemanticKITTI. However, it is more challenging because each scene contains more than *twenty times* sparse small objects (i.e., people and bicycle), while the total frames number are only 1/20 of SemanticKITTI.

On both two datasets, we firstly merge consecutive 70 frames for every single frame to generate the complete volume of the entire scans. Then we select a volume of $51.2m$ in front of the LiDAR, $25.6m$ to each side of the LiDAR, and $6.4m$ in height with the resolution of $0.2m$, which results in a volume of $256 \times 256 \times 32$ voxels for prediction. Each voxel is assigned a single label based on a majority vote over all labeled points inside a voxel. Voxels containing no point are labeled as empty voxels. Note that voxels absent in all frames will not be considered in the loss calculation and evaluation.

## Joint Learning Protocol

During the end-to-end training process of JS3C-Net, we use the Adam optimizer as our optimizer. The batch size is set to 6 for the total 50 epochs. The initial learning rate is set as $0.001$ and decreases by $30\%$ after every 5 epochs. The weighted terms $\sigma_1$ and $\sigma_2$ in Eqn. 2 are randomly initialized and are trained with $\times 10$ learning rates. For semantic segmentation, $0.05m$ grids are used to conduct voxelization in SparseConv model. We randomly rotate the input point cloud along the y-axis during the training process and randomly scale it in the range of $0.9$ to $1.1$. During the inference, we apply the general voting strategy (Thomas et al. 2019; Hu et al. 2020) to average multiple prediction results of randomly augmented point clouds. Similar data augmen-

Table 2: Semantic segmentation results on the *SemanticPOSS* benchmark. The upper, medium and bottom parts contain previous projection-based, point-based and voxel-based methods, respectively.

| Method | Selected 3 classes | | | 11 classes |
|---|---|---|---|---|
| | People | Rider | Bike | avg IoU |
| SequeezeSegV2 | 18.4 | 11.2 | 32.4 | 29.8 |
| PointNet++ | 20.8 | 0.1 | 0.1 | 20.1 |
| RandLA-Net | 69.2 | 26.7 | 43.9 | 53.5 |
| KPConv | 77.3 | 29.4 | 53.2 | 55.2 |
| SparseConv | 76.3 | 30.5 | 53.5 | 57.2 |
| **JS3C-Net (Ours)** | **80.0** | **39.1** | **59.8** | **60.2** |

tation and voting strategies are also exploited for semantic scene completion. However, due to the particularity of the input volume, we use random flip along the x-axis and z-axis, and randomly rotate along y-axis by 90 degrees. All experiments are conducted on an Nvidia Tesla V100 GPU. More network details will be described in the supplementary material.

## Semantic Segmentation

In Tab. 1, we compare our JS3C-Net with recent methods on SemanticKITTI benchmark. The upper, medium and bottom parts of the table contain projection-based, point-based and voxel-based methods, respectively. The class averaged interactions over union (mIoU) is used in evaluation.

As shown in Tab. 1, our JS3C-Net surpasses all existing methods by a large margin. Merely using the SparseConv (Graham and van der Maaten 2017), training from scratch already improves upon prior arts. Yet, using our joint-learning strategy achieves markedly better segmentation results in mIoU. Specifically, JS3C-Net achieves significant improvements on small objects (e.g., motorcycle, bicycle and etc), where these objects always lose geometric details during the LiDAR collection. Thanks to the contextual shape priors from SSC, our JS3C-Net can segment
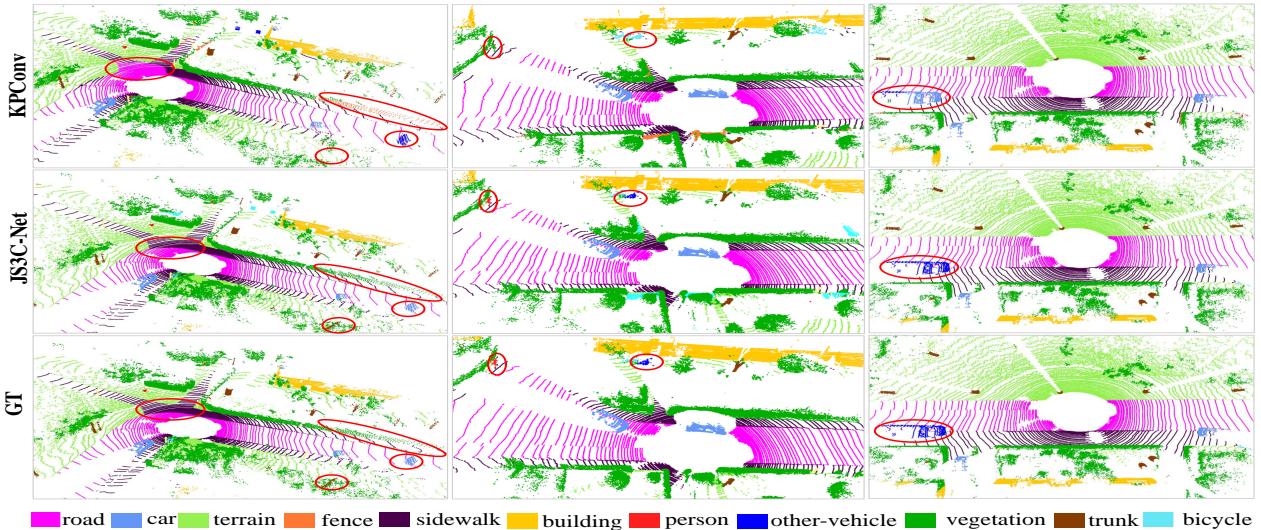
Figure 4: Qualitative results of JS3C-Net on the validation set of *SemanticKITTI* (Behley et al. 2019). Red circles show that our method performs better in many details than recent state-of-the-art KPConv (Thomas et al. 2019). Results for *SemanticPOSS* dataset are illustrated in supplementary material.

Table 3: Semantic scene completion results on the *SemanticKITTI* benchmark. Only the recent published approaches are compared.

| Method | precision | recall | IoU | mIoU |
|---|---|---|---|---|
| SSCNet | 31.7 | 83.4 | 29.8 | 9.5 |
| TS3D | 31.6 | 84.2 | 29.8 | 9.5 |
| TS3D[2] | 25.9 | **88.3** | 25.0 | 10.2 |
| EsscNet | 62.6 | 55.6 | 41.8 | 17.5 |
| TS3D[3] | **80.5** | 57.7 | 50.6 | 17.7 |
| **JS3C-Net (Ours)** | 71.5 | 73.5 | **56.6** | **23.8** |

Table 4: Ablation study on *SemanticKITTI* **validation set** for semantic segmentation (SS), semantic scene completion (SSC) and scene completion (SC).

| Model | JL | UMTL | PVI | SS | SSC | SC |
|---|---|---|---|---|---|---|
| A | | | | 63.1 | 19.4 | 51.1 |
| B | ✓ | | | 66.1 | 22.6 | 55.0 |
| C | ✓ | ✓ | | 66.4 | 23.0 | 56.1 |
| D | ✓ | ✓ | ✓ | **67.5** | **24.0** | **57.0** |

them well. Meanwhile, Fig. 4 presents some visualization results of JS3C-Net on the validation split, which demonstrates great improvements on small objects, in particular.

Tab. 2 illustrates the semantic segmentation results on SemanticPOSS dataset, where we compare result state-of-the-art methods. These results show that our proposed JS3C-Net can achieve larger improvement compared with our baseline on more challenging data with remarkable small objects.

**Semantic Scene Completion**

With the semantic guidance from segmentation network, our JS3C-Net can also make significant breakthrough in semantic scene completion (SSC) task. Tab. 3 illustrates the results of SSC on SemanticKITTI benchmark, where we compare our JS3C-Net with recent state-of-the-art methods. All methods are implemented with same settings and the detailed implementation and concrete results will be further elaborated in supplementary.

Since semantic scene completion requires to simultaneously predict the occupation status and the semantic label of a voxel, we follow the evaluation protocol of (Song et al. 2017) to compute the precision, recall and IoU for

the task of scene completion (SC) while ignoring the semantic label. Meanwhile, the mIoU over the 19 classes is also exploited for the evaluation of semantic scene completion (SSC). As shown in Tab. 3, our JS3C-Net achieves state-of-the-art results on both SC and SSC tasks. Benefit from semantic guidance, joint learning strategy and well-designed interaction module, our SSC module can generate more faithful geometric details. The results of our methods are 6% higher than previous state-of-the-art TS3D[3] (Garbade et al. 2019b; Behley et al. 2019; Liu et al. 2018) for scene completion, which uses segmentation results from DarkNet53Seg (Behley et al. 2019) as inputs as well. Fig. 5 shows visualization results of semantic scene completion.

**Design Analysis**

**Ablation Study.** The ablation results are summarized in Tab. 4. Since the limit of submission times, all ablated networks are tested on the validation set of SemanticKITTI dataset (Behley et al. 2019).

The baseline (model A) is set to learn without joint learning, i.e., train two tasks separately. The baseline only gets IoU of 63.1% on semantic segmentation (SS) and 51.1% on scene completion (SC). This convincingly confirms the effectiveness of joint learning (JL) in model B (we manually
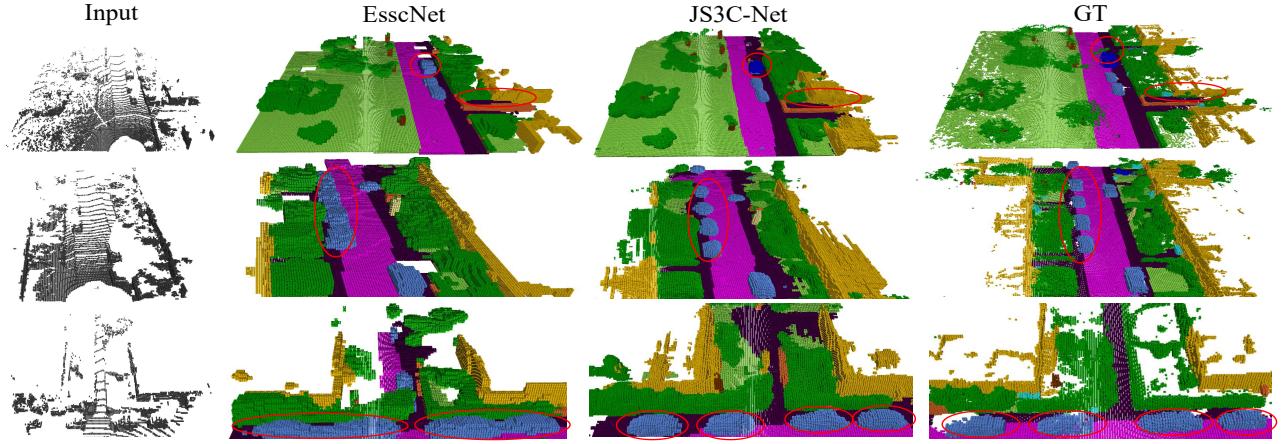
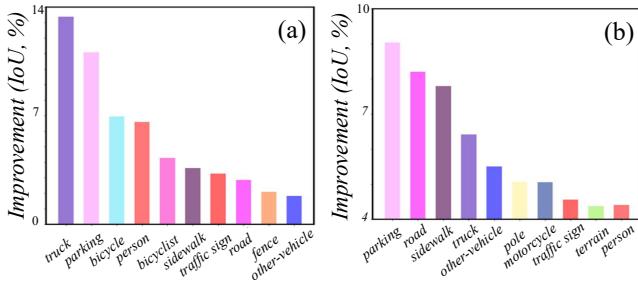Figure 5: Qualitative results of SSC task on the validation set of *SemanticKITTI* (Behley et al. 2019).



Figure 6: Top-10 mIoU gains between JS3C-Net and split-trained single task (SS or SSC) on the **validation set** of *SemanticKITTI* (Behley et al. 2019), where (a) and (b) illustrate SS and SSC respectively.

Table 5: **Complexity Analysis.** Model size and latency for different methods. Here underline correspond to the post-processing time. Parameters and time in ($\cdot$) represents the extra operation of SSC and PVI modules, which can be ignored during inference of SS.

|  | Parameters (million) | Latency (ms) |
|---|---|---|
| PointNet++ | 6.0 | 5900 |
| TangentConv | 0.4 | 3000 |
| RandLA-Net | 1.2 | 256+624 |
| KPConv | 18.3 | 1117+624 |
| SparseConv | 2.7 | 471 |
| **JS3C-Net (Ours)** | 2.7(+0.4) | **471**(+107) |

set weights of task losses as 1:0.8), which is significantly improved to 66.1% for SS and 55.0% for SC by a large margin. Correspondingly, in model C, uncertainty multi-task loss (UMTL) is used to achieve the optimal trade-off between two tasks automatically. As a result, improvements of 0.4% and 1.1% on SS and SC are further obtained through UMTL (model C). Then, with PVI module for knowledge fusion, our JS3C-Net achieves the best results on both tasks.

**Mutual Promotion.** To further study the reciprocal effects between two tasks, we conducted comparative experiments between the single task (SS or SSC) and the multiple tasks (JS3C-Net). As shown in Fig. 6, when the JS3C-Net is introduced (i.e., SS and SSC learn jointly), the performances of the two tasks are both largely enhanced, where we show the top 10 mIoU gains in Fig. 6. It shows that the IoUs of all 19 classes have been improved, especially in the case of small objects, such as trucks, bicycles and persons. The plausible explanation is that, for small objects, their raw point clouds are very sparse and usually lose local details. When SS and SSC learn jointly, SS can take advantage of the contextual shape prior of small objects from SSC, which benefits SS to classify each point more precisely. Therefore, these two tasks can benefit one from the other mutually.

**Complexity Analysis.** In this section, we evaluate the overall complexity of JS3C-Net. As shown in Tab. 5, our proposed JS3C-Net is much more light-weighted and faster than previous point-based methods (**1/5** model size and **1/3** inference time of KPConv). More importantly, since the disposable properties of our SSC module and PVI module, JS3C-Net has the same speed with segmentation backbone, which makes it more suitable for the real-time applications.

## Conclusion

In this work, we propose an single sweep LiDAR point cloud semantic segmentation framework via contextual shape priors from semantic scene completion network, named JS3C-Net. By exploiting some sophisticated pipelines, interactive modules, and reasonable loss function, our JS3C-Net model achieves state-of-the-art results on both semantic segmentation and scene completion tasks, outperforming previous methods by a large margin. We believe that our work can be applied to a wider range of other scenarios in the future, such as indoor point cloud sequence. Meanwhile, our method provides an alternative solution to the comprehension of large-scale LiDAR scenes with severe local details missing. It can improve the performance through contextual shape priors learning and interactive knowledge transferring.

## Acknowledgments

## References

Alonso, I.; Riazuelo, L.; Montesano, L.; and Murillo, A. C. 2020. 3D-MiniNet: Learning a 2D Representation from Point Clouds for Fast and Efficient 3D LIDAR Semantic Segmentation. *arXiv preprint arXiv:2002.10893* .

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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1534–1543.

Behley, J.; Garbade, M.; Milioto, A.; Quenzel, J.; Behnke, S.; Stachniss, C.; and Gall, J. 2019. SemanticKITTI: A dataset for semantic scene understanding of lidar sequences. In *Proceedings of the IEEE International Conference on Computer Vision*, 9297–9307.

Boulch, A.; Le Saux, B.; and Audebert, N. 2017. Unstructured Point Cloud Semantic Labeling Using Deep Segmentation Networks. *3DOR* 2: 7.

Choy, C.; Gwak, J.; and Savarese, S. 2019. 4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3075–3084.

Dai, A.; Chang, A. X.; Savva, M.; Halber, M.; Funkhouser, T.; and Nießner, M. 2017. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5828–5839.

Dai, A.; Ritchie, D.; Bokeloh, M.; Reed, S.; Sturm, J.; and Niener, M. 2018. ScanComplete: Large-Scale Scene Completion and Semantic Segmentation for 3D Scans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4578¨C4587.

Fan, H.; and Yang, Y. 2019. PointRNN: Point recurrent neural network for moving point cloud processing. *arXiv preprint arXiv:1910.08287* .

Garbade, M.; Chen, Y.-T.; Sawatzky, J.; and Gall, J. 2019a. Two Stream 3D Semantic Scene Completion. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.

Garbade, M.; Chen, Y.-T.; Sawatzky, J.; and Gall, J. 2019b. Two stream 3d semantic scene completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 0–0.

Graham, B.; Engelcke, M.; and van der Maaten, L. 2018. 3d semantic segmentation with submanifold sparse convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9224–9232.

Graham, B.; and van der Maaten, L. 2017. Submanifold sparse convolutional networks. *arXiv preprint arXiv:1706.01307* .

Hackel, T.; Savinov, N.; Ladicky, L.; Wegner, J. D.; Schindler, K.; and Pollefeys, M. 2017. Semantic3d. net: A new large-scale point cloud classification benchmark. *arXiv preprint arXiv:1704.03847* .

Han, L.; Zheng, T.; Xu, L.; and Fang, L. 2020. OccuSeg: Occupancy-aware 3D Instance Segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

Han, X.; Li, Z.; Huang, H.; Kalogerakis, E.; and Yu, Y. 2017. High-Resolution Shape Completion Using Deep Neural Networks for Global Structure and Local Geometry Inference. In *ECCV*, 85¨C93,.

Hu, Q.; Yang, B.; Xie, L.; Rosa, S.; Guo, Y.; Wang, Z.; Trigoni, N.; and Markham, A. 2020. RandLA-Net: Efficient semantic segmentation of large-scale point clouds. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* .

Kendall, A.; Gal, Y.; and Cipolla, R. 2018. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 7482–7491.

Landrieu, L.; and Boussaha, M. 2019. Point cloud oversegmentation with graph-structured deep metric learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 7440–7449.

Landrieu, L.; and Simonovsky, M. 2018. Large-scale point cloud semantic segmentation with superpoint graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4558–4567.

Lawin, F. J.; Danelljan, M.; Tosteberg, P.; Bhat, G.; Khan, F. S.; and Felsberg, M. 2017. Deep projective 3D semantic segmentation. In *International Conference on Computer Analysis of Images and Patterns*, 95–107. Springer.

Li, J.; Liu, Y.; Gong, D.; Shi, Q.; Yuan, X.; Zhao, C.; and Reid, I. 2019. RGBD based dimensional decomposition residual network for 3D semantic scene completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 7693–7702.

Liu, S.; Hu, Y.; Zeng, Y.; Tang, Q.; Jin, B.; Han, Y.; and Li, X. 2018. See and think: Disentangling semantic scene completion. In *Advances in Neural Information Processing Systems*, 263–274.

Liu, X.; Yan, M.; and Bohg, J. 2019. MeteorNet: Deep learning on dynamic 3D point cloud sequences. In *Proceedings of the IEEE International Conference on Computer Vision*, 9246–9255.

Liu, Y.; Li, J.; Yan, Q.; Yuan, X.; Zhao, C.; Reid, I.; and Cadena, C. 2020. 3D Gated Recurrent Fusion for Semantic Scene Completion. *arXiv preprint arXiv:2002.07269* .

Milioto, A.; Vizzo, I.; Behley, J.; and Stachniss, C. 2019. Rangenet++: Fast and accurate lidar semantic segmentation. In *Proc. of the IEEE/RSJ Intl. Conf. on Intelligent Robots and Systems (IROS)*.

Pan, Y.; Gao, B.; Mei, J.; Geng, S.; Li, C.; and Zhao, H. 2020. SemanticPOSS: A Point Cloud Dataset with Large Quantity of Dynamic Instances. *arXiv preprint arXiv:2002.09147* .

Pham, Q.-H.; Nguyen, T.; Hua, B.-S.; Roig, G.; and Yeung, S.-K. 2019. JSIS3D: joint semantic-instance segmentation of 3d point clouds with multi-task pointwise networks and multi-value conditional random fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 8827–8836.

Qi, C. R.; Su, H.; Mo, K.; and Guibas, L. J. 2017a. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 652–660.

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 *Advances in neural information processing systems*, 5099–5108.

Riegler, G.; Osman Ulusoy, A.; and Geiger, A. 2017. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3577–3586.

Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, 234–241. Springer.

Shi, H.; Lin, G.; Wang, H.; Hung, T.-Y.; and Wang, Z. 2020. SpSequenceNet: Semantic Segmentation Network on 4D Point Clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4574–4583.

Song, S.; Yu, F.; Zeng, A.; Chang, A. X.; Savva, M.; and Funkhouser, T. 2017. Semantic scene completion from a single depth image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1746–1754.

Tang, H.; Liu, Z.; Zhao, S.; Lin, Y.; Lin, J.; Wang, H.; and Han, S. 2020. Searching Efficient 3D Architectures with Sparse Point-Voxel Convolution. In *European Conference on Computer Vision*.

Tatarchenko, M.; Park, J.; Koltun, V.; and Zhou, Q.-Y. 2018. Tangent convolutions for dense prediction in 3d. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3887–3896.

Thomas, H.; Qi, C. R.; Deschaud, J.-E.; Marcotegui, B.; Goulette, F.; and Guibas, L. J. 2019. KPConv: Flexible and Deformable Convolution for Point Clouds. In *The IEEE International Conference on Computer Vision (ICCV)*.

Wang, L.; Huang, Y.; Hou, Y.; Zhang, S.; and Shan, J. 2019a. Graph attention convolution for point cloud semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 10296–10305.

Wang, P.-S.; Liu, Y.; Guo, Y.-X.; Sun, C.-Y.; and Tong, X. 2017. O-cnn: Octree-based convolutional neural networks for 3d shape analysis. *ACM Transactions on Graphics (TOG)* 36(4): 72.

Wang, X.; Liu, S.; Shen, X.; Shen, C.; and Jia, J. 2019b. Associatively segmenting instances and semantics in point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4096–4105.

Wang, Y.; Sun, Y.; Liu, Z.; Sarma, S. E.; Bronstein, M. M.; and Solomon, J. M. 2019c. Dynamic graph cnn for learning on point clouds.

Wei, J.; Lin, G.; Yap, K.-H.; Hung, T.-Y.; and Xie, L. 2020. Multi-Path Region Mining For Weakly Supervised 3D Semantic Segmentation on Point Clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

Wu, B.; Wan, A.; Yue, X.; and Keutzer, K. 2018. Squeezeseg: Convolutional neural nets with recurrent crf for real-time road-object segmentation from 3d lidar point cloud. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 1887–1893. IEEE.

Wu, B.; Zhou, X.; Zhao, S.; Yue, X.; and Keutzer, K. 2019. Squeezesegv2: Improved model structure and unsupervised domain adaptation for road-object segmentation from a lidar point cloud. In *2019 International Conference on Robotics and Automation (ICRA)*, 4376–4382. IEEE.

Wu, W.; Qi, Z.; and Fuxin, L. 2019. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 9621–9630.

Xu, C.; Wu, B.; Wang, Z.; Zhan, W.; Vajda, P.; Keutzer, K.; and Tomizuka, M. 2020. Squeezesegv3: Spatially-adaptive convolution for efficient point-cloud segmentation. *arXiv preprint arXiv:2004.01803* .

Yan, X.; Zheng, C.; Li, Z.; Wang, S.; and Cui, S. 2020. PointASNL: Robust Point Clouds Processing using Nonlocal Neural Networks with Adaptive Sampling.

Zhang, J.; Zhao, H.; Yao, A.; Chen, Y.; and Zhang, L. 2018. Efficient Semantic Scene Completion Network with Spatial Group Convolution. In *ECCV*, 733¨C749.

Zhang, Y.; Zhou, Z.; David, P.; Yue, X.; Xi, Z.; Gong, B.; and Foroosh, H. 2020. PolarNet: An Improved Grid Representation for Online LiDAR Point Clouds Semantic Segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9601–9610.

# Supplementary Material

In this supplementary material, we first show our experimental details. Besides, we provide additional discussion to further demonstrate the superiority of our model.

## Concrete Experimental Design

### Network Architectures

Our model, JS3C-Net, consists of three parts: a semantic segmentation network, a semantic scene completion (SSC) decoder, and a point-voxel interaction (PVI) module.

For semantic segmentation stage, we use Submanifold SparseConv (Graham, Engelcke, and van der Maaten 2018) as our backbone due to its superior performance. The encoder consists of 7 blocks of 3D sparse convolutions, each of which has two 3D sparse convolutions inside. Going deeper, we gradually increase the number of channels (i.e., $16, 32, 48, 64, 80, 96, 112$). We also apply a 3D sparse pooling operation after each block to reduce the spatial resolution of the feature maps. For the decoder, we use the same structure but in the reverse order and replace the 3D sparse pooling layers with unpooling operations. Furthermore, it concatenates the features from the encoder phase on the decoder features at each scale through skip-connection. Three MLPs are used to generate shape embedding for PVI module and C-category prediction.

The architecture of SSC decoder is already shown in the manuscript. In fact, we use five basic blocks and the channel dimension of each block is identically 32. As for the point-voxel interaction (PVI) module, graph-based edge learning is implemented by a three-layer MLP with 32 channels. In experiment, we set number of GCN layer $l = 1$.

### Semantic Segmentation

In this section, we illustrate in details the concrete results of semantic segmentation and implementation of compared methods. For SemanticKITTI benchmark, we directly compare our results with results on official benchmark. For SemanticPOSS dataset (see Tab.3, we use results of PointNet++ and SequeezeSegV2 in their official paper, and results of RandLA-Net and KPConv are implemented by official codes of corresponding methods.

### Semantic Scene Completion

In this section, we illustrate in details the concrete results of semantic scene completion and implementation of compared methods, which are already provided by (Behley et al. 2019). Tab. 5 depicts the results of semantic scene completion on benchmark (Behley et al. 2019). Among the methods to be compared, SSCNet (Song et al. 2017), EsscNet (Zhang et al. 2018), TS3D (Garbade et al. 2019b) directly use sparse point cloud as input, TS3D[2] and TS3D[3] are multi-stage methods, which use segmentation results of DarkNet53Seg (Behley et al. 2019) as semantic priors.

**SSCNet.** SSCNet (Song et al. 2017) is a weak baseline that directly uses processed volumes as input, then several 3D convolution layers with the same structure introduced in their paper are conducted.

Table 1: The results of different joint learning strategies on *SemanticKITTI* validation set. SS, SC and SSC mean semantic segmentation, scene completion and semantic scene completion, respectively. JL means learn two tasks jointly.

| Model | SS | SC | SSC | JL | SS | SSC | SC |
|-------|----|----|----|----|------|------|------|
| A | ✓ | | | | 63.1 | 19.4 | 51.1 |
| B | ✓ | ✓ | | ✓ | 62.5 | - | 55.7 |
| C | ✓ | | ✓ | ✓ | **67.5** | **24.0** | **57.0** |

Table 2: Effect of moving objects on *SemanticKITTI* and *SemanticPOSS* validation set for semantic segmentation, where w and w/o mean 'with' and 'without'.

| | SemanticKITTI | SemanticPOSS |
|---|---|---|
| w moving objects | 67.5 | 60.2 |
| w/o moving objects | 67.7 | 60.2 |

**TS3D.** Two Stream (TS3D) approach (Garbade et al. 2019b) makes use of the additional information (processed from pre-trained DeepLab_v2) from the RGB image corresponding to the input laser scan, and then it performs a 4 fold downsampling in a forward process but it renders them incapable of dealing with details of the scene.

**TS3D[2].** It uses the semantic segmentation results of DarkNet53Seg (Behley et al. 2019) to enhance the semantic scene completion. However, without a suitable joint learning strategy, the pre-processed feature cannot improve the result of semantic scene completion effectively and significantly.

**TS3D[3].** It replaces the backbone of the TS3D[2] with SATNet (Liu et al. 2018) without downsampling, and divide each whole volume into six equal parts for more fine-grained prediction.

**EsscNet.** We also compare our method with EsscNet (Zhang et al. 2018), which uses spatial group convolution (SGC) for memory saving. It uses the entire scene as input. However, due to the lack of priors, it still cannot perform very well. Also, the large amount of parameters shown in the manuscript makes it unable to meet the requirement for real-time applications.

## Additional Discussion

### Why Use Semantic Scene Completion?

In this section, we further discuss why our model should use semantic scene completion as the object of joint learning. Note that semantic segmentation (SS) and semantic scene completion (SSC) both include semantic labeling process, which seems to be redundant in these two tasks. Therefore, we conduct an experiment that uses scene completion (SC) instead of semantic scene completion (SSC) for the joint learning with semantic segmentation (SS).

As shown in Tab. 1, when using scene completion and semantic segmentation for joint learning (model B), although it can somehow improve the result of scene completion, the segmentation results become worse than directly conducting segmentation network (model A). This proves that directly

Table 3: Semantic segmentation results on the *SemanticPOSS* dataset.

| Method | people | rider | car | traffic sign | trunk | plants | pole | fence | building | bike | road | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PointNet++ (Qi et al. 2017a,b) | 20.8 | 0.1 | 8.9 | 21.8 | 4.0 | 51.2 | 3.2 | 6.0 | 42.7 | 0.1 | 62.2 | 20.1 |
| SequeezeSegV2 (Wu et al. 2019) | 18.4 | 11.2 | 34.9 | 11.0 | 15.8 | 56.3 | 4.5 | 25.5 | 47.0 | 32.4 | 71.3 | 29.8 |
| RandLA-Net (Hu et al. 2020) | 69.2 | 26.7 | 77.2 | 24.3 | 27.3 | 73.4 | 30.9 | 51.4 | 82.5 | 43.9 | 81.2 | 53.5 |
| KPConv (Thomas et al. 2019) | 77.3 | 29.4 | 78.2 | 23.0 | 28.1 | 71.6 | 32.4 | 51.5 | 81.8 | 53.2 | 80.6 | 55.2 |
| SparseConv (Baseline) | 76.3 | 30.5 | 80.8 | 28.3 | 29.1 | 74.9 | 39.8 | 51.8 | 83.7 | 53.5 | 80.5 | 57.2 |
| **JS3C-Net(Ours)** | **80.0** | **39.1** | **83.2** | **28.9** | **31.6** | **76.1** | **44.7** | **54.0** | **83.9** | **59.8** | **81.7** | **60.2** |

adding scene completion will introduce more bias into segmentation network.

The explanation for such results can be summarized as follows: (1) Due to the cascaded architecture of our model, the results of SS and SSC actually have spatial alignment, which allows our SSC decoder to focus more on completing the shapes of each category rather than the entire scene. This setting allows the process of back propagation to bring shape prior to each category of semantic segmentation with less noise. (2) PVI module uses both spatial coordinates and features information to construct a learnable weight of graph edges, where features from SS and SSC are mainly used to compare the differences in feature space. Therefore, it can further enhance the distinction of features in segmentation network through using 20 categories to predict the SSC with strong semantic priors instead of simply predicting whether a certain voxel has an object or not.

### Do Moving Objects Affect the Results?

As shown in the red circle of Fig. 1 (left), when using multiple frames to generate an entire dense scene volume, moving objects will be inevitably reconstructed into a long bar, which may introduce bias to the semantic segmentation. Therefore, to further test whether these moving objects will affect segmentation results, we design an experiment as follows. Firstly, we normally reconstruct moving objects, as shown in Fig. 1 (left). Then, as shown in Fig. 1 (right), we use the annotations of moving objects in both two datasets, and only use 5 frames to reconstruct each moving object (instead of 70 frames).

The quantitative results are shown in Tab. 2. It can be seen that, when we remove moving objects in semantic scene completion, there is only 0.2% improvement in semantic segmentation. This is because moving objects only account for a small proportion in the datasets. Finally, since there are also moving objects in SSC benchmark of (Behley et al. 2019), we keep moving objects in our training data in order to make a fair comparison with other methods.

### Can We Only Consider Historical Frames?

To further illustrate the advantages of using semantic scene completion, we also compare the results of only consider historical frames. Here we consider the method proposed in SpSequenceNet (Shi et al. 2020) to use attention mechanism and feature aggregation from historical frames.

Tab. 4 shows the results of the correlative ablation study. Since results of SpSequenceNet on SemanticKITTI is much lower than ours, we re-implement their framework with our

Table 4: The results on *SemanticKITTI* validation set.

| Ablation | Latency (ms) | mIoU |
|---|---|---|
| SparseConv (baseline) | 471 | 63.1 |
| SparseConv + SpSequenceNet | 972 | 64.3 |
| JS3C-Net (historical) | 471 | 66.6 |
| JS3C-Net (historical+future) | 471 | **67.5** |

baseline. Experiment results show that only using attention or feature aggregation with previous frames cannot better learn the complete shape priors from LiDAR sequences. Our results of only using historical frames to generate ground truth is still higher than their results. Note that only using historical frames will ignore the shape details of newly incoming objects in each frame, thus hampers the results. When we consider future frames of LiDAR sequence, our results will be further improved.

Furthermore, their method inevitably introduces more computational burden to segmentation, while our auxiliary components are fully discardable (see latency in table).

## More Visualization

In Fig. 2, we show more cases on SemanticPOSS dataset.

Table 5: Semantic scene completion results on the *SemanticKITTI* benchmark. Only the recent published approaches are compared.

| | Scene Completion | | | Semantic Scene Completion | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Method | precision | recall | IoU | road | sidewalk | parking | other-ground | building | car | truck | bicycle | motorcycle | other-vehicle | vegetation | trunk | terrain | person | bicyclist | motorcyclist | fence | pole | traffic sign | mIoU |
| SSCNet (Song et al. 2017) | 31.7 | 83.4 | 29.8 | 27.6 | 17.0 | 15.6 | 6.0 | 20.9 | 10.4 | 1.8 | 0.0 | 0.0 | 0.1 | 25.8 | 11.9 | 18.2 | 0.0 | 0.0 | 0.0 | 14.4 | 7.9 | 3.7 | 9.5 |
| TS3D (Garbade et al. 2019b) | 31.6 | 84.2 | 29.8 | 28.0 | 17.0 | 15.7 | 4.9 | 23.2 | 10.7 | 2.4 | 0.0 | 0.0 | 0.2 | 24.7 | 12.5 | 18.3 | 0.0 | 0.1 | 0.0 | 13.2 | 7.0 | 3.5 | 9.5 |
| TS3D[2] (Garbade et al. 2019b; Behley et al. 2019) | 25.9 | **88.3** | 25.0 | 27.5 | 18.5 | 18.9 | 6.6 | 22.1 | 8.0 | 2.2 | 0.1 | 0.0 | 4.0 | 19.5 | 12.9 | 20.2 | 2.3 | 0.6 | 0.0 | 15.8 | 7.6 | 6.7 | 10.2 |
| EsscNet (Zhang et al. 2018) | 62.6 | 55.6 | 41.8 | 43.8 | 28.1 | 26.9 | 10.3 | 29.8 | 26.4 | 5.0 | 0.3 | 5.4 | 9.1 | 35.8 | **20.1** | 28.7 | 2.9 | 2.7 | 0.1 | 23.3 | 16.4 | **16.7** | 17.5 |
| TS3D[3] (Garbade et al. 2019b; Behley et al. 2019; Liu et al. 2018) | **80.5** | 57.7 | 50.6 | 62.2 | 31.6 | 23.3 | 6.5 | 34.1 | 30.7 | 4.9 | 0.0 | 0.0 | 0.1 | 40.1 | 21.9 | 33.1 | 0.0 | 0.0 | 0.0 | 24.1 | 16.9 | 6.9 | 17.7 |
| **JS3C-Net (Ours)** | 70.2 | 74.5 | **56.6** | **64.7** | **39.9** | **34.9** | **14.1** | **39.4** | **33.3** | **7.2** | **14.4** | **8.8** | **12.7** | **43.1** | 19.6 | **40.5** | **8.0** | **5.1** | **0.4** | **30.4** | **18.9** | 15.9 | **23.8** |



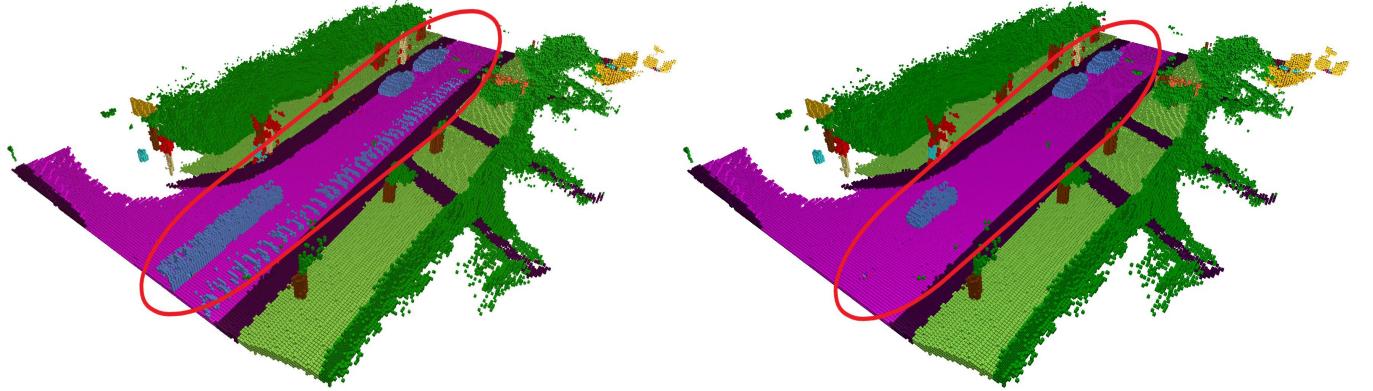Figure 1: The selected example of ground truths for SSC, where (left) and (right) illustrate reconstructed results with and without moving objects.



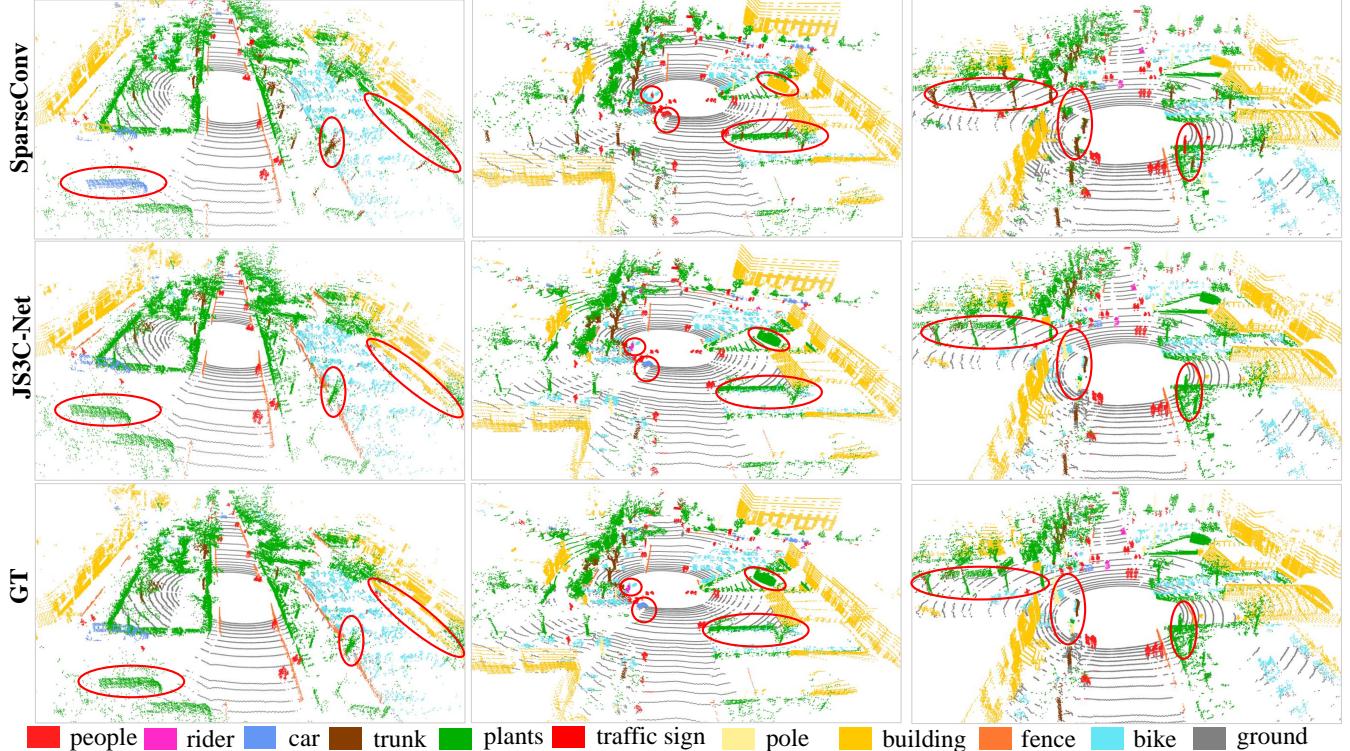| people | rider | car | trunk | plants | traffic sign | pole | building | fence | bike | ground |

Figure 2: The visualization results on *SemanticPOSS* dataset.