

# MM-3DScene: 3D Scene Understanding by Customizing Masked Modeling with Informative-Preserved Reconstruction and Self-Distilled Consistency

Mingye Xu<sup>1,3,4,\*</sup>, Mutian Xu<sup>2,5,\*</sup>, Tong He<sup>4</sup>, Wanli Ouyang<sup>4</sup>, Yali Wang<sup>1,4,†</sup>, Xiaoguang Han<sup>2,5</sup>, Yu Qiao<sup>1,4,†</sup>

<sup>1</sup> The Guangdong Provincial Key Laboratory of Computer Vision and Virtual Reality Technology, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, Shenzhen 518055, China

<sup>2</sup> SSE, CUHKSZ

<sup>3</sup> University of Chinese Academy of Sciences

<sup>4</sup> Shanghai Artificial Intelligence Laboratory

<sup>5</sup> FNii, CUHKSZ

[mingyexu.github.io/mm3dscene](https://mingyexu.github.io/mm3dscene)

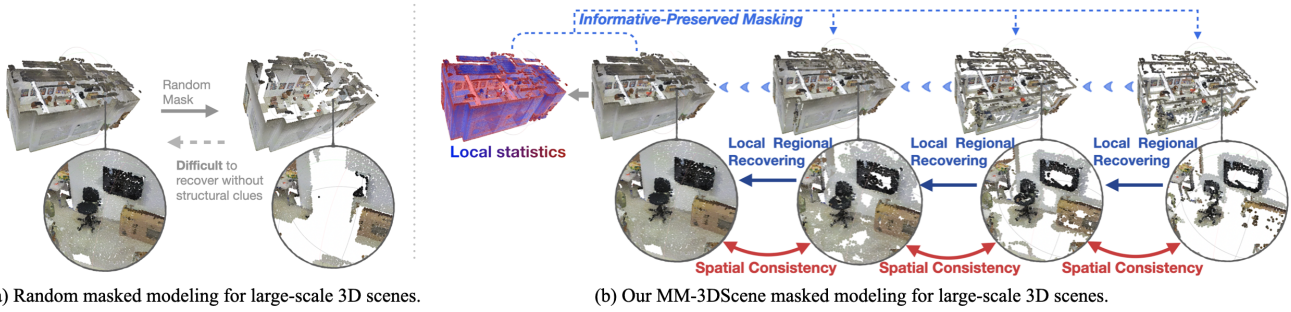


Figure 1. **How to apply masked modeling for large-scale 3D scenes?** (a) Conventional random masked modeling on 3D scenes may cause a high risk of uncertainty. In this figure, a *chair* and a *TV* are totally masked, which are extremely difficult to be recovered without any context guidance. (b) Our **MM-3DScene** exploits local statistics to discover and *preserve* representative structured points, effectively simplifying the pretext task. At each learning step, our method focuses on restoring *regional* geometry, and enjoys less ambiguity. Moreover, since unmasked areas are underexplored during reconstruction, the model is encouraged to maintain the intrinsic *spatial consistency* on unmasked points between different masking ratios, which requires the consistent understanding of unmasked areas.

## Abstract

Masked Modeling (MM) has demonstrated widespread success in various vision challenges, by reconstructing masked visual patches. Yet, applying MM for large-scale 3D scenes remains an open problem due to the data sparsity and scene complexity. The conventional random masking paradigm used in 2D images often causes a high risk of ambiguity when recovering the masked region of 3D scenes. To this end, we propose a novel informative-preserved reconstruction, which explores local statistics to discover and preserve the representative structured points, effectively enhancing the pretext masking task for 3D scene understanding. Integrated with a progressive reconstruction manner, our method can concentrate on modeling regional geometry and enjoy less ambiguity for masked reconstruction. Besides, such scenes with progressive masking ra-

tios can also serve to self-distill their intrinsic spatial consistency, requiring to learn the consistent representations from unmasked areas. By elegantly combining informative-preserved reconstruction on **masked** areas and consistency self-distillation from **unmasked** areas, a unified framework called **MM-3DScene** is yielded. We conduct comprehensive experiments on a host of downstream tasks. The consistent improvement (e.g., +6.1% mAP@0.5 on object detection and +2.2% mIoU on semantic segmentation) demonstrates the superiority of our approach.

## 1. Introduction

3D scene understanding plays an essential role in various visual applications, such as virtual reality, robot navigation, and autonomous driving. Over the past few years, deep learning has dominated 3D scene parsing tasks [27, 43, 80]. However, traditional supervised learning methods require massive annotation of 3D scene data that are extremely la-

\* Equal contribution.

† Corresponding authors.

borious to obtain [10], where millions of points or mesh vertices per scene need to be labeled.

To solve this, self-supervised learning (SSL) becomes a favorable choice since it can extract rich representations without any annotation [11, 17, 18]. Masked Modeling (MM) [17, 62], as one of the representative methods in SSL, recently draws significant attention in the vision community. Recently, It has been explored in 3D vision [34, 42, 56, 73, 76, 77], where these 3D MM approaches randomly mask local regions of point clouds, and pre-train neural networks to reconstruct the masked areas. Nevertheless, such random masking paradigms are not feasible for large-scale 3D scenes, which often causes a high risk of reconstruction ambiguity. As illustrated in Fig. 1 (a), a chair and a TV are totally masked, which are extremely difficult to be recovered without any context guidance. Such ambiguity often makes MM difficult to learn informative representation for 3D scenes. Hence, we ask a natural question: *can we customize a better way of masked modeling for 3D scene understanding?*

To tackle this question, we propose a novel informative-preserved masked reconstruction scheme in this paper. Specifically, we leverage local statistics of each point (*i.e.*, the difference between each point and its neighboring points in terms of color and shape) as guidance to discover the representative structured points which are usually located at the boundary regions in the 3D scene. We denote these points as ‘**Informative Points**’ since they provide highly useful information hints and rich semantic context for *assisting* masked reconstruction. To this end, our mask strategy is definite: to *preserve* Informative Points in a scene and mask other points. In this way, the basic geometric information of a scene is explicitly *retained*, which effectively simplifies the pretext task and reduces ambiguity.

Based on our mask strategy, a progressive masked reconstruction manner is integrated, to better model the *masked* areas. As illustrated in Fig. 1 (b), during each iteration, our method *concentrates* on reconstructing the local *regional* geometric patterns rather than rebuilding the original intact scene. In doing so, it enjoys less ambiguity and is able to restore accurate geometric information.

Moreover, we realize the information of *unmasked* areas (*i.e.*, Informative Points) is underexplored. We find that there exists point correspondence in the unmasked areas under progressive masking ratios. Accordingly, we introduce a dual-branch encoding scheme for learning such intrinsic consistency, with the ultimate goal of unearthing the consistent (*i.e.*, masking-invariant) representations from unmasked areas. This leads to a more powerful SSL framework on 3D scenes, called MM-3DScene, which elegantly combines the masked modeling on the masked and unmasked areas in 3D scenes together, while *complements* each other. It achieves superior performance in Table 7 (v).

Datasets	Complexity	Task	Gain (from scratch)
S3DIS	Entire floor, office	segmentation	(+1.5%) mIoU
ScanNet	Large rooms	segmentation	(+2.2%) mIoU
		detection	(+4.4) mAP@0.25
SUN-RGBD	Cluttered rooms	detection	(+2.9) mAP@0.25
		detection	(+4.4) mAP@0.5

Table 1. **Summary of fine-tuning MM-3DScene** on various downstream tasks and datasets for 3D understanding. Our MM-3DScene conspicuously boosts the performance of the baseline model trained from scratch.

Our contributions are motivated and comprehensive:

- We raise the concept of Informative Points – the points providing significant information hints, and indicate that preserving them is critical for assisting masked modeling on 3D scenes (Table 8).
- For masked areas, we propose an informative-preserved reconstruction scheme to focus on restoring the regional geometry in a novel progressive manner, which explicitly simplifies the pretext task.
- For unmasked areas, we introduce a self-distillation branch, which is encouraged to learn spatial-consistent representations under progressive masking ratios.
- A unified self-supervised framework, called MM-3DScene, delivers performance improvements across on various downstream tasks and datasets (Table 1).

## 2. Related Work

**3D Scene Understanding** With the rapid development of deep learning methods in point cloud analysis and the emergence of large-scale 3D datasets [1, 5, 10, 16, 63], the research focus gradually migrate from synthetic, single object analysis [35, 38, 60, 65, 66, 68] to complex large-scale scene understanding, especially scene segmentation [14, 19–21, 67, 79], 3D object detection [36, 40, 43]. PointNet [45] and its variants [32, 46, 57, 58] extract local features from neighbors through hierarchical grouping architecture to capture fine-grained representation. Thomas *et al.* [54] defines deformable kernel point convolution to capture the point cloud representation using a set of learnable kernel points. In PAConv [64], a continuous convolutional kernel is built by dynamically combining several weight banks, where the coefficients are learned from point positions. Zhao *et al.* [80] proposes to assign learnable attentional weights to local point features, and introduces a Transformer [55]-like architecture into point cloud analysis. Rather than focusing on developing deep architectural details, in this paper, we explore an effective self-supervised pre-training mechanism for scene understanding based on the basic version of Point Transformer.

**Unsupervised Scene Pre-training** Self-supervised learning (SSL) has recently achieved great success in 2D vision [3, 6, 7, 15, 17, 18, 62] and NLP tasks [4, 11, 12]. But there has been limited exploration of SSL for 3D vision. Most of the existing 3D-SSL methods aim to understand 3D point clouds, and can be divided into two types, masked modeling (MM)-based methods, and contrastive learning-based methods, respectively.

MM-based pre-training usually takes as input point clouds to recover itself as the pretext task. OcCo [56] first constructs an occlusion point cloud and applies an encoder-decoder mechanism to reconstruct the original object. Yu *et al.* [73] proposes the idea of restoring the masked proxies under the supervision of the pre-trained tokenizer. Borrowing the idea from MAE [17], Pang *et al.* [42] designs a masked auto-encoder to recover the masked parts of objects. However, most of these MM-based methods are focused on 3D shape-level pre-training, how to apply such a scheme on 3D scenes has not been fully investigated. Although some concurrent works [22, 39] voxelize the outdoor automotive point clouds and randomly mask voxels, they are only manifested to especially benefit the LiDAR-based object detection task. In this paper, we creatively apply the self-supervised pre-training directly on 3D scenes by a scene-specific masked reconstruction.

On the other hand, using the idea of contrastive scene contexts for pre-training has been shown to be effective [8, 25, 61]. Xie *et al.* [61] and Hou *et al.* [25] perform scene pre-training using contrastive learning of point features between a pair of overlapping scans. 4DContrast [8] first composites the synthetic 3D object with real-world 3D scans to create 4D sequential data, and utilizes the contrastive loss to learn 4D invariance constraints. For our MM-based framework, we leverage the momentum contrast [18] to self-distill the consistency between the same scenes under different masking degrees, for unearthing the hidden information from unmasked areas.

### 3. Method

#### 3.1. Overview

As shown in Fig. 3, we propose MM-3DScene for masked modeling on 3D scenes, which can effectively alleviate the uncertainty and unpredictability during masked reconstruction for enhancing the self-supervised pre-training on complex 3D scenes.

Firstly, we propose a Local-Statistics-Guided masking strategy in Sec. 3.2, which explores the local difference for discovering and preserving Informative Points during the scene masking. Next, our mask strategy is integrated with a progressive reconstruction manner, yielding an informative-preserved reconstruction scheme. Finally, we introduce a self-distillation branch for learning the spatial consistency

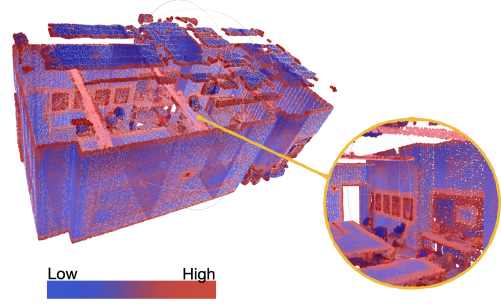


Figure 2. **The heat map of local statistics on 3D scenes.** Points with high statistics value also provide highly important semantic information for understanding or reconstructing 3D scenes.

in the unmasked areas under progressive masking ratios. By elegantly combining informative-preserved reconstruction and self-distilled consistency, MM-3DScene is yielded. The framework details are presented in Sec. 3.3.

#### 3.2. Local-Statistics-Guided Masking

As we analyzed in Sec. 1, the traditional random mask strategy is not feasible for complex large-scale 3D scenes, which often causes a high risk of ambiguity during masked reconstruction. In order to explore a more effective masked modeling mechanism for 3D scenes, we need to first design a better mask strategy, aiming to *reduce the ambiguity* during the masked reconstruction pretext task.

In 3D scenes, some representative structured points in provide highly important information hints and rich semantic context for assisting the scene understanding or reconstruction tasks. Accordingly, preserving them may help a lot to simplify the reconstruction. Thus, the first question is: *how to find such points?*

In this paper, we adopt local statistics as straightforward guidance to discover the representative points. For complex 3D scenes, we use the local difference of each point (*i.e.*, the difference between each point and its neighboring points in terms of colors and coordinates) to calculate the local statistics of each point. Specifically, for each point  $p_i$  in the scene, we first use K-Nearest Neighborhood (KNN) search to obtain its  $K$  neighboring points  $p_{ik}$  in Euclidean space, and then calculate the local difference to denote point statistics:

$$D_q(p_i) = \sum_{k=1}^K \|p_{i,q} - p_{ik,q}\|_2 \quad (1)$$

where  $p_{i,q} \in \mathbb{R}^{1 \times C_q}$  is point statistic of  $p_i$ . Specifically,  $p_{i,0}$  and  $p_{i,1}$  are the point coordinates and colors, which are normalized and accumulated together as the local statistics.  $D = \sum_{q=0}^1 (\alpha_q \times \text{norm}(D_q))$ , where  $D \in \mathbb{R}^{N \times 1}$ ,  $N$  is the number of the points in the scene.

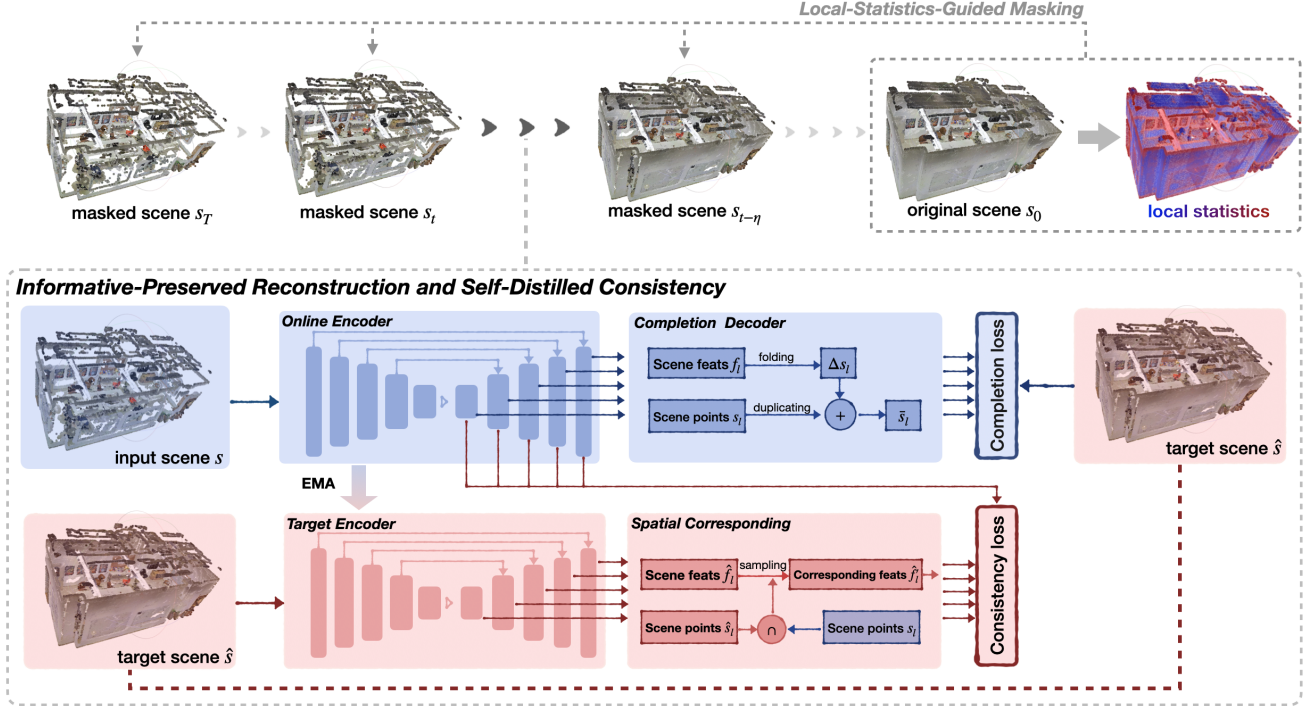


Figure 3. Overview of our customized masked modeling framework **MM-3DScene**. We first propose the Local-Statistics-Guided masking strategy to discover and preserve the representative structured points. This mask strategy yields an Informative-Preserved Reconstruction, where our method focuses on restoring the *regional* geometric patterns of **masked areas** at each learning step, enhancing the pretext masking task with less ambiguity. Moreover, since **unmasked areas** are underexplored during reconstruction, we introduce a self-distillation branch to maintain the intrinsic *spatial consistency* under progressive masking ratios, which enables MM-3DScene to learn consistent (*i.e.*, masking-invariant) representations of the unmasked areas.

For a clearer picture, we visualize the heat map of local statistics on the 3D scene in Fig. 2, it can be seen that points with high statistics value are concentrated in the foreground objects or contours of the scene (*red regions*), and these areas are relatively more informative and provide highly important information hints for understanding or reconstructing a 3D scene, denoted by ‘**Informative Points**’ in this paper. Informative Points provide highly useful information hints and rich semantic context for *assisting* masked reconstruction. To this end, our mask strategy is definite: to *preserve* Informative Points in a scene and mask other points. Through this, the basic geometric information of a scene is explicitly *retained*, which effectively reduces ambiguity in the pretext masking task.

### 3.3. MM-3DScene

After introducing the mask strategy with local statistics guidance, we illustrate the details of our MM-3DScene – a self-supervised masked modeling framework customized for 3D scene understanding. Our pre-training architecture consists of two parts: informative-preserved reconstruction and self-distillation of spatial consistency, respectively aiming at better modeling of masked areas and unmasked areas in 3D scenes. These two parts are elegantly combined while

*complementing* each other in our MM-3DScene.

**Informative-preserved reconstruction.** Based on the proposed mask strategy, we introduce an informative-preserved reconstruction manner.

To begin with, for the original scene  $s_0$ , we align all scene points in descending order according to the local statistics  $D$ . Then we use the *incremental* masking ratio  $\theta = \{\theta_1, \dots, \theta_t, \dots, \theta_T\}$  to *progressively* mask the scene  $s_0$  according to  $D$  from low to high, forming the scene sequence  $\mathcal{S} = \{s_1, \dots, s_{t-1}, s_t, \dots, s_T\}$  as shown in the Fig. 3. It can be seen that in the masked sequence  $\mathcal{S}$ , the masked regions are gradually shifted from the background surface to the foreground objects, and the representative structural points are preserved. It is also worth noting that scene  $s_t$  is the *subset* of scene  $s_{t-1}$ .

During each training iteration, we randomly select a masked scene  $s_t$  in sequence  $\mathcal{S}$  as the input, and take a *relatively* more complete one  $s_{t-\eta}$  as the target for masked reconstruction, where  $\eta$  indicates the masked gap to be recovered (to facilitate the subsequent description, we later define the input scene as  $s = s_t$  and the target as  $\hat{s} = s_{t-\eta}$ ). To encode point-wise representations from the input scene  $s$ , we utilize a well-established network  $\Phi_{OE}$  as our backbone.  $\Phi_{OE}$  is a point based feature extractor with hier-



archical structure. The hierarchical encoded features  $\mathcal{F}$  of scene  $s$  can be represented as  $\mathcal{F} = \Phi_{OE}(s)$  where  $\mathcal{F} = \{f_1, \dots, f_l, \dots, f_L\}$ , and  $L$  is the layers of  $\Phi_{OE}$ .

Then we utilize the features extracted from the online encoder to learn the coordinate variations for scene reconstruction. Note that the number of output points  $\hat{s}$  is larger than that of input points  $s$ . Before generating coordinate variations via MLP, we follow SnowflakeNet [59] to replicate the coordinates and features, producing more points for outputs. Specifically, for  $l$ -th layer, we have the displacement feature  $f_l$  and scene points  $s_l$ . Referring to FoldingNet [70], we incorporate a standard two-dimensional grid  $I$  to the displacement feature  $f_l$ , and then use two consecutive Multi-Layer Perceptrons (MLPs) [23] to generate coordinate variation  $\Delta s_l$ :

$$\Delta s_l = \Psi_2(f_l \oplus \Psi_1(f_l \oplus I)), \quad (2)$$

where  $\Psi_1$  and  $\Psi_2$  are two 3-layer MLPs, and  $\oplus$  is the concatenation operator. Then, we add the coordinate variation  $\Delta s_l$  with the duplicated input mask scene  $s$ , which generates the predicted scene  $\bar{s}_l$  in the  $l$ -th layer. Finally, in order to train the masked reconstruction task more efficiently, we choose the multi-scale symmetrical chamfer distance [53] as the reconstruction loss, with the following details:

$$\mathcal{L}_{PC} = \sum_l \left( \frac{1}{|\bar{s}_l|} \sum_{x \in \bar{s}_l} \min_{y \in \hat{s}} \|x - y\|^2 + \frac{1}{|\hat{s}|} \sum_{y \in \hat{s}} \min_{x \in \bar{s}_l} \|x - y\|^2 \right). \quad (3)$$

The proposed reconstruction scheme encourages the model to *focus* on restoring the *regional* geometry in a novel progressive manner, which enjoys less ambiguity and is able to restore accurate geometric information, so to enhance the modeling on masked areas.

**Consistency self-distillation.** Moreover, we realize the information of unmasked areas (*i.e.*, Informative Points) is underexplored. We find that there exists point correspondence in the unmasked areas under progressive masking ratios. Leveraging this unique behaviour of our method, and based on a generalized model of teacher-student mutual learning [13, 75], we introduce a self-distillation branch to maintain the intrinsic spatial consistency on unmasked areas during progressive reconstruction.

As Fig. 3 shows, we treat the online encoder as the student model and maintain a teacher model as the target encoder. It is worth mentioning that the target encoder has the same parameters and structure as the online encoder, but does not participate in the back-propagation of the gradient. The parameters of the target encoder are dynamically updated by the Exponential Moving Average (EMA) [13]:

$$\mathcal{W}_{\text{target}} \leftarrow [\beta \mathcal{W}_{\text{target}} + (1 - \beta) \mathcal{W}_{\text{online}}], \quad (4)$$

where  $\mathcal{W}_{\text{online}}$  and  $\mathcal{W}_{\text{target}}$  are the parameters of online encoder and target encoder.

Next, we feed the target scene  $\hat{s}$  to the target encoder and extract the target feature  $\hat{\mathcal{F}} = \{\hat{f}_1, \dots, \hat{f}_l, \dots, \hat{f}_L\}$ , where  $\hat{f}_l \in \mathbb{R}^{\hat{N}_l \times C}$ . Considering the input scene is the subset of target scene, where  $s \subset \hat{s}$ , we can select a subset of the target feature  $\hat{f}_l$  that have natural spatial correspondence with the online feature  $f_l$ . We define such subset of target feature as  $\hat{f}'_l \in \mathbb{R}^{N_l \times C}$ , where  $N_l < \hat{N}_l$ . Finally, we use info-NCE loss [41] to model the spatial consistency of the feature representation between the input scene and the target scene, which can be formulated as:

$$\mathcal{L}_{CSD} = - \sum_l \sum_{(i,j) \in s_l} \log \frac{\exp(f_{i,l} \cdot \hat{f}'_{j,l} / \tau)}{\sum_{(\cdot,k) \in s_l} \exp(f_{i,l} \cdot \hat{f}'_{k,l} / \tau)}, \quad (5)$$

where the input scene  $s_l$  is the set of points that can find the spatial correspondences in the target scene  $\hat{s}_l$ . For online feature  $f_{i,l}$  of point  $p_i$ , we take the corresponding point feature  $\hat{f}'_{j,l}$  in the target scene as the positive sample, and use feature  $\hat{f}'_{k,l}$  as the negative samples, where  $\exists(\cdot, k) \in s_l$  and  $k \neq j$ .

Through self-distillation, our method is able to unearth the spatial-consistent (*i.e.*, masking-invariant) representations from unmasked areas.

**MM-3DScene.** Ultimately, by elegantly combining informative- preserved reconstruction on masked areas and consistency self-distillation from unmasked areas, a unified framework called MM-3DScene is yielded. The final pre-training loss can be denoted by  $\mathcal{L} = \zeta_1 \times \mathcal{L}_{PC} + \zeta_2 \times \mathcal{L}_{CSD}$ . We find that setting both  $\zeta_1$  and  $\zeta_2$  to 1 yields the best result, which means that the masked reconstruction and the consistency distillation share *balanced importance* for our framework.

## 4. Experiments

We pre-train our MM-3DScene and demonstrate its effectiveness on a variety of generic downstream tasks for 3D scene understanding, including 3D semantic segmentation, 3D object detection, as well as data-efficient setup and linear probing evaluation.

### 4.1. Experimental Setup

**Datasets.** Following the state-of-the-art 3D self-supervised pre-training methods [8, 25, 61], we pre-train the proposed MM-3DScene on ScanNetv2 [10] dataset with 1201 train scans. As for 3D object detection, we fine-tune the pre-trained weights on ScanNetv2 training set (inner-domain) and evaluate on the corresponding validation set with 312 scenes. We also fine-tune our model on SUN RGB-D [49] training set for cross-domain evaluation. SUN RGB-D covers fully 3D object bounding box annotations for 10 categories, with 5,285 train frames and 5,050 test

Method	mAP@0.25	mAP@0.5
DSS [50]	15.2	6.8
F-PointNet [44]	19.8	10.8
GSPN [72]	30.6	17.7
3D-SIS [24]	40.2	22.5
VoteNet [43] (scratch)	58.7	35.4
RandomRooms [47] + VoteNet	61.3 (+2.6)	36.2 (+0.8)
PointContrast [61] + VoteNet	58.5 (-0.2)	38.0 (+2.6)
CSC [25] + VoteNet	-	39.3 (+3.9)
4DContrast [8] + VoteNet	-	40.0 (+4.6)
<b>MM-3DScene (w/o <math>L_{CSD}</math>) + VoteNet</b>	<b>61.9 (+3.2)</b>	<b>41.3 (+5.9)</b>
<b>MM-3DScene + VoteNet</b>	<b>63.1 (+4.4)</b>	<b>41.5 (+6.1)</b>

Table 2. 3D object detection results on ScanNetv2.

Method	mAP@0.25	mAP@0.5
DSS [50]	42.1	-
COG [48]	47.6	-
2D-driven [29]	45.1	-
F-PointNet [44]	54.0	-
VoteNet [43] (scratch)	57.7	32.9
PointContrast [61] + VoteNet	57.5 (-0.2)	34.8 (+1.9)
RandomRooms [47] + VoteNet	59.2 (+1.5)	35.4 (+2.5)
CSC [25] + VoteNet	-	36.4 (+3.5)
4DContrast [8] + VoteNet	-	38.2 (+5.3)
<b>MM-3DScene + VoteNet</b>	<b>60.6 (+2.9)</b>	<b>37.3 (+4.4)</b>

Table 3. 3D object detection results on SUN RGB-D.

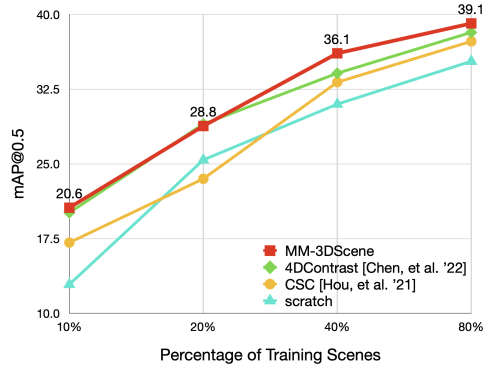


Figure 4. Data-efficient 3D object detection results on ScanNetv2.

frames. As for downstream semantic segmentation, in addition to conducting the experiments on the ScanNetv2, we also evaluate our pre-trained models on S3DIS [1], which contains 271 rooms from 3 different buildings.

**Backbone networks.** In 3D object detection, we follow [8, 25, 61] to use VoteNet [43] as the backbone. We directly perform MM-3DScene pre-training on the original PointNet++ layers, for minimizing the backbone effect on performance gains and better verifying the stand-alone effectiveness of pre-trained representations. Besides, we follow the data augmentations in the original backbone networks. For the scene semantic segmentation, we advocate using Point Transformer [80] as our backbone, considering

Method	mAcc	mIoU
PointNet++ [46]	-	53.5
PointConv [58]	-	61.0
PointASNL [69]	-	63.5
FPConv [33]	-	64.4
KPConv [54]	-	69.2
SR-UNet [9]	78.1	70.0
CSC [25] + SR-UNet	78.8 (+0.7)	70.7 (+0.7)
PointContrast [61] + SR-UNet	79.3 (+1.2)	71.3 (+1.3)
4DContrast [8] + SR-UNet	80.8 (+2.7)	72.3 (+2.3)
PointTrans [80] (scratch)	79.6	70.6
PointMAE [42] * + PointTrans	79.6 (+0.0)	70.6 (+0.0)
PointContrast [61] * + PointTrans	80.0 (+0.4)	70.9 (+0.3)
<b>MM-3DScene (w/o <math>L_{CSD}</math>) + PointTrans</b>	<b>81.2 (+1.6)</b>	<b>72.1 (+1.5)</b>
<b>MM-3DScene + PointTrans</b>	<b>82.0 (+2.4)</b>	<b>72.8 (+2.2)</b>

Table 4. 3D semantic segmentation results on ScanNetv2.

Method	mAcc	mIoU
PointNet [45]	49.0	41.1
PointCNN [32]	63.9	57.3
PointWeb [79]	66.6	60.3
PACConv [64]	73.0	66.6
KPConv [54]	72.8	67.1
PointTrans [80] (scratch)	76.5	70.4
PointMAE [42] * + PointTrans	76.4 (-0.1)	70.4 (+0.0)
PointContrast [61] * + PointTrans	76.9 (+0.4)	70.7 (+0.3)
<b>MM-3DScene + PointTrans</b>	<b>78.0 (+1.5)</b>	<b>71.9 (+1.5)</b>

Table 5. 3D semantic segmentation results on S3DIS Area-5.

its lightweight and effective attributes for 3D scene understanding (Table 6). We retain its original network architectures, and discard the output head in the pre-training process in order to obtain the fine-grained feature representation. The output head will be added to the network structure for the downstream training.

**Implementation details.** For pre-training, we use AdamW [37] optimizer with a weight decay of 0.0005. The initial learning rate is set to 0.001 and decayed at the 60% and 80% epochs. We pre-train the network for 300 epochs with a batch size of 8. The reconstruction gap  $\eta$  is set to 0.1. We set the scale parameter  $\tau$  as 1.

For downstream semantic segmentation, we follow the setting of Point Transformer [80], and utilize SGD optimizer with a momentum of 0.9 and weight decay of 0.0001. In addition to the standard data augmentation, we also use random cropping of the scene for more effective training. For S3DIS, we fine-tune our model for 150 epochs with a data loop of 30, where the voxel size is 4cm. The data loop for ScanNet is set to 6 and the scene resolution is 2cm.

For downstream object detection, the fine-tuning settings all follow 4DContrast [8], where the networks are trained for 500 epochs, the initial learning rate is set to 0.001 and decayed by a factor of 0.5 at epoch 250, 350, 450. The batch size is 6 for ScanNet and 16 for SUN RGB-D.

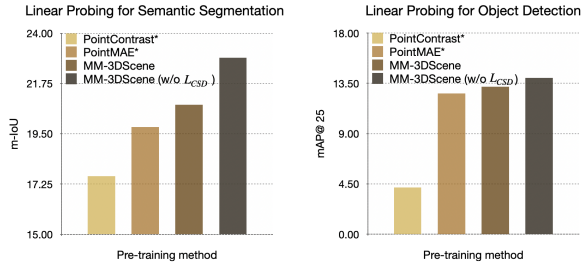


Figure 5. Linear probing on semantic segmentation and object detection.

## 4.2. 3D Object Detection

**ScanNetv2 object detection.** 3D object detection is a widely used downstream task for scene understanding. We report the fine-tuning results of object detection for ScanNetv2 in the Table 2, our MM-3DScene achieves the state-of-the-art performance among the SSL methods. With of MM-3DScene pre-training, it achieves a significant improvement of **4.4** on mAP@0.25, and **6.1** on mAP@0.5 compared to training from scratch. Meanwhile, even if we train without the consistency distillation, the result also surpasses the state-of-the-art method 4DContrast [8], which strongly demonstrates the effectiveness of our customized masked modeling method for object detection. Note that, *without* RGB, our method still gets **41.5** mAP@0.5 on ScanNetv2 detection, which is much higher than 4DC (40.0) and PC (38.0). Moreover, we conduct experiments with the updated H3DNet [78] as the backbone and achieve significant improvements, which can be found in appendix.

**SUN RGB-D object detection.** We also conduct the cross-domain experiments on SUN RGB-D dataset for object detection, which is shown in Table 3. We pre-train our method on ScanNetv2 and fine-tune on SUN RGB-D, for which our method substantially exceeds the training-from-scratch baseline and surpasses 3D-based pre-training of CSC [25], PointContrast [61] and RandomRooms [47]. *Notably*, pretraining 4DContrast [8] requires a *prerequisite* non-trivial generation of spatio-temporal correspondence, while our method produces the masked scenes *on the fly*.

**Data-efficient evaluation.** As shown in Fig. 4, We evaluate our method for fine-tuning with limited training data, where the results of CSC and 4DContrast are officially borrowed from [8]. Our MM-3DScene surpasses both of them in almost all data-efficient settings. It is worth noting that our method achieves **39.1** mAP@0.5 when finetuning with only **80%** training data, which even exceeds many methods trained on 100% of the training data in Table 2.

## 4.3. 3D Semantic Segmentation

**ScanNetv2 semantic segmentation.** Semantic segmentation is a common but challenging task. It requires models

Method	Model Size	mIoU
SR-UNet [9] (scratch)	37.85M	70.0
4DContrast (3D) [8] + SR-UNet	70.85M (+33M)	71.7 (+1.7)
4DContrast (4D) [8] + SR-UNet	75.85M (+38M)	72.3 (+2.3)
PointTrans [80] (scratch)	7.76M	70.6
MM3DScene + PointTrans	8.63M (+0.87M)	72.8 (+2.2)

Table 6. Comparisons of model parameters for different pre-training methods.

to predict the semantic classes of each point, which is involved in the fine-grained understanding of 3D scenes. Table 4 lists the downstream semantic segmentation results on ScanNetv2. For SR-UNet based methods, we directly quote the results from [8], which do not use the RGB information as input. However, for Point Transformer backbone, the RGB information is indispensable and omitting it will cause a deteriorated performance drop. Thus, we follow its original input setting (points + rgb) and apply various pre-training methods on it. We denote the methods with \* when they are adapted to Point Transformer backbone, e.g., “PointMAE [42]\* + PointTrans”. The table shows that the random masking strategy of PointMAE does not help the downstream segmentation task. We also observe that PointContrast\* [61] has a marginal improvement over the scratch model, which may indicate that contrastive learning is good at learning discriminative features but not robust to occlusions (as suggested by [26, 28]). In contrast, our MM-3DScene achieves superior results over these conventional pre-training methods. Furthermore, our MM-3DScene also outperforms 4DContrast with SRUNet backbone, while having a much lighter model size (Table 6).

**S3DIS semantic segmentation.** To further validate the effectiveness of our method, we conduct experiments on the S3DIS dataset [1], the results are shown in Table 5. We achieve a competitive result of 71.9% mIoU on S3DIS Area-5, which gets a relative improvement of 1.5% from scratch. The qualitative results are illustrated in appendix.

## 4.4. Linear Probing

Linear probing is *underexplored* in 3D scene understanding, we remedy this defect by locking the pre-trained backbone network, and only fine-tuning the segmentation head and detection head. The experiments are conducted on S3DIS Area-5 and the ScanNetv2 validation set. Fig. 5 shows results following different pre-training methods. Compared to the conventional random masking method (*i.e.*, PointMAE [42]) or contrastive-based method (*i.e.*, PointContrast [61]), our MM-3DScene performs best on the linear probing task. It is worth noting that our framework without learning spatial consistency performs better on linear probing, we guess the reason is that *fewer constraints* can enhance the generalizability.

## 4.5. Ablation Studies

**Why MM-3DScene is a better way for masked modeling on large-scale scenes?** Table 7 shows the ablation study of our MM-3DScene framework on S3DIS semantic segmentation. When we only use the informative masked modeling, (ii) achieves an improvement of 0.7% mIoU over the train-from-scratch method. On the other hand, focusing on the unmasked areas (setting iii), we only learn the spatial consistency of the unmasked informative points during each iteration, which still makes performance gains for downstream segmentation. Finally, setting iv integrates informative-preserved modeling of the masked areas and spatial consistency learning on unmasked areas together, which further boosts the result to 71.9% mIoU.

**How MM-3DScene simplifies the pretext task?** Table 8 shows the ablation study on mask strategies. Comparing the settings of (a) and (d) in Table 8, we found that our informative-preserved masked modeling can significantly improve the downstream performance, while random mask modeling does not. The possible reason is that such masking will randomly drop Informative Points, thus bringing great uncertainty to the pretext task. With the same consistency loss, PointMAE gets 71.0%, which is 0.9% lower than ours. While our approach can make the masked areas perceptible via informative-preserved masking. Instead, when we use the *informative-abandoned* masking setting (c) by masking off the informative structured points, the downstream performance is significantly reduced, demonstrating that masking the Informative Points may cause large unpredictability.

Moreover, different from (d), setting (e) utilizes the progressive manner by reconstructing a masked scene into a more complete one, and shows a significant improvement. We also apply the progressive reconstruction manner based on the random masking (setting (b)), and the results are improved accordingly.

**Memory cost.** Current 3D pre-training methods [8, 25, 61] basically adopts SR-UNet [9] as the backbone with 37.85M parameters. We advocate the much more lightweight Point Transformer [80], with just 7.76M parameters. In Table 6, we compare our model size with 4DContrast [8] which uses SR-UNet. Our MM-3DScene *only adds* 0.87M parameters to the backbone. Besides, our approach does not require the time-consuming pre-generation of contrastive point cloud pairs for pretraining used in PointContrast [61] and CSC [25]. Our masked scenes are generated *on the fly* during network training.

## 5. Conclusion

We have presented MM-3DScene, a customized masked modeling framework for 3D scene understanding. It explicitly preserves the representative structured points, which

	Pre-training method	Pre-training loss	mIoU	mAcc
i	Scratch	-	70.4	76.5
ii	MM-3DScene (MM only)	$\mathcal{L}_{PC}$	71.1 (+0.7)	77.2 (+0.7)
iii	MM-3DScene (consistency only)	$\mathcal{L}_{CSD}$	70.9 (+0.5)	77.1 (+0.6)
iv	MM-3DScene	$\mathcal{L}_{PC}, \mathcal{L}_{CSD}$	71.9 (+1.5)	78.0 (+1.5)

Table 7. Ablation study on MM-3DScene framework on S3DIS semantic segmentation.

	Masked Modeling	Mask Strategy	Progressive	mIoU	mAcc
	Scratch	-	-	70.4	76.5
(a)	PointMAE* [42]	Random	×	70.4	76.4
(b)	PointMAE* [42]	Random	✓	70.6	76.5
(c)	MM-3DScene (MM only)	Informative-abandoned	×	70.2	76.1
(d)	MM-3DScene (MM only)	Informative-preserved	×	70.9	76.9
(e)	MM-3DScene (MM only)	Informative-preserved	✓	71.1	77.2
(f)	MM-3DScene	Informative-preserved	✓	71.9	78.0

Table 8. Ablation study on mask strategies on S3DIS semantic segmentation.

provides highly useful information clues to simplify the pretext task of masked reconstruction. At each learning step, a masked scene is reconstructed in a progressive manner, so that to focus on restoring regional geometry and enjoy less ambiguity. Moreover, a self-distillation branch is integrated for maintaining the intrinsic spatial consistency on unmasked areas under the progressive masking ratios. Extensive experiments on various downstream tasks verified that our MM-3DScene significantly boosts the performance of baseline models trained from scratch.

**Limitations.** In this paper, we mainly focus on indoor scenes, following recent self-supervised methods [8, 25] for 3D scene understanding. We believe that the generic design insight of our masked modeling may inspire more researchers to solve 3D outdoor perception, 3D shape understanding, and 2D image recognition.

## Acknowledgments

This work was supported in part by the National Key R&D Program of China (NO. 2022ZD0160505), and in part by the Youth Innovation Promotion Association of Chinese Academy of Sciences (No. 2020355). It was also partially supported by NSFC62172348, Outstanding Youth Fund of Guangdong Province with No. 2023B1515020055 and Shenzhen General Project with No. JCYJ20220530143604010, the National Key R&D Program of China with grant No. 2018YFB1800800, by Shenzhen Outstanding Talents Training Fund 202002, by Guangdong Research Projects No. 2017ZT07X152 and No. 2019CX01X104, by the Guangdong Provincial Key Laboratory of Future Networks of Intelligence (Grant No. 2022B1212010001), and by Shenzhen Key Laboratory of Big Data and Artificial Intelligence (Grant No. ZDSYS201707251409055). Besides, thanks to Ji Hou for helpful suggestions on experiments.



## References

- [1] I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese. 3d semantic parsing of large-scale indoor spaces. In *CVPR*, 2016. 2, 6, 7
- [2] Matan Atzmon, Haggai Maron, and Yaron Lipman. Point convolutional neural networks by extension operators. *arXiv preprint arXiv:1803.10091*, 2018. 13
- [3] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEit: BERT pre-training of image transformers. In *ICLR*, 2022. 3
- [4] et.al Brown, Tom. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *NeurIPS*, 2020. 3
- [5] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. In *3DV*, 2017. 2
- [6] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, 2021. 3
- [7] Xinlei Chen\*, Saining Xie\*, and Kaiming He. An empirical study of training self-supervised vision transformers. In *ICCV*, 2021. 3
- [8] Yujin Chen, Matthias Nießner, and Angela Dai. 4dcontrast: Contrastive learning with dynamic correspondences for 3d scene understanding. In *ECCV*, 2022. 3, 5, 6, 7, 8
- [9] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *CVPR*, 2019. 6, 7, 8
- [10] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017. 2, 5
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 2, 3
- [12] Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence embeddings. In *EMNLP*, 2021. 3
- [13] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A survey. *IJCV*, 2021. 5
- [14] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *CVPR*, 2018. 2
- [15] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, koray kavukcuoglu, Remi Munos, and Michal Valko. Bootstrap your own latent - a new approach to self-supervised learning. In *NeurIPS*, 2020. 3
- [16] Ankur Handa, Viorica Patraucean, Vijay Badrinarayanan, Simon Stent, and Roberto Cipolla. Scenenet: Understanding real world indoor scenes with synthetic data. In *CVPR*, 2016. 2
- [17] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022. 2, 3
- [18] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 2, 3
- [19] Tong He, Dong Gong, Zhi Tian, and Chunhua Shen. Learning and memorizing representative prototypes for 3d point cloud semantic and instance segmentation. In *ECCV*, 2020. 2
- [20] Tong He, Chunhua Shen, and Anton van den Hengel. DyCo3d: Robust instance segmentation of 3d point clouds through dynamic convolution. In *CVPR*, 2021. 2
- [21] Tong He, Wei Yin, Chunhua Shen, and Anton van den Hengel. Pointinst3d: Segmenting 3d instances by points. In *ECCV*, 2022. 2
- [22] Georg Hess, Johan Jaxing, Elias Svensson, David Hagerman, Christoffer Petersson, and Lennart Svensson. Masked autoencoders for self-supervised learning on automotive point clouds. *arXiv preprint arXiv:2207.00531*, 2022. 3
- [23] Kurt Hornik. Approximation capabilities of multilayer feed-forward networks. *Neural Netw.*, 1991. 5, 12
- [24] Ji Hou, Angela Dai, and Matthias Nießner. 3d-sis: 3d semantic instance segmentation of rgb-d scans. In *CVPR*, 2019. 6
- [25] Ji Hou, Benjamin Graham, Matthias Nießner, and Saining Xie. Exploring data-efficient 3d scene understanding with contrastive scene contexts. In *CVPR*, 2021. 3, 5, 6, 7, 8
- [26] Zhicheng Huang, Xiaojie Jin, Chengze Lu, Qibin Hou, Ming-Ming Cheng, Dongmei Fu, Xiaohui Shen, and Jiashi Feng. Contrastive masked autoencoders are stronger vision learners. *arXiv preprint arXiv:2207.13532*, 2022. 7
- [27] Li Jiang, Hengshuang Zhao, Shaoshuai Shi, Shu Liu, Chi-Wing Fu, and Jiaya Jia. Pointgroup: Dual-set point grouping for 3d instance segmentation. *CVPR*, 2020. 1
- [28] Li Jing, Jiachen Zhu, and Yann LeCun. Masked siamese convnets. *arXiv preprint arXiv:2206.07700*, 2022. 7
- [29] Jean Lahoud and Bernard Ghanem. 2d-driven 3d object detection in rgb-d images. In *ICCV*, 2017. 6
- [30] Xin Lai, Jianhui Liu, Li Jiang, Liwei Wang, Hengshuang Zhao, Shu Liu, Xiaojuan Qi, and Jiaya Jia. Stratified transformer for 3d point cloud segmentation. In *CVPR*, 2022. 12, 13
- [31] Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In *CVPR*, 2018. 13
- [32] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. In *NeurIPS*, 2018. 2, 6, 13
- [33] Yiqun Lin, Zizheng Yan, Haibin Huang, Dong Du, Ligang Liu, Shuguang Cui, and Xiaoguang Han. Fpconv: Learning local flattening for point convolution. In *CVPR*, 2020. 6
- [34] Haotian Liu, Mu Cai, and Yong Jae Lee. Masked discrimination for self-supervised learning on point clouds. In *ECCV*, 2022. 2
- [35] Yongcheng Liu, Bin Fan, Shiming Xiang, and Chunhong Pan. Relation-shape convolutional neural network for point cloud analysis. In *CVPR*, 2019. 2
- [36] Ze Liu, Zheng Zhang, Yue Cao, Han Hu, and Xin Tong. Group-free 3d object detection via transformers. In *ICCV*, 2021. 2

- [37] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. [6](#)
- [38] Xu Ma, Can Qin, Haoxuan You, Haoxi Ran, and Yun Fu. Rethinking network design and local geometry in point cloud: A simple residual MLP framework. In *ICLR*, 2022. [2](#)
- [39] Chen Min, Xinli Xu, Dawei Zhao, Liang Xiao, Yiming Nie, and Bin Dai. Voxel-mae: Masked autoencoders for pre-training large-scale point clouds. *arXiv preprint arXiv:2206.09900*, 2022. [3](#)
- [40] Ishan Misra, Rohit Girdhar, and Armand Joulin. An End-to-End Transformer Model for 3D Object Detection. In *ICCV*, 2021. [2](#)
- [41] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. [5](#)
- [42] Yatian Pang, Wenxiao Wang, Francis E. H. Tay, Wei Liu, Yonghong Tian, and Li Yuan. Masked autoencoders for point cloud self-supervised learning. In *ECCV*, 2022. [2](#), [3](#), [6](#), [7](#), [8](#)
- [43] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. In *ICCV*, 2019. [1](#), [2](#), [6](#), [12](#), [13](#)
- [44] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *CVPR*, 2018. [6](#)
- [45] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, 2017. [2](#), [6](#), [13](#)
- [46] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NeurIPS*, 2017. [2](#), [6](#)
- [47] Yongming Rao, Benlin Liu, Yi Wei, Jiwen Lu, Cho-Jui Hsieh, and Jie Zhou. Randomrooms: Unsupervised pre-training from synthetic shapes and randomized layouts for 3d object detection. In *ICCV*, 2021. [6](#), [7](#)
- [48] Zhile Ren and Erik B. Sudderth. Three-dimensional object detection and layout prediction using clouds of oriented gradients. In *CVPR*, 2016. [6](#)
- [49] Shuran Song, Samuel P. Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *CVPR*, 2015. [5](#)
- [50] Shuran Song and Jianxiong Xiao. Deep sliding shapes for amodal 3d object detection in rgb-d images. In *CVPR*, 2016. [6](#)
- [51] Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-Yi Zhou. Tangent convolutions for dense prediction in 3d. In *CVPR*, 2018. [13](#)
- [52] Lyne Tchapmi, Christopher Choy, Iro Armeni, JunYoung Gwak, and Silvio Savarese. Segcloud: Semantic segmentation of 3d point clouds. In *3DV*, 2017. [13](#)
- [53] Lyne P Tchapmi, Vineet Kosaraju, Hamid Rezatofighi, Ian Reid, and Silvio Savarese. Topnet: Structural point cloud decoder. In *CVPR*, 2019. [5](#)
- [54] Hugues Thomas, Charles R. Qi, Jean-Emmanuel Deschaud, Beatriz Marcotequi, François Goulette, and Leonidas J. Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *ICCV*, 2019. [2](#), [6](#), [13](#)
- [55] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. [2](#)
- [56] Hanchen Wang, Qi Liu, Xiangyu Yue, Joan Lasenby, and Matthew J. Kusner. Unsupervised point cloud pre-training via occlusion completion. In *ICCV*, 2021. [2](#), [3](#)
- [57] 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 Trans. Graph.*, 2019. [2](#)
- [58] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *CVPR*, 2019. [2](#), [6](#)
- [59] Peng Xiang, Xin Wen, Yu-Shen Liu, Yan-Pei Cao, Pengfei Wan, Wen Zheng, and Zhizhong Han. Snowflakenet: Point cloud completion by snowflake point deconvolution with skip-transformer. In *ICCV*, 2021. [5](#)
- [60] Tiange Xiang, Chaoyi Zhang, Yang Song, Jianhui Yu, and Weidong Cai. Walk in the cloud: Learning curves for point clouds shape analysis. In *ICCV*, 2021. [2](#)
- [61] Saining Xie, Jiatao Gu, Demi Guo, Charles R. Qi, Leonidas Guibas, and Or Litany. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In *ECCV*, 2020. [3](#), [5](#), [6](#), [7](#), [8](#)
- [62] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In *CVPR*, 2022. [2](#), [3](#)
- [63] Mutian Xu, Pei Chen, Haolin Liu, and Xiaoguang Han. To-scene: A large-scale dataset for understanding 3d tabletop scenes. In *ECCV*, 2022. [2](#)
- [64] Mutian Xu, Runyu Ding, Hengshuang Zhao, and Xiaojuan Qi. Paconv: Position adaptive convolution with dynamic kernel assembling on point clouds. In *CVPR*, 2021. [2](#), [6](#)
- [65] Mutian Xu, Junhao Zhang, Zhipeng Zhou, Mingye Xu, Xiaojuan Qi, and Yu Qiao. Learning geometry-disentangled representation for complementary understanding of 3d object point cloud. In *AAAI*, 2021. [2](#)
- [66] Mingye Xu, Zhipeng Zhou, and Yu Qiao. Geometry sharing network for 3d point cloud classification and segmentation. In *AAAI*, 2020. [2](#)
- [67] Mingye Xu, Zhipeng Zhou, Junhao Zhang, and Yu Qiao. Investigate indistinguishable points in semantic segmentation of 3d point cloud. In *AAAI*, 2021. [2](#), [13](#)
- [68] Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. Spidercnn: Deep learning on point sets with parameterized convolutional filters. In *ECCV*, 2018. [2](#)
- [69] Xu Yan, Chaoda Zheng, Zhen Li, Sheng Wang, and Shuguang Cui. Pointasnl: Robust point clouds processing using nonlocal neural networks with adaptive sampling. In *CVPR*, 2020. [6](#)
- [70] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *CVPR*, 2018. [5](#)
- [71] Xiaoqing Ye, Jiamao Li, Hexiao Huang, Liang Du, and Xiaolin Zhang. 3d recurrent neural networks with context fusion for point cloud semantic segmentation. In *ECCV*, 2018. [13](#)

- [72] Li Yi, Wang Zhao, He Wang, Minhyuk Sung, and Leonidas Guibas. Gspn: Generative shape proposal network for 3d instance segmentation in point cloud. In *CVPR*, 2019. 6
- [73] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert: Pre-training 3d point cloud transformers with masked point modeling. In *CVPR*, 2022. 2, 3
- [74] Chris Zhang, Wenjie Luo, and Raquel Urtasun. Efficient convolutions for real-time semantic segmentation of 3d point clouds. In *3DV*, 2018. 13
- [75] Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *CVPR*, 2019. 5
- [76] Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hongsheng Li. Point-m2ae: Multi-scale masked autoencoders for hierarchical point cloud pre-training. In *NeurIPS*, 2022. 2
- [77] Yabin Zhang, Jiehong Lin, Chenhang He, Yongwei Chen, Kui Jia, and Lei Zhang. Masked surfel prediction for self-supervised point cloud learning. *arXiv preprint arXiv:2207.03111*, 2022. 2
- [78] Zaiwei Zhang, Bo Sun, Haitao Yang, and Qixing Huang. H3dnet: 3d object detection using hybrid geometric primitives. In *ECCV*, 2020. 7, 12, 13
- [79] Hengshuang Zhao, Li Jiang, Chi-Wing Fu, and Jiaya Jia. PointWeb: Enhancing local neighborhood features for point cloud processing. In *CVPR*, 2019. 2, 6, 13
- [80] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip Torr, and Vladlen Koltun. Point transformer. In *ICCV*, 2021. 1, 2, 6, 7, 8, 12, 13

## Appendix A: Visualization Results

Fig. 8 visualizes the masked reconstruction results of MM-3DScene. It can be observed that: **i)** For masked **input**, our mask strategy *preserves* Informative Points to provide basic geometric information, which explicitly reduces the ambiguity during masked reconstruction. **ii)** For **target**, instead of being the original intact scene, it is a relatively more complete one with a smaller masking ratio. This prompts models to *concentrate* on reconstructing the local regional 3D structures where models focus on recovering *regional* geometric patterns. **iii)** For reconstruction **result**, our model is able to recover the masked areas, suggesting it successfully learned numerous visual representations. For example, our method works well to recover details of the masked foreground objects (*e.g.* table and chair). For the background surfaces (*e.g.* floor, wall), our method can also achieve a smooth and complete recovery. In addition to the visualization of the pre-training reconstruction results, as shown in the Fig. 6, we present the visualization results of the downstream semantic segmentation task. It can be seen that compared with other methods, our method can correct the results in some areas where the prediction is inaccurate.

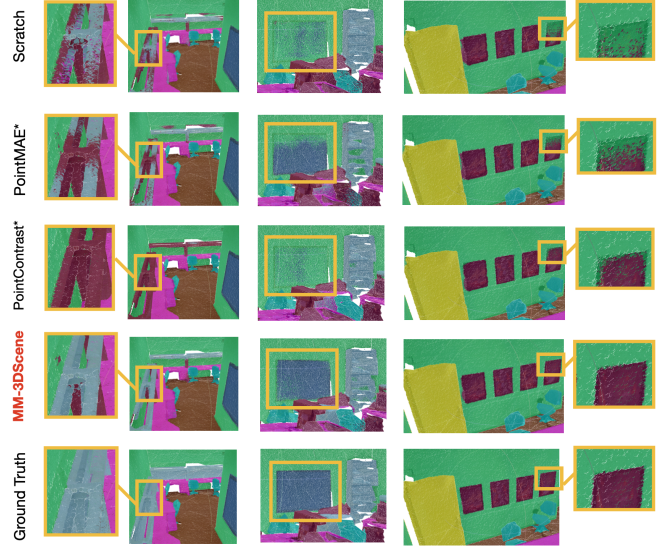


Figure 6. Qualitative results on S3DIS semantic segmentation.

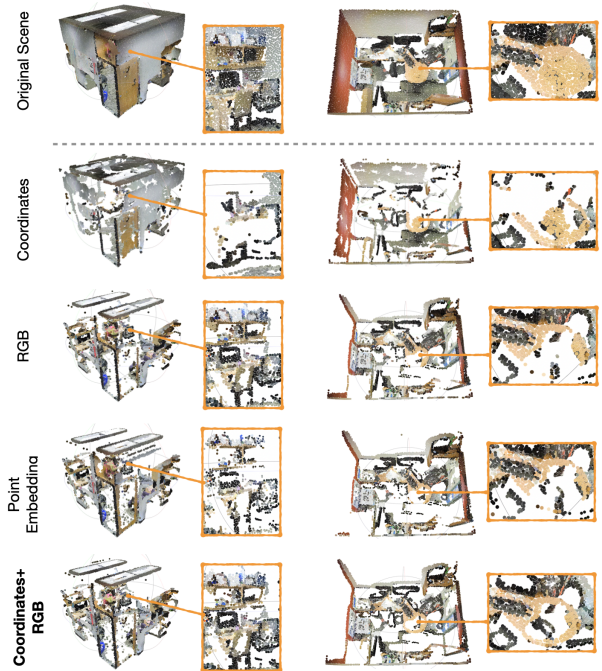


Figure 7. Visualization of masked scene guided by different formats of local statistics. In this figure, the **contour of a table** (*i.e.*, the representative geometric structures) is accurately found and preserved, when calculating the local difference of coordinates+RGB.

## Appendix B: More Ablations of Masked Reconstruction

**More formats of local statistics.** In our Local-Statistics-Guided Masking, we exploit local statistics to discover in-



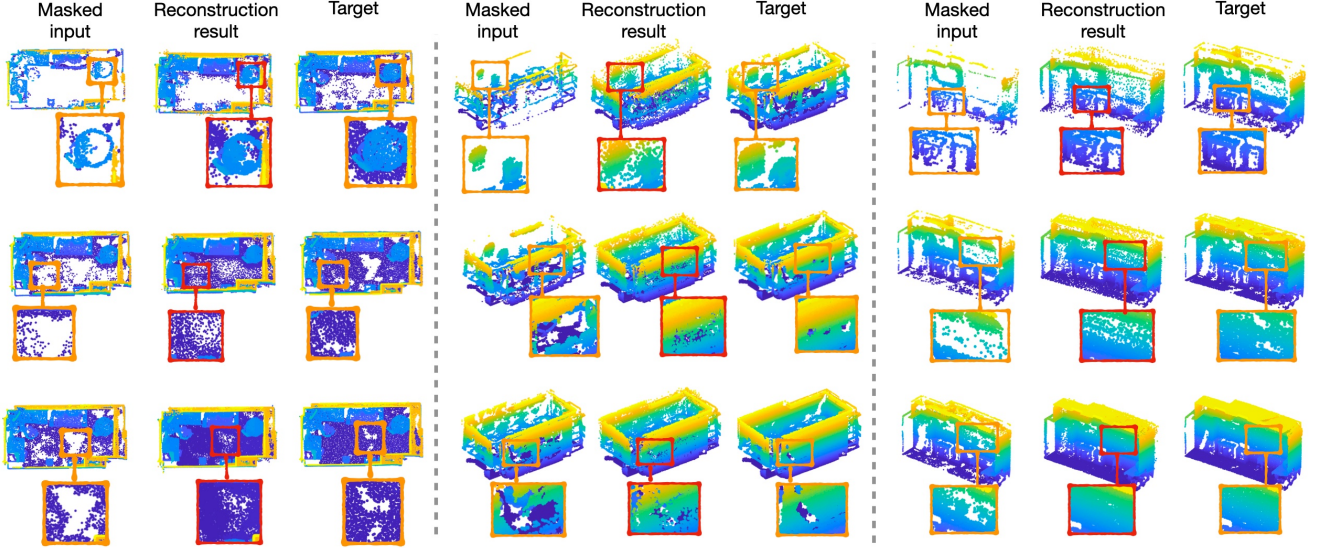


Figure 8. Qualitative results of MM-3DScene masked reconstruction pretext task. Our method can effectively reconstruct the masked areas, suggesting it learned rich visual representations for understanding 3D scenes.

formative points, aiming to *accurately* retain the representative geometric structures. The local statistics are denoted by the local difference between each point and its neighboring points in terms of coordinates and colors. Here we investigate other possible formats of the local difference, including point embedding difference (gained by applying MLPs [23] on each point), only coordinates difference, and only RGB difference. Fig. 7 shows that considering the local difference of both coordinates and RGB is the most applicable way to find informative points. For example in this figure, the *structural contour* of a table is *accurately* found and preserved, when guided by the local differences of coordinates+RGB, which provides useful information hints for recovering the masked interior area of this table. As a result, our method performs best when the masking is guided by the local statistic of coordinates+RGB, as listed in Table 9.

**Ablation studies of reconstruction gap.** Our method uses the incremental masking ratio  $\theta = \{\theta_1, \dots, \theta_t, \dots, \theta_T\}$  to progressively mask the scene. During the masked reconstruction, the masking ratio is  $\theta_t$  for the input scene, and  $\theta_{t-\eta}$  for the target scene, where  $\eta$  indicates the masked gap to be recovered and latently influences the *difficulty* of the pretext task. Fig. 9 provides the ablation study of such reconstruction gap, where our model enjoys the least difficulty and performs best under  $\theta_t - \theta_{t-\eta} = 0.1$ , and degrades when the gap becomes larger. Additionally, we also implement the random reconstruction gap, which probably causes more ambiguity, yielding 70.36% mIoU.

	Pre-training method	Local Statistic	mIoU	mAcc
i	Scratch	-	70.4	76.5
ii	MM-3DScene (w/o $L_{CSD}$ )	Coordinates	70.6 (+0.2)	76.5 (+0.0)
iii	MM-3DScene (w/o $L_{CSD}$ )	RGB	70.9 (+0.5)	77.0 (+0.5)
iv	MM-3DScene (w/o $L_{CSD}$ )	point embedding	70.9 (+0.5)	76.9 (+0.4)
v	MM-3DScene (w/o $L_{CSD}$ )	<b>Coordinates + RGB</b>	<b>71.1 (+0.7)</b>	<b>77.2 (+0.7)</b>

Table 9. MM-3DScene (w/o  $L_{CSD}$ ) guided by **different formats of local statistics** for S3DIS semantic segmentation.

Method	Model Size	Train Time	Infer Time	mAP@0.25	mAP@0.5
VoteNet [43] (scratch)	0.95M	5.9h	0.2s	58.7	35.4
MM3DScene + VoteNet	1.48M	15.6h	-	63.1 (+4.4)	41.5 (+6.1)
H3DNet [78] (scratch)	4.74M	12.3h	7.3s	64.8	47.4
MM3DScene + H3DNet	6.87M	36.1h	-	66.8 (+2.0)	48.9 (+1.5)

Table 10. 3D object detection results on ScanNetv2. The baseline results come from official code implementations. The training and inference times are evaluated with the same training settings.

Method	Model Size	Train Time	Infer Time	mIoU
PointTrans [80] (scratch)	7.76M	17.3h	4.36s	70.4
MM3DScene + PointTrans	8.63M	29.1h	-	71.9 (+1.5)
Stratified Trans* [30] (scratch)	8.02M	45.7h	11.77s	70.3
MM-3DScene + Stratified Trans*	8.89M	73.2h	-	71.6(+1.3)

Table 11. 3D semantic segmentation results on S3DIS. The baseline results come from official code implementations. The training and inference times are evaluated with the same training settings.

## Appendix C: Other Backbones with MM-3DScene

In the main paper, we adopt VoteNet [43] as the backbone for object detection, and Point Transformer [80] for semantic segmentation. In this section, we utilize other backbone networks for verifying the generalization ability of our MM-3DScene.



Method	mIoU	ceiling	floor	wall	beam	col.	win.	door	table	chair	sofa	bookc.	board	clu.
PointNet [45]	41.1	88.8	97.3	69.8	0.1	3.9	46.3	10.8	58.9	52.6	5.9	40.3	26.4	33.2
SegCloud [52]	48.9	90.1	96.1	69.9	0.0	18.4	38.4	23.1	70.4	75.9	40.9	58.4	13.0	41.6
TangentConv [51]	52.6	90.5	97.7	74.0	0.0	20.7	39.0	31.3	77.5	69.4	57.3	38.5	48.8	39.8
SPGraph [31]	58.0	89.4	96.9	78.1	0.0	42.8	48.9	61.6	84.7	75.4	69.8	52.6	2.1	52.2
PCNN [2]	58.3	92.3	96.2	75.9	0.3	6.0	69.5	63.5	65.6	66.9	68.9	47.3	59.1	46.2
RNNFusion [71]	57.3	92.3	98.2	79.4	0.0	17.6	22.8	62.1	80.6	74.4	66.7	31.7	62.1	56.7
Eff 3D Conv [74]	51.8	79.8	93.9	69.0	0.2	28.3	38.5	48.3	73.6	71.1	59.2	48.7	29.3	33.1
PointCNN [32]	57.3	92.3	98.2	79.4	0.0	17.6	22.8	62.1	74.4	80.6	31.7	66.7	62.1	56.7
PointWeb [79]	60.3	92.0	98.5	79.4	0.0	21.1	59.7	34.8	76.3	88.3	46.9	69.3	64.9	52.5
IAF-Net [67]	64.6	91.4	98.6	81.8	0.0	34.9	62.0	54.7	79.7	86.9	49.9	72.4	74.8	52.1
KPCov [54]	67.1	92.8	97.3	82.4	0.0	23.9	58.0	69.0	81.5	91.0	75.4	75.3	66.7	58.9
PointTransformer [80]	70.4	94.0	98.5	86.3	0.0	38.0	63.4	74.3	89.1	82.4	74.3	80.2	76.0	59.3
MM-3DScene(Ours)	71.9	94.6	98.6	87.1	0.0	44.2	62.9	79.2	90.7	81.7	74.3	81.4	79.3	60.3

Table 12. Semantic segmentation results on S3DIS dataset evaluated on Area 5.

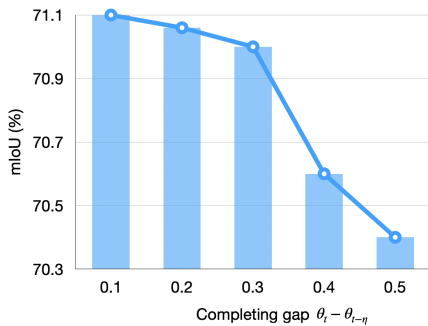


Figure 9. Ablation study of **masked reconstruction gap** on S3DIS semantic segmentation (based on MM-3DScene w/o  $L_{CSD}$ ).

**H3DNet object detection.** We apply our MM-3DScene pretrained framework on H3DNet [78] which is a more powerful network using hybrid geometric primitives based on VoteNet [43]. As shown in Table. 10, MM-3DScene improves the H3DNet with the mAP@0.25 by 2.0 and mAP@0.5 by 1.5, exceeding the performance with VoteNet as the backbone.

**Stratified Transformer semantic segmentation.** We also evaluate the performance of Stratified Transformer [30] as the backbone on S3DIS semantic segmentation. We reproduce the backbone performance using its official code and report the results in Table. 11. Our MM-3DScene surpasses Stratified Transformer by 1.3% mIoU. However, it comes with a high computational cost (2.5 times of MM3D-Scene + PT) and a long training time.

**Discussions.** Although both H3DNet [78] and Stratified Transformer [30] inherit VoteNet [43] and Point Transformer [80], and achieve decent performance, they introduce highly-engineered architectures tailored to their network-specific operations, making it difficult to evaluate the improvement made by the self-supervised frameworks.

Thus, we advocate simple and classical baselines, with the goal of minimizing the influence of network architectures to better measure the performance gain *purely* from the self-supervised pretraining framework – MM-3DScene.

Moreover, both Point Transformer [80] and VoteNet [43] stand out with conspicuously excellent **efficiency**, as reflected in *model size*, *training time*, and *inference time* of Table 10 and Table 11, which is highly important for the deployment on real applications.

## Appendix D: More fine-grained quantitative results

To provide a more comprehensive analysis, we present the segmentation results of each category in Table 12. We observe that most categories have different degrees of improvement over the Point Transformer [80] backbone that we use. For instance, we achieve 6.2% gain on column, 4.9% on door, 3.3% on board, and slight decrease on window and chair.