

Masked Discrimination for Self-Supervised Learning on Point Clouds

Haotian Liu Mu Cai Yong Jae Lee

University of Wisconsin–Madison
{lht,mucai,yongjaelee}@cs.wisc.edu

Abstract. Masked autoencoding has achieved great success for self-supervised learning in the image and language domains. However, mask based pretraining has yet to show benefits for point cloud understanding, likely due to standard backbones like PointNet being unable to properly handle the training versus testing distribution mismatch introduced by masking during training. In this paper, we bridge this gap by proposing a discriminative mask pretraining Transformer framework, *MaskPoint*, for point clouds. Our key idea is to represent the point cloud as discrete occupancy values (1 if part of the point cloud; 0 if not), and perform simple binary classification between masked object points and sampled noise points as the proxy task. In this way, our approach is robust to the point sampling variance in point clouds, and facilitates learning rich representations. We evaluate our pretrained models across several downstream tasks, including 3D shape classification, segmentation, and real-world object detection, and demonstrate state-of-the-art results while achieving a significant pretraining speedup (e.g., $4.1\times$ on ScanNet) compared to the prior state-of-the-art Transformer baseline.¹

1 Introduction

Learning rich feature representations without human supervision, also known as self-supervised learning, has made tremendous strides in recent years. We now have methods in NLP [41,12,40] and computer vision [21,5,20,8,2] that can produce stronger features than those learned on labeled datasets.

In particular, masked autoencoding, whose task is to reconstruct the masked data from the unmasked input (e.g., predicting the masked word in a sentence or masked patch in an image, based on surrounding unmasked context) is the dominant self-supervised learning approach for text understanding [12,26,27,63] and has recently shown great promise in image understanding [2,20] as well. Curiously, for point cloud data, masked autoencoding has not yet been able to produce convincing results [53,65,62]. Self-supervised learning would be extremely beneficial for point cloud data, as obtaining high-quality annotations is both hard and expensive, especially for real-world scans. At the same time, masked autoencoding should also be a good fit for point cloud data, since each point (or group of points) can easily be masked or unmasked.

We hypothesize that the primary reason why masked autoencoding has thus far not worked well for point cloud data is because standard point cloud back-

¹ Code will be publicly available at <https://github.com/haotian-liu/MaskPoint>.

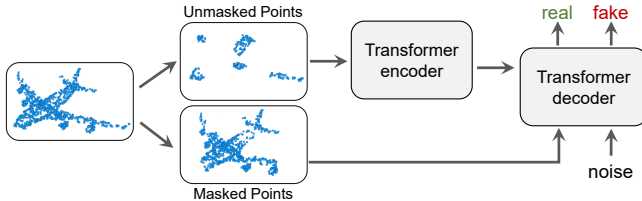


Fig.1: **Main Idea.** We randomly partition the point cloud into masked and unmasked sets. We only feed the visible portion of the point cloud into the encoder. Then, a set of real query points are sampled from the masked points, and a set of fake query points are randomly sampled from 3D space. We train the decoder so that it distinguishes between the real and fake points. After pre-training, we discard the decoder and use the encoder for downstream tasks.

bones are unable to properly handle the distribution mismatch between training and testing data introduced by masking. Specifically, PointNet type backbones [39,37,38] leverage local aggregation layers that operate over local neighborhoods (e.g., k -nearest neighbors) of each point. The extent of the local neighborhoods can change drastically with the introduction of masking, creating a discrepancy between the distribution of local neighborhoods seen on masked training scenes versus unmasked test scenes.

Transformers [52], on the other hand, can perform self-attention (a form of aggregation) on either all or selective portions of the input data. This means that it has the ability to only process the unmasked portions of the scene in the training data, without being impacted by the masked scene portions. This property suggests that Transformers could be an ideal backbone choice for self-supervised masked autoencoding for point clouds.

For image understanding, the state-of-the-art masked autoencoding Transformer approach MAE [20] masks out a large random subset of image patches, applies the Transformer encoder to the *unmasked* patches, and trains a small Transformer decoder that takes in the positional encodings of the *masked* patches to reconstruct their original pixel values. However, this approach cannot be directly applied to point cloud data, because the raw representation of each 3D point is its spatial xyz location. Thus, training the decoder to predict the xyz coordinates of a masked point would be trivial, since its positional encoding would leak the correct answer. In this case, the network would simply take a shortcut and not learn meaningful features.

To address this, we propose a simple binary point classification objective as a new pretext task for point cloud masked autoencoding. We first group points into local neighborhoods, and then mask out a large random subset of those groups. The Transformer encoder takes in the unmasked point groups, and encodes each group through self-attention with the other groups. The Transformer decoder takes in a set of real query points and fake query points, where the real queries are sampled from the *masked* points, while the fake queries are randomly sampled from the full 3D space. We then perform cross attention between the decoder

queries and encoder outputs. Finally, we apply a binary classification head to the decoder’s outputs and require it to distinguish between the real and fake queries. We find this design to be simple yet effective, as it creates a difficult and meaningful pretext task that requires the network to deduce the shape of the object from only a small amount of visible point groups.

Importantly, we find that a much higher masking ratio (e.g., 90%) is required for point cloud data compared to the image domain (75% in [20]). We identify two possible reasons. First, the additional depth dimension in point clouds can help disambiguate object instances (i.e., make the task easier) – with the same masking ratio, a 3D patch in point cloud data is less likely to overlap multiple objects than a 2D patch in images. Second, local point groups often already contain unique categorical information, e.g. a wheel-shaped local point patch can yield a high probability of a car-shaped object. These two observations demand stronger masking ratios in order for the pretext task to be hard enough to make the model learn useful feature representations.

Among existing self-supervised point cloud approaches, Point-BERT [65] is the most related. It trains a discrete Variational AutoEncoder (dVAE) [43] to encode the input point cloud into discrete point token representations, and performs BERT-style pretraining over them. To aid training, it uses point patch mixing augmentation, together with an auxiliary MoCo [21] loss. However, the dependency on a pretrained dVAE together with other auxiliary techniques, creates a significant computational overhead in pretraining – our experiments show that its pre-training is significantly slower (e.g., $4.1\times$ on ScanNet [9]) than ours even without taking into account the training time of the dVAE module. The large speedup is also due to our design of having a high masking rate and the Transformer encoder processing only the unmasked points.

In sum, our main contributions are: (1) A novel masked point classification Transformer, *MaskPoint*, for self-supervised learning on point clouds. (2) Our approach is simple and effective, achieving state-of-the-art performance on a variety of downstream tasks, including object classification on ModelNet40 [58] / ScanObjectNN [51], part segmentation on ShapeNetPart [64], object detection on ScanNet [9], and few-shot object classification on ModelNet40 [58]. (3) Notably, for the first time, we show that a standard Transformer architecture can outperform sophisticatedly designed point cloud backbones.

2 Related Work

Transformers. Transformers were first proposed to model long-term dependencies in sequential data [52], and have achieved great success in natural language processing [52,12,41]. More recently, they have also shown promising performance on various computer vision tasks, including image classification [14,50,46], object detection [3], semantic segmentation [54], image generation [25], and multi-modal learning [40]. There have also been attempts to adopt transformers to 3D point cloud data. PCT [19] and Point Transformer [71] propose new attention mechanisms for point cloud feature aggregation. 3DETR [34] uses Transformer blocks and the parallel decoding strategy from DETR [3] for 3D object detec-

tion. However, it is still hard to get promising performance using the standard Transformer. For example, in 3D object detection, there is a large performance gap between 3DETR [34] and state-of-the-art point based [69] and convolution based [10] methods. In this paper, we propose a novel masked autoencoding Transformer for self-supervised learning on point clouds.

Self-supervised Learning. Self-supervised learning (SSL) aims to learn meaningful representations from the data itself, to better serve downstream tasks. Traditional methods typically rely on pretext tasks, such as image rotation prediction [18], image colorization [67], and solving jigsaw puzzles [35]. Recent methods based on contrastive learning (e.g., MoCo [21], SimCLR [5], SimSiam [8]) have achieved great success in the image domain, sometimes producing even better downstream performance compared to supervised pretraining on ImageNet [11].

Self-supervised learning has also begun to be explored for point cloud data. Pretext methods include deformation reconstruction [1], geometric structure prediction [47], and orientation estimation [36]. Contrastive learning approaches include PointContrast [60], which learns corresponding points from different camera views, and DepthContrast [68], which learns representations by comparing transformations of a 3D point cloud/voxel. OcCo [53] learns an autoencoder to reconstruct the scene from the occluded input. However, due to the sampling variance of the underlying 3D shapes, explicitly reconstructing the original point cloud will inevitably capture such variance. In this paper, we explore a simple but effective discriminative classification pretext task to learn representations that are robust to the sampling variance.

Mask based Pretraining. Masking out content has been used in various ways to improve model robustness including as a regularizer [45,17], data augmentation [13,44,72], and self-supervised learning [7,12,20]. For self-supervised learning, the key idea is to train the model to predict the masked content based on its surrounding context. The most successful approaches are built upon the Transformer [52], due in part to its token-based representation and ability to model long-range dependencies.

In masked language modeling, BERT [12] and its variants [26,27] achieve state-of-the-art performance across nearly all NLP downstream tasks by predicting masked tokens during pretraining. Masked image modeling works [2,20] adopt a similar idea for image pretraining. BEiT [2] maps image patches into discrete tokens, then masks a small portion of patches, and feeds the remaining visible patches into the Transformer to reconstruct the tokens of the masked patches. Instead of reconstructing tokens, the recent Masked AutoEncoder (MAE) [20] reconstructs the masked patches at the pixel level, and with a much higher mask ratio of $\geq 70\%$. Following works try to predict high-level visual features such as HoG [56], or improve the representation capability of the encoder by aligning the feature from both visible patches and masked patches [7]. To our knowledge, the only self-supervised mask modeling Transformer approach for point clouds is Point-BERT [65], which adopts a similar idea as BEiT [2]. However, to obtain satisfactory performance, it requires a pretrained dVAE and other auxiliary techniques (e.g., a momentum encoder [21]), which slow down training. We

propose a simple and effective masked autoencoding Transformer approach for point clouds, which largely accelerates training ($4.1\times$ faster than Point-BERT) while achieving state-of-the-art performance for various downstream tasks.

3 Approach

The goal is to learn semantic feature representations without human supervision that can perform well on downstream point cloud recognition tasks. We motivate our self-supervised learning design with a qualitative example. Fig. 1 “Unmasked Points” shows a point cloud with a large portion (90%) of its points masked out. Still, based on our prior semantic understanding of the world, we as humans are able to say a number of things about it: (1) it might be an airplane; (2) if so, it should consist of the head, body, tail, and wing; and even (3) roughly where these parts should be present. In other words, because we already know what airplanes are, we can recover the missing information from the small visible subset of the point cloud. In a similar way, training a model to recover information about the masked portion of the point cloud given the visible portion could force the model to learn object semantics.

However, even as humans, it can be difficult or impossible to precisely reconstruct all missing points, since there are several ambiguous factors; e.g., the precise thickness of the wings or the precise length of the airplane. If we are instead given a sampled 3D point in space, and are asked to answer whether it likely belongs to the object or not, we would be more confident in our answer. This discriminative point classification task is much less ambiguous than the reconstruction task, yet still requires a deep understanding of object semantics in order to deduce the masked points from the small number of visible points.

3.1 Masked Point Discrimination

Our approach works as follows. We randomly partition each input point cloud $\mathcal{P} \in \mathbb{R}^{N \times 3}$ into two groups: masked \mathcal{M} and unmasked \mathcal{U} . We use the Transformer encoder to model the correlation between the sparsely-distributed unmasked tokens \mathcal{U} via self-attention. Ideally, the resulting encoded latent representation tokens \mathcal{L} should not only model the relationship between the unmasked points \mathcal{U} , but also recover the latent distribution of masked points \mathcal{M} , so as to perform well on the pretraining task. We next sample a set of *real* query points \mathbf{Q}_{real} and a set of *fake* query points \mathbf{Q}_{fake} . The real query points are sampled from the masked point set \mathcal{M} , while the fake query points are randomly sampled from the full 3D space. We then perform cross attention between each decoder query $\mathbf{q} \in \{\mathbf{Q}_{real}, \mathbf{Q}_{fake}\}$ and the encoder outputs: $\text{CA}(\mathbf{q}, \mathcal{L})$, to model the relationship between the masked query point and the unmasked points. Finally, we apply a binary classification head to the decoder’s outputs and require it to distinguish between the real and fake queries.

We show in our experiments that our approach is both simple and effective, as it creates a pretext task that is difficult and meaningful enough for the model to learn rich semantic point cloud representations.

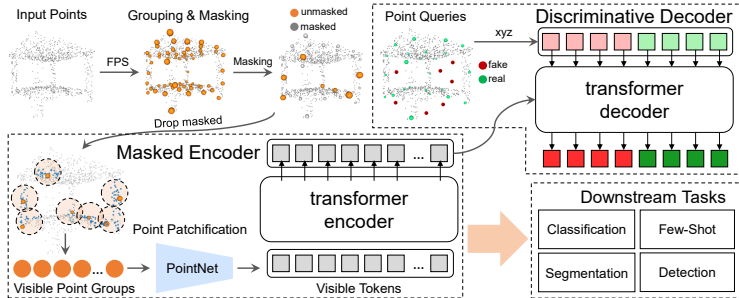


Fig. 2: **MaskPoint architecture.** We first uniformly sample point groups from the point cloud, and partition them to masked and unmasked. We patchify the visible point groups to token embeddings with PointNet and feed these visible tokens into the encoder. Then, a set of real query points are sampled from the masked points, and a set of fake query points are randomly sampled from 3D space. We train the decoder so that it distinguishes between the real and fake points. After pre-training, we discard the decoder and use the encoder for downstream tasks. See Sec. 3.1 for details.

Discarding Ambiguous Points. Since we sample fake query points uniformly at random over the entire space, there will be some points that fall close to the object’s surface. Such points can cause training difficulties since their target label is ‘fake’ even though they are on the object. In preliminary experiments, we find that such ambiguous points can lead to vanishing gradients in the early stages of training. Thus, to stabilize training, we simply remove all fake points $\hat{\mathbf{p}} \in \mathbf{Q}_{fake}$ whose euclidean distance is less than γ to any object (masked or unmasked) point $\mathbf{p}_i \in \mathcal{P}$: $\min_i \|\hat{\mathbf{p}} - \mathbf{p}_i\|_2 < \gamma$. To address the size variance of the input point cloud, γ is dynamically selected per point cloud \mathcal{P} : $\hat{\mathcal{P}} = \text{FPS}(\mathcal{P})$, $\gamma = \min_{j \neq i} \|\hat{\mathcal{P}}_i - \hat{\mathcal{P}}_j\|_2$.

3D Point Patchification. Feeding every single point into the Transformer encoder can yield an unacceptable cost due to the quadratic complexity of self-attention operators. Following [65, 14], we adopt a patch embedding strategy that converts input point clouds into 3D point patches.

Given the input point cloud $\mathcal{P} \in \mathbb{R}^{N \times 3}$, S points $\{\mathbf{p}_i\}_{i=1}^S$ are sampled as patch centers using farthest point sampling [39]. We then gather the k nearest neighbors for each patch center to generate a set of 3D point patches $\{\mathbf{g}_i\}_{i=1}^S$. A PointNet [38] is then applied to encode each 3D point patch $\mathbf{g}_i \in \mathbb{R}^{k \times 3}$ to a feature embedding $\mathbf{f}_i \in \mathbb{R}^d$. In this way, we obtain S tokens and their corresponding features $\{\mathbf{f}_i\}_{i=1}^S$ and center coordinates $\{\mathbf{p}_i\}_{i=1}^S$.

Note that our patchification strategy can generate overlapping patches (e.g., if two point centers are sampled to be close to each other). Although a prior study [59] showed that such overlapping patches can stabilize training in vision transformers, in our case, they could undesirably leak information that would allow the network to take a shortcut solution; e.g., a portion of a masked patch being part of a neighboring unmasked patch. In practice, we set a very high

masking ratio (e.g., 90%), and consequently such overlap rarely happens and does not hinder training.

Transformer Architecture. Our network architecture is shown in Fig. 2. We adopt the standard Transformer encoder [52] as the encoding backbone, where each Transformer encoder block consists of a multi-head self-attention (MSA) layer and a feed forward network (FFN). As noted in Section 3.1, we construct patch-wise features $\{\mathbf{f}_i\}_{i=1}^M$ from the input point cloud $\mathcal{P} \in \mathbb{R}^{N \times 3}$. Following [65], we apply the MLP positional embedding $\{\mathbf{pos}_i\}_{i=1}^M$ to the patch features $\{\mathbf{f}_i\}_{i=1}^M$. Then, the class token $\mathbf{E}[s]$, which will be used for downstream classification tasks, is stacked to the top of the patch features $\{\mathbf{f}_i\}_{i=1}^M$; i.e., the input to the Transformer encoder is $I_0 = \{\mathbf{E}[s], \mathbf{f}_1 + \mathbf{pos}_1, \mathbf{f}_2 + \mathbf{pos}_2, \dots, \mathbf{f}_M + \mathbf{pos}_M\}$. After n Transformer blocks, we get the feature embedding for each group $I_n = \{\mathbf{E}^n[s], \mathbf{f}_1^n, \mathbf{f}_2^n, \dots, \mathbf{f}_M^n\}$.

During the decoding stage, N_q *real* query points \mathbf{Q}_{real} and N_q *fake* query points \mathbf{Q}_{fake} are sampled. We pass the encoder output I_n and its positional embedding $\{\mathbf{pos}_i\}_{i=1}^M$, \mathbf{Q}_{real} , \mathbf{Q}_{fake} and their positional embedding $\{\mathbf{pos}_i^{\mathbf{Q}}\}_{i=1}^{2N}$ into a one-layer Transformer decoder. Cross attention is only performed between the queries and encoder keys/values, but not between different queries. Finally, the decoder output goes through an MLP classification head, which is trained with the binary focal loss [29], since there can be a large imbalance between positive and negative samples.

For downstream tasks, the point patchification module and Transformer encoder will be used with their pretrained weights as initialization.

An Information Theoretic Perspective. Here, we provide an information theoretic perspective to our self-supervised learning objective, using mutual information. The mutual information between random variables X and Y , $I(X; Y)$, measures the amount of information that can be gained about random variable X from the knowledge about the other random variable Y .

Ideally, we would like the model to learn a rich feature representation of the point cloud: the latent representation \mathcal{L} from our encoder \mathcal{E} should contain enough information to recover the original point cloud \mathcal{P} , i.e., we would like to maximize the mutual information $I(\mathcal{P}; \mathcal{L})$. However, directly estimating $I(\mathcal{P}; \mathcal{L})$ is hard since we need to know the exact probability distribution of $P(\mathcal{P}|\mathcal{L})$. Following [6], we instead use auxiliary distribution Q to approximate it:

$$\begin{aligned}
 I(\mathcal{P}; \mathcal{L}) &= -H(\mathcal{P}|\mathcal{L}) + H(\mathcal{P}) \\
 &= \mathbb{E}_{x \sim \mathcal{L}}[\mathbb{E}_{p' \sim P(\mathcal{P}|\mathcal{L})}[\log P(p'|x)]] + H(\mathcal{P}) \\
 &= \mathbb{E}_{x \sim \mathcal{L}}[\underbrace{D_{\text{KL}}(P(\cdot|x) \| Q(\cdot|x))}_{\geq 0} + \mathbb{E}_{p' \sim P(\mathcal{P}|\mathcal{L})}[\log Q(p'|x)]] + H(\mathcal{P}) \quad (1) \\
 &\geq \mathbb{E}_{x \sim \mathcal{L}}[\mathbb{E}_{p' \sim P(\mathcal{P}|\mathcal{L})}[\log Q(p'|x)]] + H(\mathcal{P})
 \end{aligned}$$

Lemma 3.1. For random variables X, Y and function $f(x, y)$ under suitable regularity conditions: $\mathbb{E}_{x \sim X, y \sim Y|x}[f(x, y)] = \mathbb{E}_{x \sim X, y \sim Y|x, x' \sim X|y}[f(x', y)]$.

Therefore, we can define a variational lower bound, $L_I(Q, \mathcal{L})$, of the mutual information, $I(\mathcal{P}; \mathcal{L})$:

$$\begin{aligned} L_I(Q, \mathcal{L}) &= \mathbb{E}_{p \sim P(\mathcal{P}), x \sim \mathcal{L}} [\log Q(p|x)] + H(\mathcal{P}) \\ &= \mathbb{E}_{x \sim \mathcal{L}} [\mathbb{E}_{p' \sim P(\mathcal{P}|\mathcal{L})} [\log Q(p'|x)]] + H(\mathcal{P}) \\ &\leq I(\mathcal{P}; \mathcal{L}) \end{aligned} \quad (2)$$

Therefore, we have:

$$\max I(\mathcal{P}; \mathcal{L}) \iff \max L_I(Q, \mathcal{L}) \iff \max \mathbb{E}_{x \sim \mathcal{L}} [\mathbb{E}_{p' \sim P(\mathcal{P}|\mathcal{L})} [\log Q(p'|x)]] \quad (3)$$

Previous works use the Chamfer distance to approximate such auxiliary function Q , but it has the disadvantage of being sensitive to point sampling variance (discussed in detail in Sec. 3.2). Thus, we instead represent the point cloud distribution with occupancy values within the tightest 3D bounding box of the point cloud: $\mathcal{B} \in \{x, y, z, o\}^L$, where $(x, y, z) \in \mathbb{R}^3$, $o \in \{0, 1\}$, and L is the number of densely sampled points. We let the output of Q denote the continuous distribution of the occupancy value \hat{o} , where $Q(\cdot) \in [0, 1]$. In our implementation, as discussed in Sec. 3.1, we construct a set of real query points and fake query points, assign them with the corresponding occupancy labels, and optimize the probability outputs from the model with a binary classification objective.

3.2 Why not reconstruction, as in MAE?

In this section, we delve into the details on why a reconstruction objective (i.e., reconstructing the original point cloud from the unmasked points) as used in the related Masked AutoEncoder (MAE) [20] approach for images would not work for our point cloud setting.

First, in MAE, the self-supervised learning task is to reconstruct the masked patches, based on the input image’s unmasked (visible) patches. Specifically, given the 2D spatial position for each masked image patch query, the objective is to generate its RGB pixel values. In our case, the analogue would be to generate the spatial xyz values for a masked 3D point patch query – which would be trivial for the model since the query already contains the corresponding spatial information. Such a trivial solution will result in perfect zero-loss, and prevent the model from learning meaningful feature representations.

Another issue with the reconstruction objective for point clouds is that there will be point sampling variance. Specifically, the true 3D shape of the object will be a continuous surface, but a point cloud will a discrete sampling of it. Suppose we sample two such point clouds, and denote the first set as the “ground truth” target, and the second set as the prediction of a model. Although both sets reflect the same geometric shape of the object, the Chamfer distance (which can be used to measure the shape difference between the two point sets) between them is non-zero (i.e., there would be a loss). Thus, minimizing the Chamfer distance would force the model to generate predictions that exactly match the first set. And since the first set is just one sampling of the true underlying distribution, this can be an unnecessarily difficult optimization problem that leads to suboptimal model performance.

4 Experiments

We evaluate the pre-trained representation learned by the proposed model on a variety of point cloud downstream tasks, including object classification, part segmentation, object detection, and few-shot object classification. We also visualize the reconstruction results from masked point clouds, to qualitatively study the effect of our pretraining. Finally, we perform ablation studies on masking strategies and decoder designs.

Pretraining Datasets. (1) ShapeNet [4] has 50,000 unique 3D models from 55 common object categories, and is used as our pre-training dataset for object classification, part segmentation, and few-shot classification. For ShapeNet pre-training, we sample 1024 points from each 3D model as the inputs. We follow [65] to sample 64 point groups, each containing 32 points.

(2) We also use single-view depth map videos from the popular ScanNet [9] dataset, which contains around 2.5 million RGBD scans. We do not use its RGB information in this paper. We adopt similar pre-processing steps as DepthContrast [68], but we generate a smaller subset of the dataset than in [68], which we call ‘ScanNet-Medium’, to accelerate pretraining. ScanNet-Medium is generated by sampling every 10-th frame from ScanNet, resulting in ~ 25 k samples. We use ScanNet-Medium (only geometry information) as the pre-training dataset for 3D object detection. For pretraining, we sample 20k points from each 3D scene scan as the input. We follow [34] to sample 2048 groups, each containing 64 points.

Transformer Encoder. We construct a 12-layer standard Transformer encoder, named PointViT, for point cloud understanding. Following Point-BERT [65], we set the hidden dimension of each encoder block to 384, number of heads to 6, FFN expansion ratio to 4, and drop rate of stochastic depth [23] to 0.1.

Transformer Decoder. We use a single-layer Transformer decoder for pre-training. The configuration of the attention block is identical to the encoder.

Training Details. Following [65], we pretrain with the AdamW [33] optimizer with a weight decay of 0.05, and a learning rate of 5×10^{-4} decayed with the cosine schedule. The model is trained for 300 epochs with a batch size of 128, with random scaling and translation data augmentation. Only for ModelNet40 experiments, we also pretrain with a MoCo loss [21], following [65]. However, we do not use it for any other datasets, unlike [65], which always uses it. For finetuning and additional training details, please see supp.

4.1 3D Object Classification

Datasets. We compare the performance of object classification on two datasets: the synthetic ModelNet40 [58], and real-world ScanObjectNN [51]. ModelNet40 [58] consists of 12,311 CAD models from 40 classes. We follow the official data splitting scheme in [58]. We evaluate the overall accuracy (OA) over all test samples. ScanObjectNN [51] is a more challenging point cloud benchmark that

Method	SSL #point	OA
PointNet [38]	1k	89.2
PointNet++ [39]	1k	90.7
PointCNN [28]	1k	92.2
SpiderCNN [61]	1k	92.4
PointWeb [70]	1k	92.3
PointConv [57]	1k	92.5
DGCNN [55]	1k	92.9
KPConv [48]	1k	92.9
DensePoint [30]	1k	93.2
PosPool [32]	5k	93.2
RSCNN [31]	5k	93.6
[T] Point Trans. [16]	1k	92.8
[T] Point Trans. [71]	–	93.7
[T] PCT [19]	1k	93.2
[ST] PointViT	1k	91.4
[ST] PointViT-OcCo [53]	✓ 1k	92.1
[ST] Point-BERT [65]	✓ 1k	93.2
[ST] MaskPoint (Ours)	✓ 1k	93.8

Table 1: **Shape Classification on ModelNet40 [58]**. With a standard Transformer backbone, our approach significantly outperforms training-from-scratch baselines and SOTA pretraining methods. It even outperforms Point-Transformer [71], which uses an attention operator specifically designed for point clouds. *SSL: Self-supervised pretraining. [T]: Transformer-based networks with special designs for point clouds. [ST]: Standard Transformer network.

consists of 2902 unique objects in 15 categories collected from noisy real-world scans. It has three splits: OBJ (object only), BG (with background), PB (with background and manually added perturbations). We evaluate the overall accuracy (OA) over all test samples on all three splits.

ModelNet40 Results. Table 1 shows ModelNet40 [58] results. With 1k points, our approach achieves a significant 2.4% OA improvement compared to training from scratch (PointViT). It also brings a 1.7% gain over OcCo [53] pretraining, and 0.6% gain over Point-BERT [65] pretraining. The significant improvement over the baselines indicates the effectiveness of our pre-training method. Notably, for the first time, with 1k points, a standard vision transformer architecture produces competitive performance compared to sophisticatedly designed attention operators from PointTransformer [71] (93.8% vs 93.7%).

Method	OA		
	OBJ	BG	PB
PointNet [38]	79.2	73.3	68.0
PointNet++ [39]	84.3	82.3	77.9
PointCNN [28]	85.5	86.1	78.5
SpiderCNN [61]	79.5	77.1	73.7
DGCNN [55]	86.2	82.8	78.1
BGA-DGCNN [51]	–	–	79.7
BGA-PN++ [51]	–	–	80.2
PointViT	80.6	79.9	77.2
PointViT-OcCo [53]	85.5	84.9	78.8
Point-BERT [65]	88.1	87.4	83.1
MaskPoint (Ours)	89.3	88.1	84.3

Table 2: **Shape Classification on ScanObjectNN [51]**. OBJ: object-only; BG: with background; PB: BG with manual perturbation.

Method	mIoU	
	cat.	ins.
PointNet [38]	80.4	83.7
PointNet++ [39]	81.9	85.1
DGCNN [55]	82.3	85.2
PointViT	83.4	85.1
PointViT-OcCo [53]	83.4	85.1
Point-BERT [65]	84.1	85.6
MaskPoint (Ours)	84.4	86.0

Table 3: **Part Segmentation on ShapeNetPart [64]**. Our method also works well on dense prediction tasks like segmentation.

Method	5-way		10-way	
	10-shot	20-shot	10-shot	20-shot
DGCNN [53]	91.8 \pm 3.7	93.4 \pm 3.2	86.3 \pm 6.2	90.9 \pm 5.1
DGCNN-OcCo [53]	91.9 \pm 3.3	93.9 \pm 3.2	86.4 \pm 5.4	91.3 \pm 4.6
PointViT	87.8 \pm 5.3	93.3 \pm 4.3	84.6 \pm 5.5	89.4 \pm 6.3
PointViT-OcCo [53]	94.0 \pm 3.6	95.9 \pm 2.3	89.4 \pm 5.1	92.4 \pm 4.6
Point-BERT [65]	94.6 \pm 3.1	96.3 \pm 2.7	91.0 \pm 5.4	92.7 \pm 5.1
MaskPoint (Ours)	95.0 \pm 3.7	97.2 \pm 1.7	91.4 \pm 4.0	93.4 \pm 3.5

Table 4: **Few-shot classification on ModelNet40 [58].**

ScanObjectNN Results. We next conduct experiments using the real-world scan dataset ScanObjectNN [51]. Table 2 shows the results. Our approach achieves SOTA performance on all three splits. On the hardest PB split, our approach achieves a large 7.1% OA improvement compared to training from scratch (Point-ViT). It achieves a 5.5% gain over OcCo [53] pretraining, and a 1.2% gain over Point-BERT [65] pretraining. The large improvement over the baselines highlights the transferability of our model’s self-supervised representation, as there is a significant domain gap between the clean synthetic ShapeNet [4] dataset used for pretraining and the noisy real-world ScanObjectNN [51] dataset.

We believe the performance gain over OcCo [53] and Point-BERT [65] is mainly because of our discriminative pretext task. OcCo suffers from the sampling variance issue (Sec. 3.2) as it uses a reconstruction-based objective in pretraining. Compared to Point-BERT, we do not use the point patch mixing technique, which mixes two different point clouds. This could introduce unnecessary noise and domain shifts to pretraining and harm downstream performance.

4.2 3D Part Segmentation

Dataset. ShapeNetPart [64] consists of 16,880 models from 16 shape categories, with 14,006 models for training and 2,874 for testing. It contains 50 different parts in total, and the number of parts for each category is between 2 and 6. We use the sampled point sets produced by [39] for a fair comparison with prior work. We report per-category mean IoU (cat. mIoU) and mean IoU averaged over all test instances (ins. mIoU).

Results. Table 3 shows the results (per-category IoU is in supp). Our approach outperforms the training from the scratch (PointViT) and OcCo-pretraining baselines by 1.0%/0.9% in cat./ins. mIoU. It also produces a 0.3%/0.4% gain compared to Point-BERT. Thanks to our dense discriminative pretraining objective, in which we densely classify points over the 3D space, we are able to obtain good performance when scaling to dense prediction tasks.

4.3 Few-shot Classification

We conduct few-shot classification experiments on ModelNet40 [58], following the settings in [65]. The standard experiment setting is the “ K -way N -shot”

configuration, where K classes are first randomly selected, and then $N + 20$ objects are sampled for each class. We train the model on $K \times N$ samples (support set), and evaluate on the remaining $K \times 20$ samples (query set). We compare our approach with OcCo and Point-BERT, which are the current state-of-the-art.

We perform experiments with 4 settings, where for each setting, we run the train/evaluation on 10 different sampled splits, and report the mean and std over the 10 runs. Table 4 shows the results. Our approach achieves the best performance for all settings. It demonstrate an absolute gain of 7.2%/3.8%/4.6%/2.5% over the PointViT training from the scratch baseline. When comparing to pre-training baselines, it outperforms OcCo by 1.0%/1.3%/1.5%/1.0%, and outperforms Point-BERT by 0.4%/0.9%/0.4%/0.7%. It also clearly outperforms the DGCNN baselines. Our state-of-the-art performance on few-shot classification further demonstrates the effectiveness of our pretraining approach.

4.4 3D Object Detection

Our most closely related work, Point-BERT [65] showed experiments only on object-level classification and segmentation tasks. In this paper, we evaluate a model’s pretrained representation on a more challenging scene-level downstream task: 3D object detection on ScanNetV2 [9], which consists of real-world richly-annotated 3D reconstructions of indoor scenes. It comprises 1201 training scenes, 312 validation scenes and 100 hidden test scenes. Axis-aligned bounding box labels are provided for 18 object categories. For this experiment, we adopt 3DETR [34] as the downstream model for both our method and Point-BERT. 3DETR is an end-to-end transformer-based 3D object detection pipeline. During finetuning, the input point cloud is first downsampled to 2048 points via a VoteNet-style Set Aggregation (SA) layer [37,39], which then goes through 3-layer self-attention blocks. The decoder is composed of 8-layer cross-attention blocks. For a fair comparison with the 3DETR train-from-scratch baseline, we strictly follow its architecture of SA layer and encoder during pretraining, whose weights are transferred during finetuning. Our pretraining dataset is ScanNet Medium, as described in Sec. 4.

Table 5 shows that our method surpasses the 3DETR train-from-scratch baseline by a large margin (+1.3mAP@0.25 and +2.7mAP@0.50). Interestingly, Point-BERT brings nearly no improvement compared to training from scratch. The low mask rate and discrete tokens learned from dVAE may impede Point-BERT from learning meaningful representations for detection. Also, the 3DETR paper [34] found that increasing the number of encoding layers in 3DETR brings only a small benefit to its detection performance. Here we increase the number of layers from 3 to 12, which leads to a large performance improvement (+2.1mAP@0.25 and +4.1mAP@0.50) for our approach compared to training from scratch. This result demonstrates that by pre-training on a large unlabeled dataset, we can afford to increase the model’s encoder capacity to learn richer representations. Finally, note that we also include VoteNet based methods at the top of Table 5 as a reference, but the numbers are not directly comparable as they are using a different (non Transformer-based) detector.

Methods	SSL	Pretrained Input	\mathbf{AP}_{25}	\mathbf{AP}_{50}
VoteNet [37]		-	58.6	33.5
STRL [24]	✓	Geo	59.5	38.4
Implicit Autoencoder [62]	✓	Geo	61.5	39.8
RandomRooms [42]	✓	Geo	61.3	36.2
PointContrast [60]	✓	Geo	59.2	38.0
DepthContrast [68]	✓	Geo	61.3	-
DepthContrast [68]	✓	Geo + RGB	64.0	42.9
3DETR [34]		-	62.1	37.9
Point-BERT [65]	✓	Geo	61.0	38.3
MaskPoint (Ours)	✓	Geo	63.4	40.6
MaskPoint (Ours, 12 Enc)	✓	Geo	64.2	42.0

Table 5: 3D object detection results on ScanNet validation set. The backbone of our pretraining model and Point-BERT [65] is 3DETR [34]. All other methods use VoteNet [38] as the finetuning backbone. Only geometry information is fed into the downstream task. “Input” column denotes input type for the pretraining stage. “Geo” denotes geometry information. Note that DepthContrast (Geo + RGB) model uses a heavier backbone (PointNet 3x) for downstream tasks.

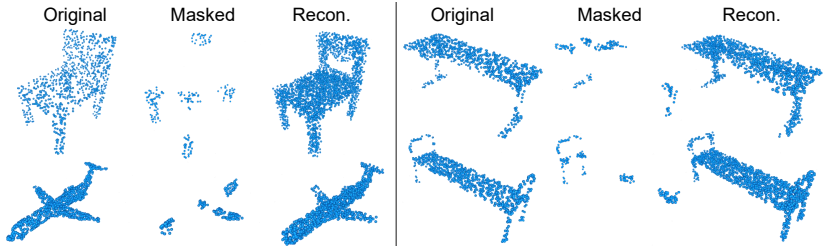


Fig. 3: **Reconstruction results.** By reformulating reconstruction as a discriminative occupancy classification task, we achieve a similar learning objective to generative reconstruction while being robust to point sampling variance. Even with a high 90% mask ratio, our approach recovers the overall shape of the original point cloud, without overfitting.

4.5 Qualitative Reconstructions

Although our model is not trained with the reconstruction objective, we can still reconstruct the point cloud with our decoder by classifying the densely sampled points from the point cloud’s full 3D bounding box space: $\mathcal{P}_{rec} = \{x | \mathcal{D}(x | \mathcal{L}) = 1\}$. Fig. 3 shows the reconstruction results with masking ratio of 90%. Even with such a large masking rate, our model is able to reconstruct the overall shape of the original point cloud, without overfitting. Fig. 3 (bottom right) shows how our model uses its learned prior to perform reasonable reconstructions: a shoe shelf where the stand on the right is masked. We hypothesize that it first recognizes it as a shoe shelf (without any labels), and learns that shoe shelves are usually symmetric from left-to-right, to reconstruct the right part.

ratio	OA	ratio	OA	# queries	OA	# dec.	OA
0.25	83.2	0.25	82.4	64	83.7	1	84.3
0.50	83.7	0.50	83.8	256	84.3	3	83.7
0.75	84.1	0.75	83.7	1024	83.9	6	83.9
0.90	84.3	0.90	84.1				

(a) **Mask rate** (random) (b) **Mask rate** (block) (c) **# dec. queries** (d) **# dec. layers**

Table 6: **Ablations** on ScanObjectNN [51] (PB split). Our findings are: a larger masking ratio generally yields better performance; random masking is slightly better than block masking; 256-query provides a good balance between information and noise; a thin decoder ensures rich feature representation in the encoder and benefits downstream performance.

4.6 Ablation Studies

Masking Strategy. We show the influence of different masking strategies in Table 6a, 6b. First, we observe that a higher masking ratio generally yields better performance, regardless of sampling type. This matches our intuition that a higher masking ratio creates a harder and more meaningful pretraining task. Further, with a higher masking ratio, random masking is slightly better than block masking. Therefore, we use a high mask rate of 90% with random masking.

Pretraining Decoder Design. We study the design of the pretraining decoder in Table 6c, 6d, by varying the number of decoder queries and layers. The number of decoder queries influence the balance between the classification of the real points and fake points. We find that 256-query is the sweet spot, where more queries could introduce too much noise, and fewer queries could result in insufficient training information.

The modeling power of the decoder affects the effectiveness of the pretraining: ideally, we want the encoder to only encode the features, while the decoder only projects the features to the pre-training objective. Any imbalance in either way can harm the model’s performance on downstream tasks. We find that a single-layer decoder is sufficient for performing our proposed point discrimination task, and having more decoder layers harms the model’s performance.

5 Conclusion

We proposed a discriminative masked point cloud pretraining framework, which facilitates a variety of downstream tasks while significantly reducing the pre-training time compared to the prior Transformer-based state-of-the-art method. We adopted occupancy values to represent the point cloud, forming a simpler yet effective binary pretraining objective function. Extensive experiments on 3D shape classification, detection, and segmentation demonstrated the strong performance of our approach. Currently, we randomly mask local point groups to partition the point cloud into masked and unmasked sets. It could be interesting to explore ways to instead learn how to mask the points. We hope our research can raise more attention to mask based self-supervised learning on point clouds.

Supplementary Material

A More results

ShapeNetPart In Table 7, we compare the categorical mIoU on ShapeNet-Part with other methods. With a PointViT backbone, we get the highest class mIoU at 84.4% and the highest instance mIoU at 86.0%, outperforming previous self-supervised learning approaches (OcCo [53] and Point-BERT [65]). It also outperforms standard train-from-scratch point cloud backbones like PointNet++ [39] and DGCNN [55]. For all categories, our method either has the highest accuracy or is among the best. Thanks to our dense discriminative pre-training objective, in which we densely classify points over the 3D space, we are able to obtain good performance when scaling to dense prediction tasks like part segmentation.

Methods	cls.	ins.	aero	bag	cap	car	chair	earp.	guit.	knif.	lamp	lapt.	mot.	mug	pist.	rock.	skt.	table
PointNet [38]	80.4	83.7	83.4	78.7	82.5	74.9	89.6	73.0	91.5	85.9	80.8	95.3	65.2	93.0	81.2	57.9	72.8	80.6
PN++ [39]	81.9	85.1	82.4	79.0	87.7	77.3	90.8	71.8	91.0	85.9	83.7	95.3	71.6	94.1	81.3	58.7	76.4	82.6
DGCNN [55]	82.3	85.2	84.0	83.4	86.7	77.8	90.6	74.7	91.2	87.5	82.8	95.7	66.3	94.9	81.1	63.5	74.5	82.6
PointViT	83.4	85.1	82.9	<u>85.4</u>	87.7	78.8	90.5	<u>80.8</u>	91.1	87.7	<u>85.3</u>	95.6	73.9	<u>94.9</u>	83.5	61.2	74.9	80.6
OcCo [53]	83.4	85.1	83.3	85.2	88.3	<u>79.9</u>	90.7	74.1	91.9	87.6	84.7	95.4	75.5	94.4	84.1	63.1	75.7	80.8
PN-BERT [65]	<u>84.1</u>	<u>85.6</u>	84.3	84.8	88.0	79.8	<u>91.0</u>	81.7	91.6	87.9	85.2	95.6	<u>75.6</u>	94.7	<u>84.3</u>	<u>63.4</u>	<u>76.3</u>	<u>81.5</u>
MaskPoint (Ours)	84.4	86.0	<u>84.2</u>	85.6	<u>88.1</u>	80.3	91.2	79.5	91.9	<u>87.8</u>	86.2	95.3	76.9	95.0	85.3	64.4	76.9	81.8

Table 7: Part segmentation results on ShapeNetPart [64]. **Bold** and underline numbers denote best and second best performance, respectively.

ScanNet Table 8 reports per-class average precision on 18 classes of ScanNetV2 with a 0.25 box IoU threshold. Relying on purely geometric information, our method exceeds 3DETR [34] in detecting objects like curtain, garbagebin, table, desk, *etc*, where geometry is a strong cue for recognition. These results indicate that our mask based discriminative pretraining framework is effective in learning strong geometric representations. More importantly, our model outperforms 3DETR on classes where it has relatively low AP, *e.g.*, picture, door, curtain, refrigerator, *etc*, which demonstrates the usefulness of pretraining: with the pretrained knowledge relevant to those hard classes, the model is able to make more accurate predictions than the training from the scratch baseline.

Model	AP ₂₅	cab.	bed	cha.	sofa	tab.	door	win.	boo.	pic.	cou.	desk	cur.	ref.	sho.	toi.	sink	bat.	gar.
3DETR [34]	62.2	50.2	87.0	86.0	87.1	61.6	46.6	40.1	54.5	9.1	62.8	69.5	48.4	50.9	68.4	97.9	67.6	85.9	45.8
Ours	63.4	51.8	82.5	85.9	86.8	69.8	50.9	36.9	47.3	10.7	59.6	76.3	65.9	55.6	66.4	99.1	61.5	83.7	49.8
Ours (12×)	64.2	49.5	81.0	87.2	86.3	65.2	51.3	42.6	56.7	16.2	56.8	73.8	59.6	56.0	77.0	97.8	66.6	85.0	47.7

Table 8: 3D object detection scores per category on the ScanNetV2 dataset, evaluated with bbox mIoU 0.25. Ours (12×): 12 encoder blocks.

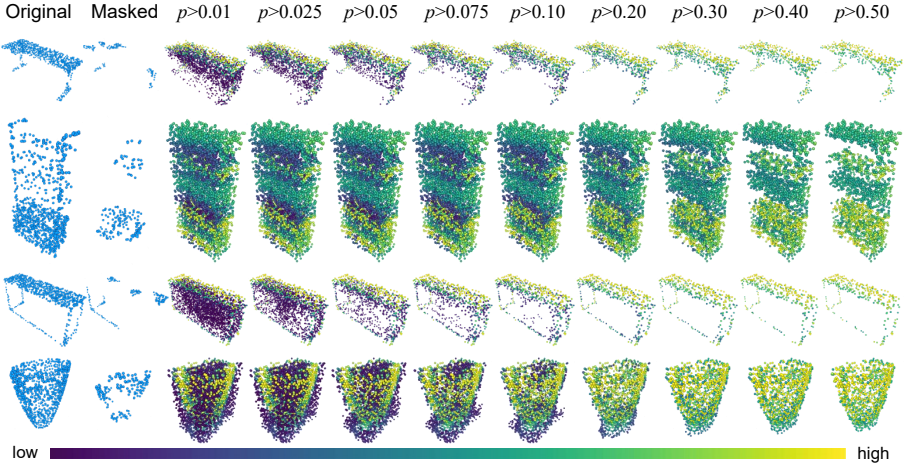


Fig. 4: **Reconstruction results.** We densely perform the discriminative occupancy classification task in 3D space, and visualize the predicted occupancy probability. By varying the confidence threshold \hat{p} , we show that our model is able to predict a continuous probability distribution of the occupancy function.

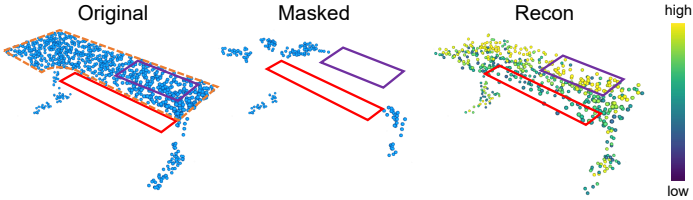


Fig. 5: **A closer look at occupancy distribution.** Although there are no points present in both **red** and **purple** regions of the masked point cloud, the reconstructed probability distribution correctly reflects that of the original point cloud: a lower occupancy in **red** region, and a higher occupancy in **purple** region.

More reconstruction visualizations We densely perform the discriminative occupancy classification task in 3D space, and visualize in Fig. 4 the predicted occupancy probability. In different columns, we vary the occupancy threshold τ , and only show the points with occupancy probability prediction that is higher than the given threshold. We can see that our model is able to output a continuous probability distribution of the occupancy function, even if it is only trained with discrete occupancy values from the sampled points.

When we take a closer look at the occupancy distribution, we find several interesting clues on how the model is modeling the probability distribution impressively well. We show our findings in Fig. 5. There are no points present in both **red** and **purple** regions of the masked point cloud, while in the original point cloud, there are points present in the **purple** region, and no points are in the **red**

region. In the reconstructed probability distribution, the model predicts a low occupancy probability in the **red** region, and a high occupancy in the **purple** region.

We find such predictions align with how a human might understand the scene. First, although there are no points in the purple region of the masked point cloud, given the partial view of the top-left region of the desk top, and the regions where the desk legs are present, it is very likely that there are points present in the purple region (top-right region of the desk top). As for the red region, the model’s prediction can be interpreted as follows: usually desk tops are rectangle-shaped; however, there do exist desks whose surface shrinks inside the region where the person sits. Given that there is not a decisive evidence that indicates how this particular desk instance is shaped, the model produces predictions with probability around 0.7, which is lower than other regions that are more certain (yellow points in Fig. 5 “Recon”, with $p > 0.9$).

These two intriguing and encouraging visualizations suggest that our pre-trained model is capable of modeling a continuous occupancy probability distribution, and it has learned a deep understanding of the input scene.

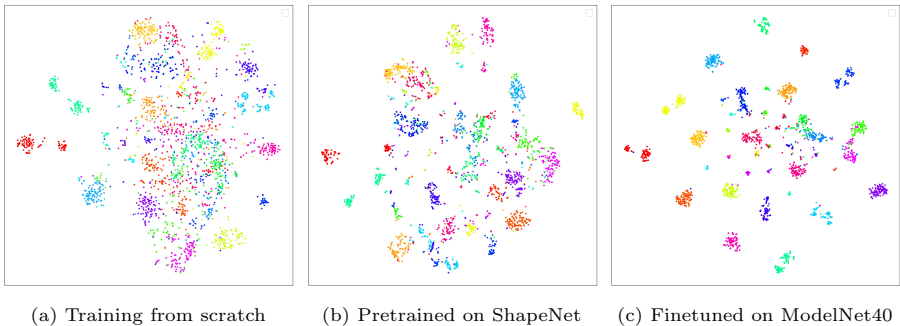


Fig. 6: t-SNE visualization of the encoder features for ModelNet40 under three settings: (a) training from scratch, (b) pretrained on ShapeNet, and (c) finetuned on ModelNet40.

t-SNE visualizations We show the t-SNE visualizations of the extracted feature vectors from our approach in Fig. 6. We use the class token from the encoder output as the high dimensional feature representation for t-SNE. Three setting are adopted here: (a) training from scratch, (b) pretraining on ShapeNet [4], and (c) finetuning on ModelNet40 [58].

When training the ModelNet40 classification model from scratch, the resulting features from different categories become heavily entangled, which can leads to less interpretable and robust predictions for new test-time inputs. In contrast, when pretraining the model on ShapeNet using our proposed MaskPoint, the features are much more distinguishable from each other. Furthermore, after finetuning on ModelNet40, the projected features from different classes become clearly separable from each other, which indicates the effectiveness of our approach. Interestingly, the feature clusters in our approach are quite tight. Such

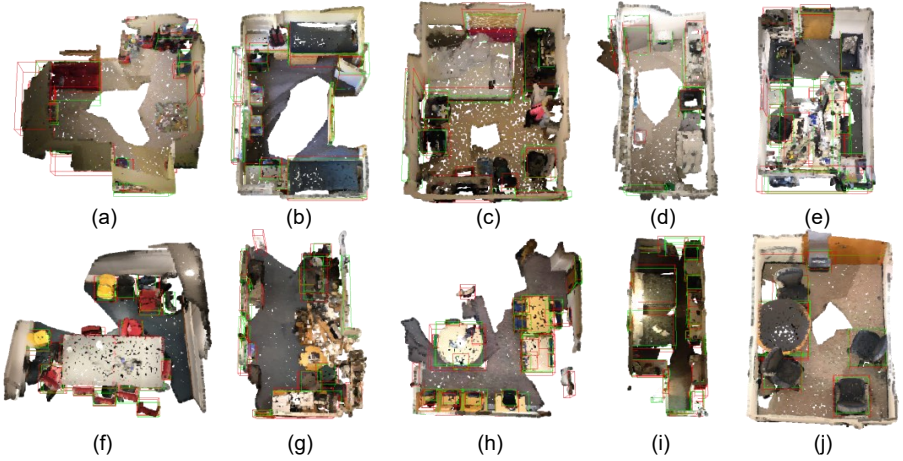


Fig. 7: **Qualitative results of 3D object detection on ScanNetV2 [9].** We show ground truth in **green** and predictions in **red** bounding boxes.

feature layout indicates that we can learn a more *compact and disjoint decision boundary*, which has been evaluated to be critical in machine learning applications such as mixup [66,49] and uncertainty estimation in deep learning [15].

3D object detection visualizations We show 3D object detection visualizations of ScanNetV2 in Figure 7 with **green** ground truth bounding boxes and **red** predicted bounding boxes. Our model is capable of precisely localizing the object (Fig. 7b and Fig. 7j). Results indicate that our masked based discriminative pretraining can not only produce high-quality bounding boxes for the previously annotated objects, but also discover objects that are not annotated. For example, in Fig. 7c, our model produces the a bounding box for the bookshelf in the lower region; in Fig. 7f, it correctly locates the sofa in the center of the room.

B Additional Implementation Details

B.1 Pretraining

Transformer Encoder. We follow the standard Transformer design in [52,65] to construct our point cloud Transformer backbone, PointViT. It consists of a linear stack of 12 Transformer blocks, where each Transformer block contains a multi-head self-attention (MHSA) layer and a feed-forward network (FFN). LayerNorm (LN) is adopted in both layers. Following [65], we use the MLP-based positional embedding. We set the Transformer hidden dimension to 384, MHSA head number to 6, expansion rate of FFN to 4, stochastic drop path [23] rate to 0.1.

Feed-forward network (FFN). Following [65], we use a two-layer MLP with ReLU and dropout as the feed-forward network. Dropout rate is set to 0.1.

Module	Block	C_{in}	C_{out}	k	N_{out}	C_{middle}
Positional Embed.	MLP	3	384			128
Point Classify Head	MLP	384	2			64
Classification Head	MLP	768	N_{cls}			512, 256
Segmentation Head	MLP	387	384			<u>384×4</u>
	DGCNN	384	512	4	128	
	DGCNN	512	384	4	128	
	DGCNN	384	512	4	256	
	DGCNN	512	384	4	256	
	DGCNN	384	512	4	512	
	DGCNN	512	384	4	512	
	DGCNN	384	512	4	2048	
	DGCNN	512	384	4	2048	

Table 9: Detailed module design of MaskPoint. C_{in}/C_{out} denotes the input/output channels, C_{middle} denotes the hidden channels of MLP modules, N_{out} denotes the cardinality of the output point/feature set, k is the number of neighbors used in the k -NN operator.

Positional Embeddings. Following [65], we use a two-layer MLP with GELU [22] as the positional embedding module. All Transformer modules share the same positional embedding MLP module. Detailed configuration is shown in Table 9.

Point Classification Head. We use a simple two-layer MLP with GELU [22] for the point classification pretext task in pretraining. We use the binary focal loss [29] to balance the information from positive and negative samples. Detailed configuration is shown in Table 9.

ScanNet-Medium Pretraining Note that for ScanNet-Medium pretraining, we use the encoder with 3 Transformer blocks, where each block still consists of a MHSA layer and a FFN layer. LN and MLP positional embedding are also utilized in the encoder. Following the downstream architecture of 3DETR [34], we set the hidden dimension to be 256, the number of MHSA heads to be 4, and Dropout rate to be 0.1 for the Transformer. The hidden dimension is set to be 128 for the FFN layer.

For other settings such as positional embedding and classification head, the setting is exactly the same as the ShapeNet pretraining setting.

B.2 Finetuning

Classification We use a three-layer MLP with dropout for the classification head. The input feature to the classification head consists of two parts from the Transformer encoder: (1) the CLS token; (2) the max-pooled feature of other output features. These two features are concatenated together and fed into the classification head. Detailed configuration is shown in Table 9.

config	value
epochs	300
optimizer	AdamW
learning rate	5e-4
weight decay	5e-2
LR schedule	cosine decay
warmup epochs	3
augmentation	Scale/Translate
batch size	128
# points	1024
# patches	64
patch size	32
mask ratio	0.90
mask type	random

Table 10: Pretraining setting on ShapeNet [4].

config	value
epochs	300
optimizer	AdamW
learning rate	5e-4
weight decay	5e-2
LR schedule	cosine decay
warmup epochs	10
augmentation	Scale/Translate
batch size	32(cls), 16(seg)
# points	1024(cls), 2048(seg)
# patches	64(cls), 128(seg)
patch size	32

Table 11: Finetuning setting on classification (cls) and segmentation (seg).

Part Segmentation The standard Transformer only has a single-scale feature output, which is not suitable for common head designs for dense prediction tasks like segmentation. Following [65], after getting the feature outputs from the Transformer encoder, we perform segmentation in two steps: (1) generating a multi-scale feature pyramid from the Transformer encoder outputs; (2) applying a standard feature propagation head for point cloud segmentation on the generated multi-scale feature maps to generate dense predictions.

We obtain the feature maps $f_{\{4,8,12\}} \in \mathbb{R}^{N_3 \times d}$ from the 4th, 8th, 12th layer, and our goal is to convert them to a feature pyramid with different cardinality $N_{\{0,1,2,3\}}$, where N_0 is the cardinality of the original point cloud \mathcal{P} , and $N_{\{1,2,3\}}$ are the desired cardinality of the feature maps at different scales; in our case, $N_{\{0,1,2,3\}} = \{2048, 512, 256, 128\}$.

First, we use furthest point sampling (FPS) to downsample the original point cloud \mathcal{P}_0 to different resolutions: $\mathcal{P}_{\{1,2,3\}} \in \mathbb{R}^{N_{\{1,2,3\}} \times 3}$, then a feature propagation module is used to upsample the feature maps $f_{\{4,8,12\}}$ to the corresponding cardinality $f_{\{4,8,12\}}^{up} \in \mathbb{R}^{N_{\{1,2,3\}} \times d}$.

After obtaining the multi-scale feature maps, we then apply the DGCNN module to propagate the features through different scales, $\hat{f}_4 = \text{DGCNN}(f_{\{4,8,12\}}^{up})$. Another feature propagation layer is then applied on \hat{f}_4 for upsampling to the highest resolution $\hat{f}_0 \in \mathbb{R}^{N_0 \times d}$.

Finally, we apply a pointwise MLP classifier on the features at the highest resolution \hat{f}_0 to obtain the segmentation results. Detailed configuration is shown in Table 9.

3D Object Detection We strictly follow the setting of the original 3DETR [34] model as the downstream 3D object detector. The points are first downsampled to 2048 points using a Set-Aggregation (SA) layer. The encoder is composed of 3 standard Transformer blocks. The decoder is comprised of 8 Transformer blocks using cross attention. During finetuning, only the weights of the SA layer

and the encoder are transferred to the downstream tasks. The finetuning epoch number is 1080, the optimizer is AdamW with learning rate of 5×10^{-4} and weight decay of 0.1, the batch size is 8.

References

1. Achituve, I., Maron, H., Chechik, G.: Self-supervised learning for domain adaptation on point clouds. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. pp. 123–133 (2021)
2. Bao, H., Dong, L., Wei, F.: Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254* (2021)
3. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: *European conference on computer vision*. pp. 213–229. Springer (2020)
4. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012* (2015)
5. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: III, H.D., Singh, A. (eds.) *Proceedings of the 37th International Conference on Machine Learning*. *Proceedings of Machine Learning Research*, vol. 119, pp. 1597–1607. PMLR (13–18 Jul 2020), <https://proceedings.mlr.press/v119/chen20j.html>
6. Chen, X., Duan, Y., Houthoofd, R., Schulman, J., Sutskever, I., Abbeel, P.: InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. *Advances in neural information processing systems* **29** (2016)
7. Chen, X., Ding, M., Wang, X., Xin, Y., Mo, S., Wang, Y., Han, S., Luo, P., Zeng, G., Wang, J.: Context autoencoder for self-supervised representation learning. *arXiv preprint arXiv:2202.03026* (2022)
8. Chen, X., He, K.: Exploring simple siamese representation learning. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 15750–15758 (2021)
9. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 5828–5839 (2017)
10. Danila Rukhovich, Anna Vorontsova, A.K.: Fcaf3d: Fully convolutional anchor-free 3d object detection. *arXiv preprint arXiv:2112.00322* (2021)
11. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: *2009 IEEE conference on computer vision and pattern recognition*. pp. 248–255. Ieee (2009)
12. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (Jun 2019). <https://doi.org/10.18653/v1/N19-1423>, <https://aclanthology.org/N19-1423>
13. DeVries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552* (2017)

14. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR* (2021)
15. Du, X., Wang, X., Gozum, G., Li, Y.: Unknown-aware object detection: Learning what you don't know from videos in the wild. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022)
16. Engel, N., Belagiannis, V., Dietmayer, K.: Point transformer. *IEEE Access* **9**, 134826–134840 (2021)
17. Ghiasi, G., Lin, T.Y., Le, Q.V.: Dropblock: A regularization method for convolutional networks. In: *NeurIPS* (2018)
18. Gidaris, S., Singh, P., Komodakis, N.: Unsupervised representation learning by predicting image rotations. In: *International Conference on Learning Representations* (2018), <https://openreview.net/forum?id=S1v4N210->
19. Guo, M.H., Cai, J.X., Liu, Z.N., Mu, T.J., Martin, R.R., Hu, S.M.: Pct: Point cloud transformer. *Computational Visual Media* **7**(2), 187–199 (2021)
20. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377* (2021)
21. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 9729–9738 (2020)
22. Hendrycks, D., Gimpel, K.: Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415* (2016)
23. Huang, G., Sun, Y., Liu, Z., Sedra, D., Weinberger, K.Q.: Deep networks with stochastic depth. In: *European conference on computer vision*. pp. 646–661. Springer (2016)
24. Huang, S., Xie, Y., Zhu, S.C., Zhu, Y.: Spatio-temporal self-supervised representation learning for 3d point clouds. *arXiv preprint arXiv:2109.00179* (2021)
25. Jiang, Y., Chang, S., Wang, Z.: Transgan: Two pure transformers can make one strong gan, and that can scale up. *Advances in Neural Information Processing Systems* **34** (2021)
26. Joshi, M., Chen, D., Liu, Y., Weld, D.S., Zettlemoyer, L., Levy, O.: Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics* **8**, 64–77 (2020)
27. Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., Soricut, R.: Albert: A lite bert for self-supervised learning of language representations. In: *International Conference on Learning Representations* (2020), <https://openreview.net/forum?id=H1eA7AEtvS>
28. Li, Y., Bu, R., Sun, M., Wu, W., Di, X., Chen, B.: Pointcnn: Convolution on x-transformed points. *Advances in neural information processing systems* **31** (2018)
29. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollar, P.: Focal loss for dense object detection. In: *ICCV* (2017)
30. Liu, Y., Fan, B., Meng, G., Lu, J., Xiang, S., Pan, C.: Densepoint: Learning densely contextual representation for efficient point cloud processing. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 5239–5248 (2019)
31. Liu, Y., Fan, B., Xiang, S., Pan, C.: Relation-shape convolutional neural network for point cloud analysis. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 8895–8904 (2019)
32. Liu, Z., Hu, H., Cao, Y., Zhang, Z., Tong, X.: A closer look at local aggregation operators in point cloud analysis. In: *European Conference on Computer Vision*. pp. 326–342. Springer (2020)

33. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017)
34. Misra, I., Girdhar, R., Joulin, A.: An End-to-End Transformer Model for 3D Object Detection. In: ICCV (2021)
35. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: European conference on computer vision. pp. 69–84. Springer (2016)
36. Poursaeed, O., Jiang, T., Qiao, H., Xu, N., Kim, V.G.: Self-supervised learning of point clouds via orientation estimation. In: 2020 International Conference on 3D Vision (3DV). pp. 1018–1028. IEEE (2020)
37. Qi, C.R., Litany, O., He, K., Guibas, L.J.: Deep hough voting for 3d object detection in point clouds. In: Proceedings of the IEEE International Conference on Computer Vision (2019)
38. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 652–660 (2017)
39. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. arXiv preprint arXiv:1706.02413 (2017)
40. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning. pp. 8748–8763. PMLR (2021)
41. Radford, A., Sutskever, I.: Improving language understanding by generative pre-training. In: arxiv (2018)
42. Rao, Y., Liu, B., Wei, Y., Lu, J., Hsieh, C.J., Zhou, J.: Randomrooms: Unsupervised pre-training from synthetic shapes and randomized layouts for 3d object detection. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3283–3292 (2021)
43. Rolfe, J.T.: Discrete variational autoencoders. In: ICLR (2017)
44. Singh, K.K., Lee, Y.J.: Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization. In: 2017 IEEE international conference on computer vision (ICCV). pp. 3544–3553. IEEE (2017)
45. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* **15**(56), 1929–1958 (2014), <http://jmlr.org/papers/v15/srivastava14a.html>
46. Steiner, A., Kolesnikov, A., Zhai, X., Wightman, R., Uszkoreit, J., Beyer, L.: How to train your vit? data, augmentation, and regularization in vision transformers. arXiv preprint arXiv:2106.10270 (2021)
47. Thabet, A., Alwassel, H., Ghanem, B.: Self-supervised learning of local features in 3d point clouds. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 938–939 (2020)
48. Thomas, H., Qi, C.R., Deschaut, J.E., Marcotegui, B., Goulette, F., Guibas, L.J.: Kpconv: Flexible and deformable convolution for point clouds. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 6411–6420 (2019)
49. Thulasidasan, S., Chennupati, G., Bilmes, J.A., Bhattacharya, T., Michalak, S.: On mixup training: Improved calibration and predictive uncertainty for deep neural networks. *Advances in Neural Information Processing Systems* **32** (2019)
50. Tolstikhin, I., Hounsby, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A., Keysers, D., Uszkoreit, J., Lucic, M., Dosovitskiy, A.: Mlp-mixer: An all-mlp architecture for vision. arXiv preprint arXiv:2105.01601 (2021)

51. Uy, M.A., Pham, Q.H., Hua, B.S., Nguyen, T., Yeung, S.K.: Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 1588–1597 (2019)
52. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*. vol. 30. Curran Associates, Inc. (2017), <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>
53. Wang, H., Liu, Q., Yue, X., Lasenby, J., Kusner, M.J.: Unsupervised point cloud pre-training via occlusion completion. In: ICCV (2021)
54. Wang, H., Zhu, Y., Adam, H., Yuille, A., Chen, L.C.: MaX-DeepLab: End-to-end panoptic segmentation with mask transformers. In: CVPR (2021)
55. Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M.: Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)* **38**(5), 1–12 (2019)
56. Wei, C., Fan, H., Xie, S., Wu, C.Y., Yuille, A., Feichtenhofer, C.: Masked feature prediction for self-supervised visual pre-training. *arXiv preprint arXiv:2112.09133* (2021)
57. Wu, W., Qi, Z., Fuxin, L.: Pointconv: Deep convolutional networks on 3d point clouds. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9621–9630 (2019)
58. Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1912–1920 (2015)
59. Xiao, T., Dollar, P., Singh, M., Mintun, E., Darrell, T., Girshick, R.: Early convolutions help transformers see better. *Advances in Neural Information Processing Systems* **34** (2021)
60. Xie, S., Gu, J., Guo, D., Qi, C.R., Guibas, L., Litany, O.: Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In: European conference on computer vision. pp. 574–591. Springer (2020)
61. Xu, Y., Fan, T., Xu, M., Zeng, L., Qiao, Y.: Spidercnn: Deep learning on point sets with parameterized convolutional filters. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 87–102 (2018)
62. Yan, S., Yang, Z., Li, H., Guan, L., Kang, H., Hua, G., Huang, Q.: Implicit autoencoder for point cloud self-supervised representation learning. *arXiv preprint arXiv:2201.00785* (2022)
63. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R.R., Le, Q.V.: Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems* **32** (2019)
64. Yi, L., Kim, V.G., Ceylan, D., Shen, I.C., Yan, M., Su, H., Lu, C., Huang, Q., Sheffer, A., Guibas, L.: A scalable active framework for region annotation in 3d shape collections. *ACM Transactions on Graphics (ToG)* **35**(6), 1–12 (2016)
65. Yu, X., Tang, L., Rao, Y., Huang, T., Zhou, J., Lu, J.: Point-bert: Pre-training 3d point cloud transformers with masked point modeling. *arXiv preprint arXiv:2111.14819* (2021)
66. Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In: ICCV (2019)
67. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: European conference on computer vision. pp. 649–666. Springer (2016)

68. Zhang, Z., Girdhar, R., Joulin, A., Misra, I.: Self-supervised pretraining of 3d features on any point-cloud. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 10252–10263 (October 2021)
69. Zhang, Z., Sun, B., Yang, H., Huang, Q.: H3dnet: 3d object detection using hybrid geometric primitives. In: European Conference on Computer Vision. pp. 311–329. Springer (2020)
70. Zhao, H., Jiang, L., Fu, C.W., Jia, J.: Pointweb: Enhancing local neighborhood features for point cloud processing. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 5565–5573 (2019)
71. Zhao, H., Jiang, L., Jia, J., Torr, P.H., Koltun, V.: Point transformer. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 16259–16268 (2021)
72. Zhong, Z., Zheng, L., Kang, G., Li, S., Yang, Y.: Random erasing data augmentation. arXiv preprint arXiv:1708.04896 (2017)