

DiffusionPointLabel: Annotated Point Cloud Generation with Diffusion Model

Abstract

Point cloud generation aims to synthesize point clouds that do not exist in supervised datasets. Generating a point cloud with certain semantic labels remains an under-explored problem. In this paper, we propose a formulation called DiffusionPointLabel, which conducts point-label pair generation based on DDPM generative model (denoising diffusion probabilistic model). Specifically, we use a point cloud diffusion generative model and aggregate the intermediate features of the generator. On top of this, we propose Feature Interpreter that transforms intermediate features into semantic labels. Furthermore, we employ an uncertainty measure to filter unqualified point-label pairs for a better generated point cloud dataset. Coupling these two designs enables us to automatically generate annotated point clouds, especially when supervised point-labels pairs are scarce. Our method extends the application of point cloud generation models and surpasses state-of-the-art models.

CCS Concepts

• **Methods and Applications** → Point-Based Methods;

1. Introduction

In recent years, Deep Neural Network has dominated point cloud processing and understanding tasks, such as object detection [YZK21], robot manipulation [YKH*19], depth estimation [WCG*19], and semantic segmentation [JST*21]. Though substantial progress has been made, point cloud application based on modern deep learning suffer from a practical limitation. It usually requires large amounts of annotated data to optimize all parameters of the network.

Unfortunately, creating point cloud datasets with point labels such as semantic or instance segmentation is labor-intensive and expensive. This is because labeling a complex shape usually involves the help of a human auxiliary to rotate and look through different angles to identify an object from incomplete or occluded point cloud data. In this scenario, creating a point cloud dataset of the scale as we desire is still a challenge [QYW*19].

One solution is to build generative models [MWYG20, GBZCO21, YWZJ21, SPK19] to synthesize expressive point clouds while having control of the structure. In these methods, point clouds and point-wise semantic labels are bred from key structural points. However, due to the irregular distribution and high complexity of 3D point clouds, existing generative models often struggle with explicit structural controllability and realistic-looking shape. Our approach goes beyond existing solutions, in that it generate point-wise labels without holding up the shapes generation results, because the semantic information is obtained from the intermediate feature of the generator.

Recently, Denoising Diffusion Probabilistic Models (DDPM)

emerging as a new class of generative models and have achieved impressive performance on point cloud generation [LH21, LKX*21]. DDPM defines a consecutive point-wise mapping between two point clouds in the diffusion process and characterizes it as a Markov chain. Our approach is motivated by the observation that the generator of diffusion model primarily recovers the structure of the point cloud at the early stage of the diffusion process, while gradually enriching surface details at the later stage. However, the output of the generator alongside the diffusion process should be a set of independent and identically distributed random variables. Based on this observation, we assume that the discriminability of point representations of the diffusion model varies along the diffusion process.

Driven by this assumption and inspiration of [BRV*21], we investigate the intermediate features of the point cloud diffusion generative model to figure out how the discriminability changes and whether it has semantically interpretable potential. On top of that, we aggregate the intermediate features of the diffusion generator and conduct a simple MLP, called feature interpreter, to transform the intermediate features into point-wise semantic labels. Under this design, our paradigm enables point-wise label annotation without affecting the quality of point cloud generation. Although our approach can generate plausible annotated point clouds, we find that it still generates some unqualified point-label pairs. To eliminate this error, following the setup of [ZLG*21], we adopt an uncertainty measurement to estimate the quality of annotated point-label pairs. Specifically, we train a committee of feature interpreters and compute uncertainty score for point-label pairs via the entropy of

this committee of feature interpreters. Then uncertainty score can be used to filter unqualified point clouds.

The main contributions of our work can be summarized as follows.

1. Different from previous point cloud generation methods that focus on structure-aware point cloud generation or breeding point labels from key structural points, we are the first to propose a point-label pairs generation framework based on point cloud diffusion generative model, termed as DiffusionPointLabel. Our method can simultaneously generate point cloud and the corresponding point-wise semantic labels.

2. We experimentally exhibit that the intermediate features of the point cloud diffusion generative model are interpretable at the semantic level and have the potential to help 3D understanding.

3. Experimental results demonstrate that the proposed method has great advantages in efficiency and effectiveness to obtain annotated point clouds, especially when the supervised labeled examples are scarce, which surpasses state-of-the-art performances.

2. Related Works

In this section, we briefly describe the existing research lines relevant to our work.

2.1. Learning Representations on Point Cloud

Deep representation learning has been developed for many years. Since point cloud data has irregular structure, [MS15a, WSK*15, ROUG17, MS15b] quantized the 3D space and transformed points into regular voxels so that the convolution neural network can process 3D data. However, since the 3D point cloud is a sparse and discretely distributed representation, convolution operators are inefficient and computationally expensive for 3D data.

Qi et al. [QSMG17] proposed PointNet, a notable landmark for point-based deep learning work, which works by leveraging weight-shared multi-layers perceptrons and a point-wise max-pooling layer to learn the features of the point cloud. The max-pooling layer can address the irregular structure of the point cloud, while it may neglect local information. Subsequent works have been proposed to tackle this issue. PointNet++ [QYSG17] defined a hierarchical architecture, which is effective for capturing local features of increasing contextual scale point set and improved semantic segmentation performance. Its design includes a deformable convolutional kernel to adapt to the local geometry and be robust to varying densities. Following them, [LCL18, ZJFJ19, WSL*19a] adopted a wider neighborhood to enhance local region features; [WQF19, TQD*19] designed a flexible kernel-based convolution operator and [WSL*19a, QLJ*17, WHH*19] regarded point cloud as an undirected graphs to group points to enrich latent features.

Recently, random walks [MBST21, XZS*21] have been used for 3D model representations and achieved the state-of-the-art performance. Xiang et al. [XZS*21] proposed to use shape curves to analyze point cloud feature, which is initialized based on a given set of rules and heuristics. Mesika and Ben-Shabat [MBST21] presented

a technique that imposes structure on the point set by multiple random walks to aggregate point features.

Following the recent success of Transformers [VSP*17] in various vision tasks, there is some work that uses this network architecture for point cloud understanding tasks. With the recent success of applying Transformer [VSP*17] in vision tasks, many works [ZWL*21, ZJJ*21, GCL*21, HJCX21, YTR*21] have proposed their transformer network frameworks for point clouds. These models focused on reducing the cost of point cloud annotation. But a point cloud dataset consists of point clouds and its corresponding label. In this paper, we propose a pipeline that can generate a point cloud dataset.

2.2. Generative Models for 3D Point Cloud

In the past few years, plenty of works have extended the generative model to point clouds. Current point cloud generative works can be generally classified into three categories: Autoregressive-based, flow-based, and GAN-based.

PointGrow [SWL*20] is one of the notable works of **Autoregressive-based** methods. It estimates the probability of samples one-by-one based on previously generated points. However, this method is restricted to generating a fixed-dimension point cloud because it assumes a determinate order of point cloud.

GAN-based generative models explore adversarial learning to train the shape generator with the help of a discriminator. Shu et al. [SPK19] combined tree-structure and graph to perform convolution on the point cloud. It demonstrated that tree-GAN can edit point clouds on the semantic level without prior knowledge, but the precision of the label falls short of expectations. Gal et al. [GBZCO21] extended it into multi-roots version. The node of the mutil-roots can generate and control different parts of point cloud. But there is no clear classification boundary between different parts, that is they do not have clear semantic definition. Wang et al. [YWZJ21] draw inspiration from S^2 -GANs and proposed using enhance controllability and point-level label accuracy. However, the label accuracy will inevitably be affected because the semantic label of a point is inherited by the structure points. Compared with the above GAN-based approaches, we incorporate a pre-trained generator and use its intermediate features to generate point-wise label, which improves the accuracy.

For **Flow-based** generative models [KLL*20, YHH*19, KBV20, HXX*20], the basic idea is to train an invertible parameterized transformation that can characterize the distribution of samples. This transformation can output a target shape by moving points from a prior distribution at one time.

Recently, denoising diffusion probabilistic models have shown superior performance in terms of generative fidelity and diversity in 2D dataset generation [GRS*20]. For 3D generation, Luo et al. [LH21] applied a diffusion model to point clouds and achieved competitive results compared to state-of-art. Zhou et al. [ZDW21] used conditional DDPM for point cloud completion by training a point-voxel CNN. Lyu et al. [LKK*21] applied a diffusion model to point clouds completion task. It proposed an adaption network architecture for point cloud and added a denoise module, which en-

hances the precision of results. However, existing point-based generative approaches mainly focus on 3D geometry while neglecting the implicit feature information that is complementary to 3D understanding.

More recently, SetVAE [KYLH21] learned to generate set-structured data such as point clouds with a tree-structure. However, its latent variable are not explicitly trained to segregate semantic part.

3. Methodology

Figure 1 shows the structural details of our method. Given a random latent code $z \sim \mathcal{N}(0, I)$ and a random noise of point cloud $X \in \mathbb{R}^{N \times 3} \sim \mathcal{N}(0, I)$, we aim to generate a point cloud $X \in \mathbb{R}^{N \times 3}$ and its corresponding semantic labels $SL \in \mathbb{R}^{N \times 4}$. To this end, firstly we extract the intermediate features of X in different layers of the diffusion generator θ_G at certain time steps $t = \{t_i | i = 1, \dots, T\}$. Formally, $C_{i,j} \in \mathbb{R}^{N \times C_{out_i}}$ denotes the intermediate feature of the i -th layer at time step $t = j$. We concatenate several $C_{i,j}$ as C^* . Finally we use a feature interpreter to transform C^* into point-label pairs.

We revisit the point cloud diffusion generative model in Section 3.1. In Section 3.2, we analyze how discriminability of intermediate features of the diffusion model changes along the diffusion process. We use K-means to verify our assumption that the intermediate features of the diffusion model are interpretable at semantic level. Then in Section 3.3, we introduce the feature interpreter which transforms the intermediate features of the diffusion generator into the point-wise label. In Section 3.4, we propose assembling a set of feature interpreters as a committee to compute the uncertainty score of annotated point clouds which can be used to filter unqualified point-label pairs.

3.1. Denoising Diffusion Probabilistic Models for Point Cloud

Denoising Diffusion Probabilistic Model regards the process of point cloud generation as a Markov chain. The parameterized Markov chain maps noise (i.e. 3D Gaussian) to shape by recursively perturbing the input point cloud. The forward diffusion process transforms shape to noise in an unconditional way. The reverse process generates the desired shape from Gaussian noise that is conditioned on a latent variable of global shape. Both processes have a fixed time step, denoted by T . Ho et al. [HJA20] utilized variational inference to solve the parameterized diffusion process.

The Diffusion Process. We use superscript to denote the diffusion step t . The forward diffusion process q transforms shape X^0 to noise $X^T \sim \mathcal{N}(0, 1)$. We assume $p_{data}(X^{(0)})$ to be the distribution of the point cloud X in the ground-truth dataset. Given N points in a point cloud $X = \{x_i | i = 1, \dots, N\} \in \mathbb{R}^{N \times 3}$, the distribution of each point in forward diffusion process can be formulated to:

$$q(x_i^{(1:T)} | x^{(0)}) := \prod_{t=1}^T q(x_i^{(t)} | x_i^{(t-1)}), \quad (1)$$

where $q(x^{(t)} | x^{(t-1)}) := \mathcal{N}(x^{(t)}; \sqrt{1 - \beta_t} x^{(t-1)}, \beta_t \mathbf{I})$

where the

According to [HJA20], the hyperparameter β_t is a fixed monotonic increasing list.

Note that in the forward diffusion process, the shape sample X^0 gradually loses its geometric features as the time step t . Eventually when $T \rightarrow \infty$, X^T is equivalent to an isotropic Gaussian distribution.

The Reverse Process. The reverse process p is to predict a 3D shape from a latent code z . Conversely to forward process, points are recursively moved from a prior noise distribution $p(x_i^{(T)})$ to approximate $q(x_i^{(0)})$. We use a network θ to estimate the denoise movement of every point at time step t . This process can be formulated as:

$$p_\theta(x^{(0:T)} | z) := p(x^T) \prod_{t=1}^T p_\theta(x^{(t-1)} | x^{(t)}, z) \quad (2)$$

$$p_\theta(x^{t-1} | x^t, z) := \mathcal{N}(x^{t-1}; s_\theta(x^t, t), \Sigma_\theta(x^t, t))$$

Since the points in a point cloud are independently sampled from a distribution, the probability of the whole point cloud is simply the product of the probability of each point:

$$q(\mathbf{X}^{(1:T)} | \mathbf{X}^0) = \prod_{i=1}^N q(x_i^{(1:T)} | x_i^{(0)}) \quad (3)$$

$$p_\theta(\mathbf{X}^{(0:T)} | z) = \prod_{i=1}^N p_\theta(x_i^{(0:T)} | z)$$

Training. Training objective can be simplified by optimizing the variational bound on negative log-likelihood:

$$\mathcal{L} = \mathbb{E}_q \left[\sum_{t=2}^T D_{KL}(q(\mathbf{X}^{(t-1)} | \mathbf{X}^{(t)}, \mathbb{X}^{(0)}) \| p_\theta(\mathbf{X}^{(t-1)} | \mathbf{X}^{(t)}, z)) \right] \quad (4)$$

$$- \log p_\theta(\mathbf{X}^{(0)} | \mathbf{X}^{(1)}, z) + D_{KL}(q_\phi(z | \mathbf{X}^{(0)}) \| p(z))$$

Ho et al. [HJA20] showed that the training objective of the diffusion model θ can be simplified in a closed-form by a parameterization trick. Let $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_1^t \alpha_i$, then the training objective becomes:

$$\mathcal{L}(\theta) := \mathbb{E}_{t, x^{(0)}, \varepsilon} \|\varepsilon - \varepsilon_\theta(\sqrt{\bar{\alpha}_t} x^{(0)} + \sqrt{1 - \bar{\alpha}_t} \varepsilon, t)\|^2 \quad (5)$$

where t is uniform distributed between 1 and T , $\varepsilon \sim \mathcal{N}(0, I)$, ε_θ is a point cloud diffusion generative network that predicts injected noise ε at each time step. DDPM does not require CD or EMD loss in training, because it defines a consecutive and invertible point-wise mapping. Note that the network is non-autoregressive, its prediction is only determined by the predecessor.

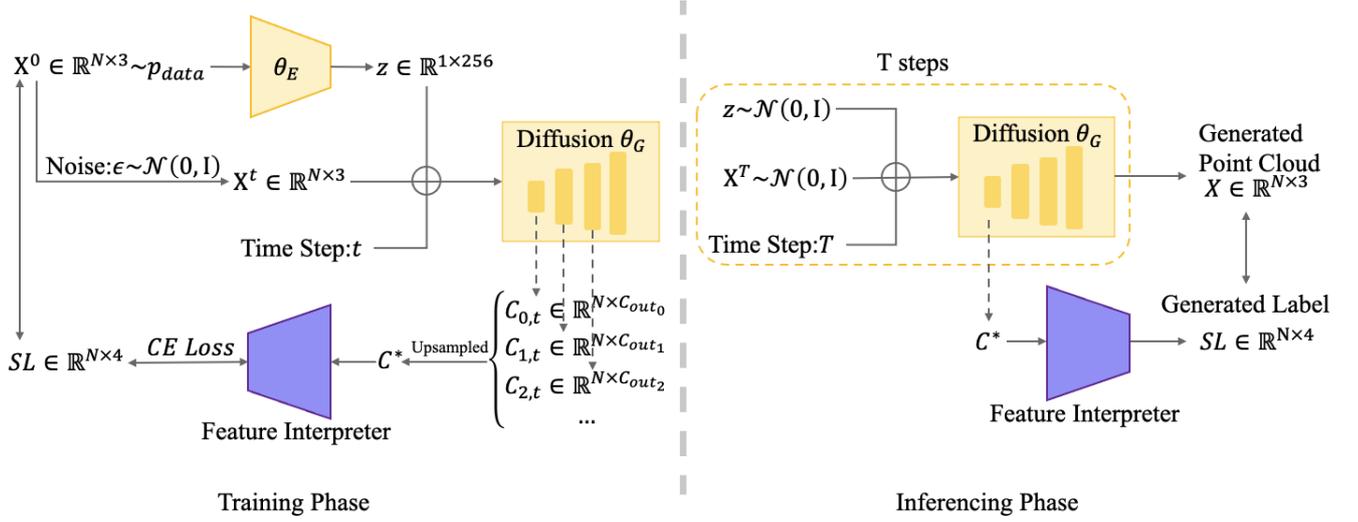


Figure 1: Illustration of the proposed point-annotation pairs generation method.

3.2. Intermediate Feature Analysis

The point cloud generation process based on the diffusion generative model is shown in the top row of Figure 2. As shown in Equation 5, the output of the diffusion generator at each step should be independent and identically distributed random variables. However, in the early and later stages of the diffusion process, the change tendency of the point cloud is different: structural features in the former and surface details in the latter. Therefore, we assume that the point representation of the diffusion generator has different discriminability alongside the diffusion process.

To prove our assumption, we use the K-means cluster to analyze the intermediate features of the diffusion generator at different time steps. In practice, we freeze a diffusion model θ and take a point cloud $X \in \mathbb{R}^{N \times 3}$ and a time step t as the inputs of θ . Then we extract the intermediate features of one layer of the generator θ_G . Because the generator we used is an MLP, the intermediate features of every layer at each time step t can be denoted by $C_{i,t} \in \mathbb{R}^{N \times out_i}$. We use the K-means clustering algorithm to estimate the cluster label of each point from the intermediate features of X and visualize the results, as shown in the bottom row of Figure 2. In K-means clustering, we use the number of ground truth labels as the cluster number. For example, the airplane in Figure 2 has 4 semantic parts and then we set K-means cluster number as 4. The results of K-means demonstrate that the discriminability of intermediate layer features gradually increases with decreasing time steps and is interpretable at the semantic level.

To further determine which layer of features we should extract or at which time steps, we quantitatively compute the results of K-means clustering. If the K-means clustering effect is very good (the cluster labels of the close points are very similar, and the cluster labels of the distant points are not the same), it means that the discriminability of this intermediate feature is very high, otherwise it is very low and cannot be used for further learning. We extract and cluster the intermediate features of each layer of θ_G from time $t = 0$

to T . We compute the clustering results with the Calinski-Harabasz Index algorithm [CH74] and the results is shown in Figure 3. The Calinski-Harabasz Index algorithm is used to measure the quality of the cluster model without the ground-truth label. The Calinski-Harabasz score is defined as the ratio of between-group dispersion to within-group dispersion, and be formulated as:

$$s = \frac{SS_B}{k-1} / \frac{SS_W}{N-k} \quad (6)$$

where k denotes the number of clusters, N denotes the number of all data, SS_B denotes the variance between different clusters, SS_W denotes the variance between one cluster. The lower the score, the better the clustering effect of K-means.

As shown in Figure 3, this Calinski-Harabasz Index score increases abruptly around $t = T/4$, so we only extract intermediate feature at $t < T/4$ in experiments. It is worth mentioning that this score becomes stable at $t \rightarrow T$. The attribution of it is that the discriminability of intermediate features tends to be consistent at this time, all point features represent the category or overall shape information of the airplane.

3.3. Feature Interpreter

The feature interpreter takes the intermediate features as the input, aiming to generate explicit semantic labels for generated point cloud. An MLP is implemented to realize the label prediction. Similar to K-means clustering process, we freeze a diffusion model θ and take a point cloud $X \in \mathbb{R}^{N \times 3}$ and a time step t as the inputs of θ . Based on the analysis in section 3.1, we sample $C_{0,t}$ across different time steps, where $t < T/4$. Then the intermediate features $C_{0,t}$ of the θ_G are upsampled and concatenated to form $C^* \in \mathbb{R}^{N \times 1024}$. In practice, we use a three-layers MLPs to predict the semantic label for each point from the C^* . The feature interpreter is optimized by cross-entropy loss.

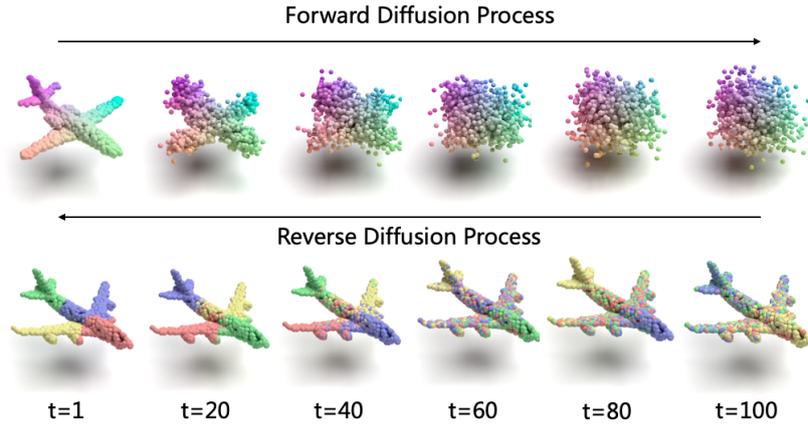


Figure 2: Visualization of the diffusion process and corresponding K -means cluster of intermediate features. The top row represents the effect of point cloud diffusion in the time variable. The bottom row represents the evolution of corresponding point-wise K -means features based on the diffusion model.

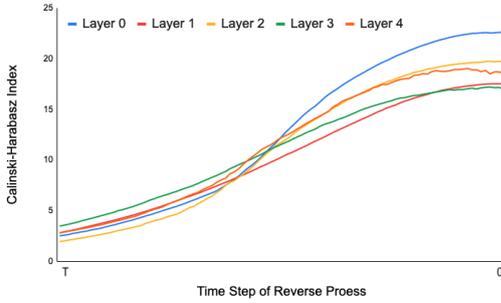


Figure 3: Calinski-Harabasz Index of feature space cluster. The score represents the quality of cluster. Different colors represent different layers of the θ_G

3.4. Uncertainty Measurement

We found that the point cloud generative models occasionally generate meaningless point clouds in experiments. Since we don't want to involve any human labor in this task, we need to filter these point clouds out before we collect the final dataset. Following [ZLG*21, GRS*20], we adopted the Jensen-Shannon (JS) divergence [KHY*18] to compute the uncertainty measure for each point-label pair. Specifically, we train a committee of feature interpreters in the same way. And then we estimate the label likelihood $LS \in \mathbb{R}^{N \times 4}$ for the point cloud $X \in \mathbb{R}^{N \times 3}$. Formally, the uncertainty measurement is denoted by $\mathcal{JS} \in \mathbb{R}^N$. The computation can be formulated as:

$$\mathcal{JS} = H\left(\frac{1}{N} \sum_i LS_i\right) - \frac{1}{N} \sum_i H(LS_i) \quad (7)$$

where N denotes the number of feature interpreters in one committee; LS_i denotes the label likelihood of the i -th feature interpreter for point cloud; H denotes the entropy function. We use the top 10% of \mathcal{JS} as the uncertainty score of each point cloud in the im-

plementation. The uncertainty score can be used to filter unqualified point cloud.

4. Evaluation and Discussion

In section 4.1, we evaluate the effectiveness of our method in two aspects: a) the usefulness of our generated dataset; b) the effectiveness of the feature interpreter. In section 4.2, we discuss the discriminability of intermediate features in three popular point cloud autoencoder networks. In section 4.3, we compare our method with a close generative model CPCGAN [YWZJ21]. We conduct an ablation study of our methods in section 4.4.

Dataset and Implementation Details. We evaluate the proposed method on ShapeNet-Partseg dataset [FSG17]. This dataset includes 16881 shapes from 16 object categories. In our evaluation, we mainly work with the "chair", "airplane", "guitar", "table", "lamp", "car" and "bag" categories.

Our method was conducted on a pre-trained point cloud diffusion generative model. We used the Luo et al. [LH21] as the backbone of our method and trained the model θ according to their specifications. We trained the feature interpreter for 100 epochs with batch size 128, starting with a learning rate of 0.02 which decayed by 0.5 every 10 epochs. Cross-entropy loss is minimized during training.

Evaluation Metrics and Baselines. 3D semantic segmentation is evaluated using mean Intersection over Union (mIoU) on point and accuracy, refed to mIoU and mAcc. We include PointNet [QSMG17], PointNet++ [QYSG17], DGCNN [WSL*19b] as baselines that verify the useful of the generated dataset.

4.1. Validation of Generated Datasets and Feature Interpreter

In this section we evaluate our method in two settings.

The usefulness of the generated dataset: We first train a network [QYSG17] for semantic segmentation using the train set of

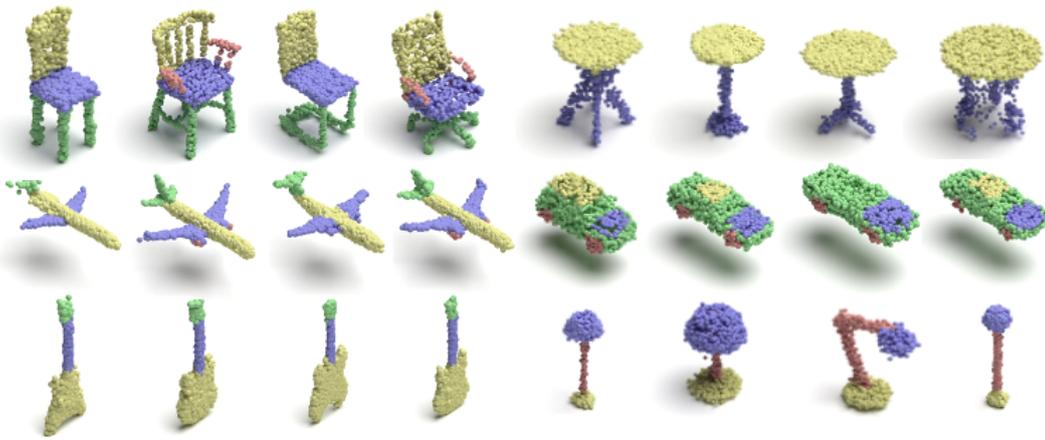


Figure 4: Examples of generated point-label pairs.

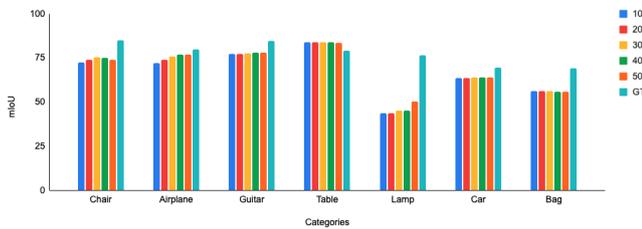


Figure 5: Comparison between our generated datasets and ground-truth ones. Different colors denote different filter ratios.

ground truth dataset. After training, we use this network to validate on our generated dataset and the ground truth test set separately. When producing our dataset, we generated 10,000 point clouds for each category, and then filtered samples based on their uncertainty scores. Figure 5 shows the quantitative comparison of the generated dataset produced by our method. Some examples from our generated dataset are visualized in Figure 4. For most categories, our generated datasets show competitive results compared to GT datasets. For Lamp, the performance of our generated dataset is much lower than the GT dataset. But the visualized results of Lamp are plausible. We attribute this result to the fact that because the Lamps in the GT dataset are few, the segmentation network has not fully learned the accurate features of the Lamps.

The effectiveness of the feature interpreter: We believe that one of the future application scenarios of our method is to generate point cloud datasets in a new category. Since the cost of point cloud annotating is too high, we can use few-shot samples of point-label pairs and generate large-scale point cloud datasets. Therefore, it is important to verify that our method can still generate results with high segmentation accuracy when there are only a few samples. We conduct few-shot segmentation to verify the effectiveness of our method. We compare the evaluation results with baseline [QYSG17] as shown in Table 1. The comparison results

demonstrate that our method is capable of generating compelling semantic labels in a few-shot setting.

4.2. Validation of Representation Effectiveness

Since the intermediate features of our analysis and learning are extracted from a diffusion generative network that has an autoencoder frame, we naturally question whether these learnable intermediate features have nothing to do with the diffusion process, but only benefit from the autoencoder framework? Therefore, we conduct an experiment to find out whether other methods capable of extracting intermediate features from point clouds can achieve the same effect. To the best of our knowledge, this is the first work to find out the discriminability of intermediate features in a point cloud generative model.

As in our method, we first collect and cluster latent feature space of existing generative models: CapsNetwork [ZBDT19], PointFlow [YHH*19], and FoldingNet [YFST18]. The cluster effect is shown in Figure 6.

The comparison results answer our question: the discriminability of the intermediate features benefits from the diffusion process, and not all point cloud autoencoders networks have similar discriminability. The possible explanations could be that, (a) these model use CD-Loss to optimize the parameter of network, which calculate the overall structural similarity; (b) these model train the network in a one-shot discriminative way. Therefore, their intermediate feature space does not contain fine-grained information.

4.3. Comparing with Related Methods

The work most closely related to ours is CPCGAN [YWZJ21]. CPCGAN [YWZJ21], which proposed a two-stage GAN to generate point clouds in structure controllable manner, and can be trained on ShapeNet-Partseg dataset as well. The first stage generator generates the key structural points and its label and the second stage

| Category | Model | k=1 | k=3 | k=5 | k=10 | k=16 | k=32 |
|----------|----------|-------------|-------------|-------------|-------------|-------------|-------------|
| Airplane | Baseline | 20.9 | 47.2 | 29.4 | 43.3 | 59.6 | 64.6 |
| | Ours | 58.1 | 62.8 | 63.9 | 64.8 | 66.0 | 67.2 |
| Chair | Baseline | 33.9 | 63.8 | 50.0 | 64.8 | 79.5 | 81.6 |
| | Ours | 66.2 | 67.9 | 72.1 | 74.7 | 77.6 | 78.2 |

Table 1: Few-shot segmentation on the ShapeNet dataset. The number of samples is denoted by k for each category. Our method demonstrates comparable performance to baseline when trained on a few samples.

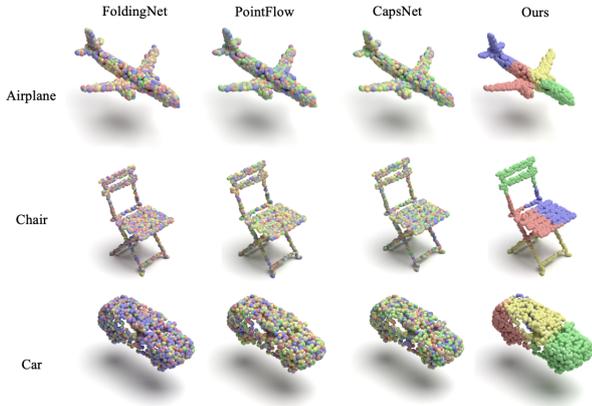


Figure 6: Qualitative comparison of intermediate features based on different baseline using K-means cluster.

| Class | Model | mIoU%(\uparrow) | mAcc%(\uparrow) |
|----------|--------|---------------------|---------------------|
| Chair | CPCGAN | 57.1 | 83.6 |
| | Ours | 72.0 | 86.3 |
| Airplane | CPCGAN | 67.8 | 82.6 |
| | Ours | 74.2 | 89.2 |

Table 2: Comparison of point cloud and label generation performance. mIoU and mACC is multiplied by 10^2

generator generates the point cloud by expanding the key structural points into complete point cloud. The semantic label is bred from the first stage. Because it is hard to annotate ground truth label for a generate point cloud, following their setting, we train a PointNet++ [QYSG17] for semantic segmentation task. We use this pre-trained segmentation network to evaluate our generated dataset. The quantitative comparison is shown in Table 2. From the results shown in the table, we can see that our generated point cloud (airplane and chair) outperforms their method consistently on the two evaluation metrics for both mIoU and mACC by a large margin.

Moreover, we visualized comparison results as shown in Figure 7. Here, we pick some random point clouds generated by both methods and categorize the semantic labels by color, using the same color for the same label across models. From these results, we can see that the point clouds with semantic labels generated by our method exhibit more accurate labels, whereas [YWZJ21] tends to generate noisy semantic labels.

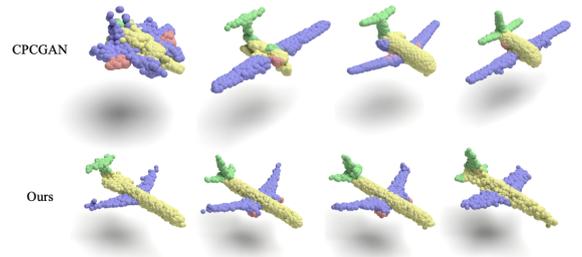


Figure 7: Visualized results of [YWZJ21] and ours.

| | $C_{(0,0)}$ | Upsample $C_{(0,0)}$ | $C_{(2,0)}$ |
|----------|-------------|----------------------|-------------|
| Airplane | 75.7 | 76.0 | 72.8 |

Table 3: Evaluation of the different intermediate feature extraction variations for part segmentation.

4.4. Ablation Study

Intuitively, there are two deterministic factors of representation discriminability: The intermediate features with the highest dimension have better discriminability because they may contain the most information. The second is that we tend to consider the features of the shallow layer because the feature of the deeper layer is closer to the estimated noise of the diffusion process, while the shallows contain abstract information, such as semantics.

Then we compare it to the following settings: a), the features of the shallow layers are upsampled to the highest dimension; b) the features of the layer that has the highest dimension.

The results are proved in Table 3. The intermediate features within the highest dimension slightly underperform than features of the shallow layers. The experimental results confirmed that the discriminability of the intermediate features of the shallow layer is better.

5. Conclusions

To conclude, this paper presents a paradigm called DiffusionPoint-Label which is a simple and useful method for point-label pairs generation. A feature interpreter is applied to transform intermediate features of the point cloud diffusion generative model into the semantic label. Uncertainty measurement is introduced to enhance the quality of the generated point cloud dataset. We further show

the effectiveness and efficiency of our method within scarced supervised labeled examples. In the future, we plan to explore more of the potential power of the point cloud generative model, such as fine-grain point-label pairs generation.

References

- [BRV*21] BARANCHUK D., RUBACHEV I., VOYNOV A., KHRULKOV V., BABENKO A.: Label-efficient semantic segmentation with diffusion models. *arXiv preprint arXiv:2112.03126* (2021). 1
- [CH74] CALIŃSKI T., HARABASZ J.: A dendrite method for cluster analysis. *Communications in Statistics-theory and Methods* 3, 1 (1974), 1–27. 4
- [FSG17] FAN H., SU H., GUIBAS L. J.: A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 605–613. 5
- [GBZCO21] GAL R., BERMANO A., ZHANG H., COHEN-OR D.: Mrgan: Multi-rooted 3d shape representation learning with unsupervised part disentanglement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 2039–2048. 1, 2
- [GCL*21] 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 (2021), 187–199. 2
- [GRS*20] GADELHA M., ROYCHOWDHURY A., SHARMA G., KALOGERAKIS E., CAO L., LEARNED-MILLER E., WANG R., MAJI S.: Label-efficient learning on point clouds using approximate convex decompositions. In *European Conference on Computer Vision* (2020), Springer, pp. 473–491. 2, 5
- [HJA20] HO J., JAIN A., ABBEEL P.: Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851. 3
- [HJCX21] HAN X.-F., JIN Y.-F., CHENG H.-X., XIAO G.-Q.: Dual transformer for point cloud analysis. *arXiv preprint arXiv:2104.13044* (2021). 2
- [HXX*20] HUI L., XU R., XIE J., QIAN J., YANG J.: Progressive point cloud deconvolution generation network. In *European Conference on Computer Vision* (2020), Springer, pp. 397–413. 2
- [JST*21] JIANG L., SHI S., TIAN Z., LAI X., LIU S., FU C.-W., JIA J.: Guided point contrastive learning for semi-supervised point cloud semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 6423–6432. 1
- [KBV20] KLOKOV R., BOYER E., VERBEEK J.: Discrete point flow networks for efficient point cloud generation. In *European Conference on Computer Vision* (2020), Springer, pp. 694–710. 2
- [KHY*18] KUO W., HÄNE C., YUH E., MUKHERJEE P., MALIK J.: Cost-sensitive active learning for intracranial hemorrhage detection. In *International Conference on Medical Image Computing and Computer-Assisted Intervention* (2018), Springer, pp. 715–723. 5
- [KLK*20] KIM H., LEE H., KANG W. H., LEE J. Y., KIM N. S.: Soft-flow: Probabilistic framework for normalizing flow on manifolds. *Advances in Neural Information Processing Systems* 33 (2020), 16388–16397. 2
- [KYLH21] KIM J., YOO J., LEE J., HONG S.: Setvae: Learning hierarchical composition for generative modeling of set-structured data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 15059–15068. 3
- [LCL18] LI J., CHEN B. M., LEE G. H.: So-net: Self-organizing network for point cloud analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2018), pp. 9397–9406. 2
- [LH21] LUO S., HU W.: Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 2837–2845. 1, 2, 5
- [LKX*21] LYU Z., KONG Z., XU X., PAN L., LIN D.: A conditional point diffusion-refinement paradigm for 3d point cloud completion. *arXiv preprint arXiv:2112.03530* (2021). 1, 2
- [MBST21] MESIKA A., BEN-SHABAT Y., TAL A.: Cloudwalker: 3d point cloud learning by random walks for shape analysis. *arXiv preprint arXiv:2112.01050* (2021). 2
- [MS15a] MATURANA D., SCHERER S.: 3d convolutional neural networks for landing zone detection from lidar. In *2015 IEEE international conference on robotics and automation (ICRA)* (2015), IEEE, pp. 3471–3478. 2
- [MS15b] MATURANA D., SCHERER S.: Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (2015), IEEE, pp. 922–928. 2
- [MWWG20] MO K., WANG H., YAN X., GUIBAS L.: Pt2pc: Learning to generate 3d point cloud shapes from part tree conditions. In *European Conference on Computer Vision* (2020), Springer, pp. 683–701. 1
- [QLJ*17] QI X., LIAO R., JIA J., FIDLER S., URTASUN R.: 3d graph neural networks for rgb-d semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision* (2017), pp. 5199–5208. 2
- [QSMG17] 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* (2017), pp. 652–660. 2, 5
- [QYSG17] QI C. R., YI L., SU H., GUIBAS L. J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems* 30 (2017). 2, 5, 6, 7
- [QYW*19] QIN C., YOU H., WANG L., KUO C.-C. J., FU Y.: Pointdan: A multi-scale 3d domain adaption network for point cloud representation. *Advances in Neural Information Processing Systems* 32 (2019). 1
- [ROUG17] RIEGLER G., OSMAN ULUSOY A., GEIGER A.: Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 3577–3586. 2
- [SPK19] SHU D. W., PARK S. W., KWON J.: 3d point cloud generative adversarial network based on tree structured graph convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 3859–3868. 1, 2
- [SWL*20] SUN Y., WANG Y., LIU Z., SIEGEL J., SARMA S.: Pointgrow: Autoregressively learned point cloud generation with self-attention. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (2020), pp. 61–70. 2
- [TQD*19] THOMAS H., QI C. R., DESCHAUD 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* (2019), pp. 6411–6420. 2
- [VSP*17] VASWANI A., SHAZEER N., PARMAR N., USZKOREIT J., JONES L., GOMEZ A. N., KAISER Ł., POLOSUKHIN I.: Attention is all you need. *Advances in neural information processing systems* 30 (2017). 2
- [WCG*19] WANG Y., CHAO W.-L., GARG D., HARIHARAN B., CAMPBELL M., WEINBERGER K. Q.: Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2019), pp. 8445–8453. 1
- [WHH*19] WANG L., HUANG Y., HOU Y., ZHANG S., SHAN J.: Graph attention convolution for point cloud semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2019), pp. 10296–10305. 2
- [WQF19] 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* (2019), pp. 9621–9630. 2

- [WSK*15] 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* (2015), pp. 1912–1920. [2](#)
- [WSL*19a] 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 (2019), 1–12. [2](#)
- [WSL*19b] 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 (2019), 1–12. [5](#)
- [XZS*21] XIANG T., ZHANG C., SONG Y., YU J., CAI W.: Walk in the cloud: Learning curves for point clouds shape analysis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 915–924. [2](#)
- [YFST18] YANG Y., FENG C., SHEN Y., TIAN D.: Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2018), pp. 206–215. [6](#)
- [YHH*19] YANG G., HUANG X., HAO Z., LIU M.-Y., BELONGIE S., HARIHARAN B.: Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 4541–4550. [2](#), [6](#)
- [YKH*19] YAN X., KHANSARI M., HSU J., GONG Y., BAI Y., PIRK S., LEE H.: Data-efficient learning for sim-to-real robotic grasping using deep point cloud prediction networks. *arXiv preprint arXiv:1906.08989* (2019). [1](#)
- [YTR*21] 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). [2](#)
- [YWZJ21] YANG X., WU Y., ZHANG K., JIN C.: Cpcgan: A controllable 3d point cloud generative adversarial network with semantic label generating. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2021), vol. 35, pp. 3154–3162. [1](#), [2](#), [5](#), [6](#), [7](#)
- [YZK21] YIN T., ZHOU X., KRÄHENBÜHL P.: Multimodal virtual point 3d detection. *Advances in Neural Information Processing Systems 34* (2021). [1](#)
- [ZBDT19] ZHAO Y., BIRDAL T., DENG H., TOMBARI F.: 3d point capsule networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2019), pp. 1009–1018. [6](#)
- [ZDW21] ZHOU L., DU Y., WU J.: 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 5826–5835. [2](#)
- [ZJFJ19] 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* (2019), pp. 5565–5573. [2](#)
- [ZJJ*21] ZHAO H., JIANG L., JIA J., TORR P. H., KOLTUN V.: Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 16259–16268. [2](#)
- [ZLG*21] ZHANG Y., LING H., GAO J., YIN K., LAFLECHE J.-F., BARRIUSO A., TORRALBA A., FIDLER S.: Datasetgan: Efficient labeled data factory with minimal human effort. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 10145–10155. [1](#), [5](#)
- [ZWL*21] ZHANG C., WAN H., LIU S., SHEN X., WU Z.: Pvt: Point-voxel transformer for 3d deep learning. *arXiv preprint arXiv:2108.06076* (2021). [2](#)