# SSEditor: Controllable Mask-to-Scene Generation with Diffusion Model

Haowen Zheng, YanyanLiang
Macau University of Science and Technology
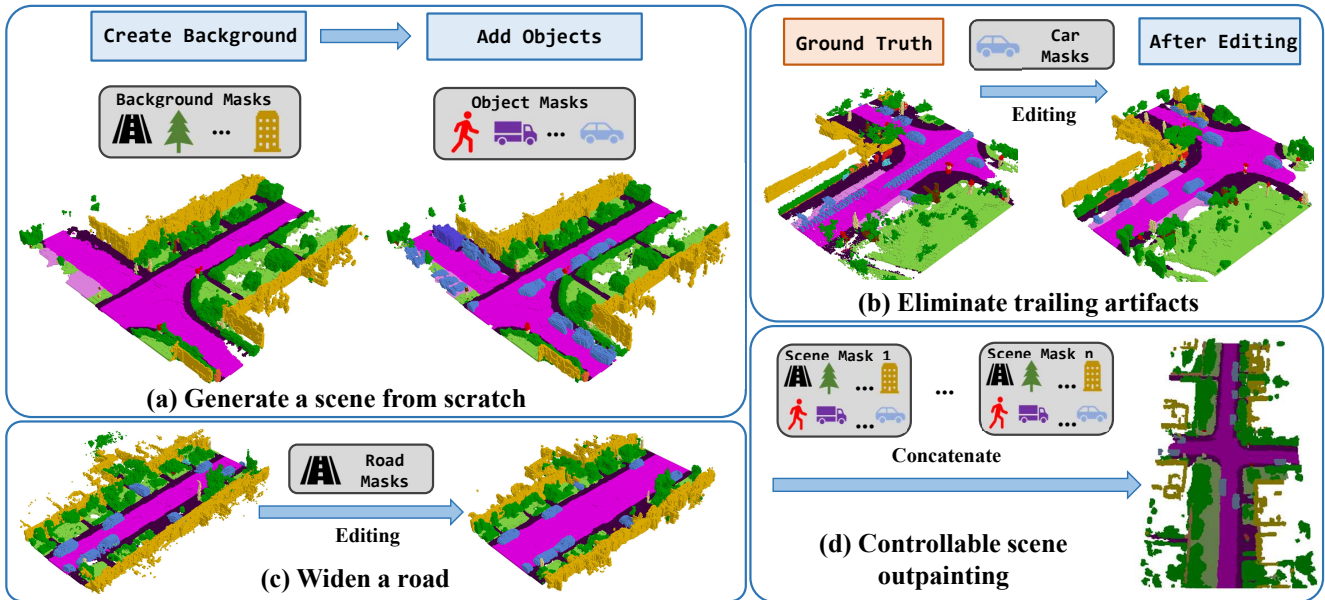zhengnayin@gmail.com, yyliang@must.edu.mo

Figure 1. Controllable 3D semantic scene generation by SSEditor. The proposed SSEditor enables users to customize the generation or editing of 3D scenes using pre-built mask assets: (a) create a background scene and generate objects on it; (b) eliminate trailing artifacts of dynamic objects in SemanticKITTI [2]; (c) modify roads, such as expanding a two-lane road to a four-lane road; (d) concatenate masks from various scenes to produce a larger-scale 3D scene.

## Abstract

*Recent advancements in 3D diffusion-based semantic scene generation have gained attention. However, existing methods rely on unconditional generation and require multiple resampling steps when editing scenes, which significantly limits their controllability and flexibility. To this end, we propose **SSEditor**, a controllable **S**emantic **S**cene **Editor** that can generate specified target categories without multiple-step resampling. SSEditor employs a two-stage diffusion-based framework: (1) a 3D scene autoencoder is trained to obtain latent triplane features, and (2) a mask-conditional diffusion model is trained for customizable 3D semantic scene generation. In the second stage, we introduce a geometric-semantic fusion module that enhance the model's ability to learn geometric and semantic information. This ensures that objects are generated with cor-*
*rect positions, sizes, and categories. Extensive experiments on SemanticKITTI and CarlaSC demonstrate that SSEditor outperforms previous approaches in terms of controllability and flexibility in target generation, as well as the quality of semantic scene generation and reconstruction. More importantly, experiments on the unseen Occ-3D Waymo dataset show that SSEditor is capable of generating novel urban scenes, enabling the rapid construction of 3D scenes.*

## 1. Introduction

In recent years, 3D diffusion models have made notable achievements in generating both indoor [13, 32, 40] and outdoor [15, 16, 19, 26, 35] environments, as well as a single object [14, 31, 43]. Compared to indoor scenes and individual objects, outdoor scenes present more challenges due to their sparser and more complex representations. For

1

instance, voxel-based representations of outdoor environments often contain a significant number of empty voxels. Moreover, outdoor environments contain smaller targets, such as pedestrians and cyclists, further complicating the generation process. While voxel-based representations [15, 19, 26, 35] provide a straightforward approach to modeling 3D semantic scenes, they suffer from redundancy in empty regions and high computational cost. To mitigate these issues, the triplane representation [5] is utilized to reduce unnecessary information in 3D outdoor scenes [16]. Although these methods have shown promising results, they still face several limitations.

The primary limitation lies in their weak controllability. Unconditional generation restricts the ability to guide the creation of 3D scenes, while conditioning on the entire scene (e.g., scene refinement based on ground truth) is overly rigid. This lack of flexible control leads to another drawback: editing specific local regions, such as adding or removing objects, necessitates masking non-target areas and employing a multi-step resampling process for repainting [20]. It significantly increases generation time. Despite the use of this resampling strategy, repainting remains uncontrollable and often fails to produce the desired results.

To address the aforementioned challenges, we propose SSEditor, a flexible and controllable two-stage framework for semantic scene generation based on the latent diffusion model (LDM) [28]. In the first stage, we train a 3D scene autoencoder to learn triplane features via semantic scene reconstruction. In the second stage, we train a mask conditional diffusion model on the triplane features. Specifically, to enable the customizable generation of 3D semantic scenes, we present a Geometric-Semantic Fusion Module (GSFM), which consists of a geometric branch and a semantic branch. The geometric branch encodes 3D masks that represent an object's position, size, and orientation, while the semantic branch processes semantic labels and tokens for providing coarse and fine-grained semantic information. The semantic tokens are generated from the features of a specific category. These features are then aggregated and integrated into the cross-attention module of the diffusion model, enhancing its perception of both geometric and semantic information. Benefiting from the above design, SSEditor effectively accomplishes the mask-to-semantic scene generation task.

In addition, we create a 3D mask asset library encompassing various categories to facilitate custom scene generation during inference. The 3D masks in the library are stored in the form of trimasks, which are composed of three orthogonal 2D planes derived from the decomposition of the 3D mask. As shown in Fig. 1, users can choose from a range of assets, such as cross-shaped roads, vehicles, pedestrians, and cyclists, to generate their desired 3D semantic scenes. The assets can also be edited to simulate more urban sce-

narios, such as expanding a two-lane road to four or more lanes.

Our contributions can be summarized into three points:
- We propose SSEditor, a controllable mask-to-scene generation framework that enables users to easily customize and generate 3D semantic scenes using various assets.
- We propose GSFM to integrate geometric and semantic information. In GSFM, the geometric branch encodes 3D masks as embeddings to accurately control the position, size, and orientation of objects, while the semantic branch processes semantic labels and tokens for improved class control of the generated targets.
- Experiments on outdoor datasets demonstrate that our proposed method achieves superior generation quality and reconstruction performance. Furthermore, qualitative results indicate that SSEditor can controllably perform various downstream tasks, such as scene inpainting, resource expansion, novel urban scene generation, and removal of trailing artifacts.

## 2. Related Work

**Controllable Diffusion Models.** Denoising diffusion probabilistic models (DDPM) [11] inspires a series of diffusion-based controllable generation approaches. Text-guided image generation shows strong capabilities in image editing tasks, such as inpainting [1, 24, 25] and outpainting [29]. In addition, several studies incorporate more control signals, such as layouts [42], semantic maps [8, 28, 36, 41], to facilitate image generation. Building on these advancements, controllable diffusion models have been further extended to the 3D domain. These models can leverage images [6, 39], text [17, 21], partial point clouds [23] or multi-modal conditions (e.g., text-image or text-voxels) [22, 34] to guide the generation of a single 3D object. However, the aforementioned controllable generative models can only be applied to 2D images or individual 3D objects, making it challenging for them to handle complex large-scale 3D scenes.

**3D Semantic Scene Generation.** 3D semantic scene generation can be categorized into indoor and outdoor scene generation. CommonScenes [40] generates indoor scenes based on scene graphs. DiffuScene [32] performs indoor scene generation and completion based on a text prompt or incomplete 3D targets. InstructScene [18] incorporates user instructions into semantic graph priors and decodes them into 3D indoor scenes. Build-A-Scene [7] enables users to flexibly create indoor scenes by adjusting layouts. In contrast, outdoor scene generation is more complex, which features diverse objects, more occlusions, and varying distances. [15] generates 3D multi-object scenes in simulated outdoor environments, while PDD [19] employs a coarse-to-fine strategy to further improve generation quality. For more complex real-world outdoor scenes, SemCity [16] uses triplane diffusion to achieve unconditional generation
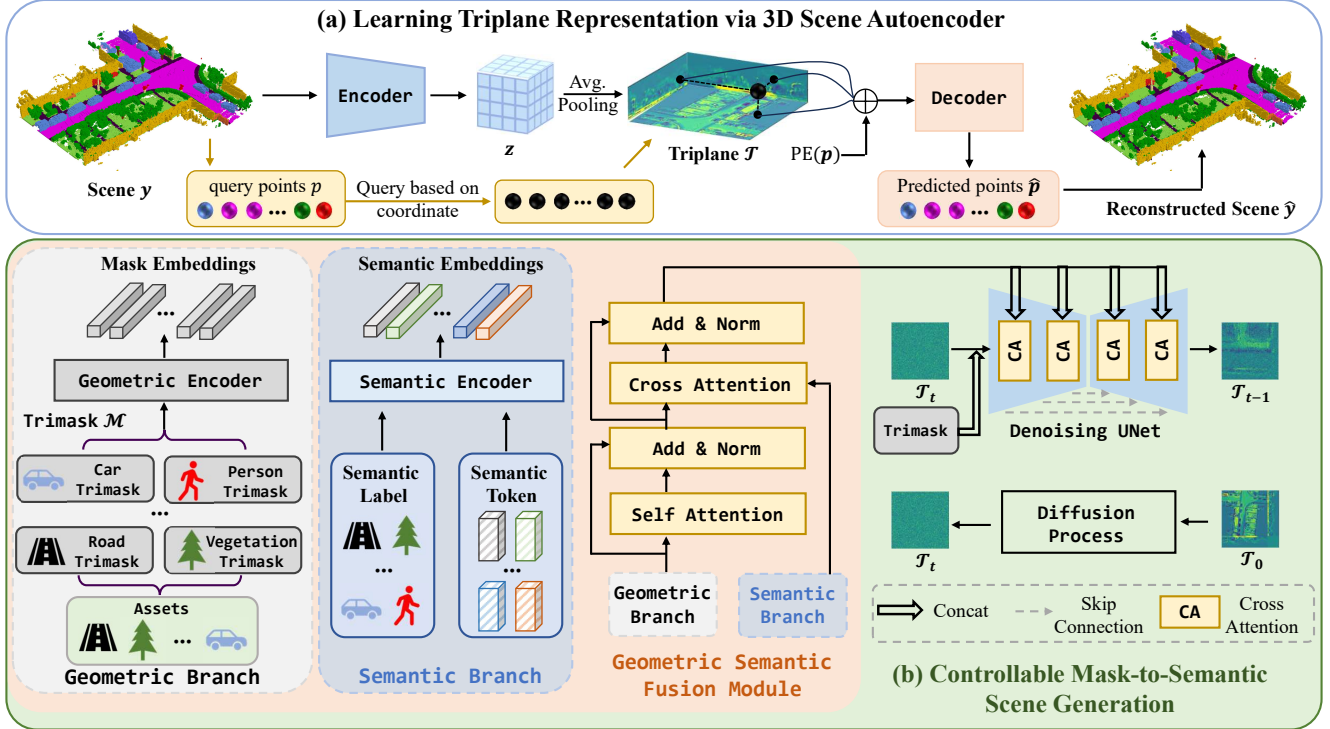
Figure 2. Illustration of our SSEditor framework. It comprises two main processes: (a) a 3D autoencoder learns the triplane representation via scene reconstruction, and (b) controllable semantic scene generation is achieved through masks, semantic labels, and tokens. The Geometric-Semantic Fusion Module is essential for the diffusion model to effectively learn both geometric and semantic information.

or conditional 3D occupancy refinement.

Due to the significant differences between indoor and outdoor environments, these controllable indoor scene generation methods [7, 18, 32] are difficult to apply to outdoor scenes. For outdoor environments, [16, 19] can only refine scenes by conditionally inputting the entire 3D layout. Moreover, when conducting scene inpainting, SemCity [16] requires multiple-step resampling [20] and lacks one-step sampling capability. Additionally, it can not control the categories of the generated regions. This lack of flexible control prevents users from generating their desired scenes. In this paper, our proposed SSEditor overcomes these limitations and enables users to generate large-scale outdoor scenes from masks with traditional DDPM sampling [11].

## 3. Method

In this paper, we propose our SSEditor, as illustrated in Fig. 2. The primary objective of SSEditor is to enable users to generate 3D outdoor semantic scenes with flexibility and controllability. To achieve this goal, we first leverage a 3D scene autoencoder to learn the triplane representation (Sec. 3.1) and then create an asset library for storing 3D masks (Sec. 3.2). To enhance the accuracy for generating the positions, sizes, and categories of target objects, we implement a geometry-semantic fusion module that improves the model's understanding of geometric and semantic information, facilitating our controllable mask-to-scene generation. (Sec. 3.3). During inference, users can flexibly select or create assets to customize 3D scene construction, such as controllable inpainting, novel urban scene generation and trailing artifacts removal (Sec. 3.4).

### 3.1. 3D Scene Autoencoder with Triplane

Fig. 2(a) illustrates that the 3D scene autoencoder learns the triplane representation through scene reconstruction. We employ an encoder composed of 3D convolutions to encode a given scene $\mathbf{y} \in \mathbb{R}^{X \times Y \times Z}$ into $\mathbf{z} \in \mathbb{R}^{C_z \times \frac{X}{d} \times \frac{Y}{d} \times \frac{Z}{d_z}}$, where $C_z$, $X$, $Y$ and $Z$ denote the number of channel and the resolution of 3D voxel space, while $d$ and $d_z$ indicate the down-sampling factors. Axis-wise average pooling is then applied across the three dimensions of $\mathbf{z}$ to derive the triplane representation $\mathcal{T} = [\mathcal{T}^{xy}, \mathcal{T}^{xz}, \mathcal{T}^{yz}]$. In addition, we sample query points $\mathbf{p}$ from the scene voxels and aggregate the corresponding triplane features based on their coordinates, which can be represented as $\mathcal{T}(\mathbf{p}) = \mathcal{T}^{xy}(\mathbf{p}^{xy}) + \mathcal{T}^{xz}(\mathbf{p}^{xz}) + \mathcal{T}^{yz}(\mathbf{p}^{yz})$. The aggregated triplane features, combined with positional embedding, are decoded to obtain the predicted points $\hat{\mathbf{p}}$. The predicted points re-
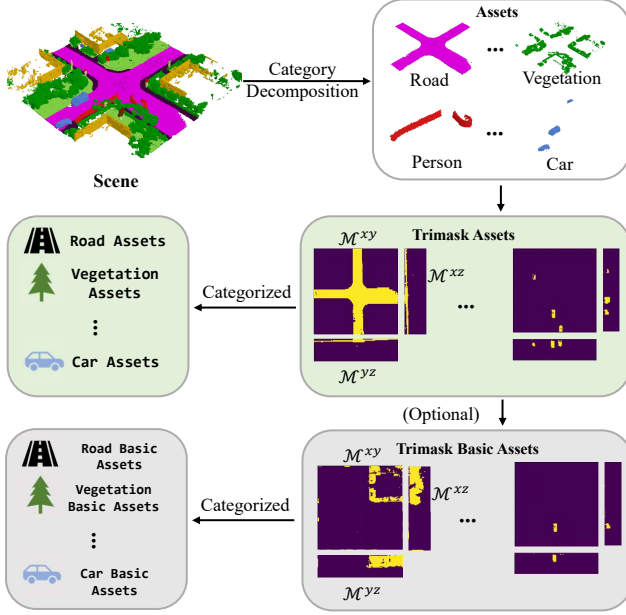
3

Figure 3. Pipeline of building 3D mask assets. The 3D mask is stored in the corresponding category in the form of a trimask.

construct the scene $\hat{y}$ based on the original coordinate information. The autoencoder is trained with a joint loss $\mathcal{L}_{AE}$, including the cross-entropy loss $\mathcal{L}_{CE}$ [27] on the points, and the Lovász-softmax loss $\mathcal{L}_{Lov}$ [3] on the reconstructed scene:

$$\mathcal{L}_{AE} = \mathcal{L}_{CE}(\hat{\boldsymbol{p}}, \mathbf{p}) + \alpha \mathcal{L}_{Lov}(\hat{\boldsymbol{y}}, \mathbf{y}) \tag{1}$$

where $\alpha$ is a loss weight.

## 3.2. 3D Mask Assets

To achieve a customizable generation of 3D scenes, controlling conditions need to be user-friendly inputs that can accurately reflect information such as target position and size. A 3D mask effectively serves this purpose. By utilizing the triplane representation, as illustrated in Fig. 3, we compress the 3D voxel mask into three 2D orthogonal planes, forming a trimask. The trimask can be represented as $\mathcal{M} = [\mathcal{M}^{xy}, \mathcal{M}^{xz}, \mathcal{M}^{yz}]$. All categories in the scene are decomposed into trimasks and stored in corresponding asset libraries. In addition to these scene-level assets, we also provide a basic version of the assets, which contains individual or segmented assets. This allows users to more conveniently utilize the basic assets to customize and construct scene-level assets. More importantly, users can also draw masks directly within the basic assets or scene-level assets. For example, the assets collected in the dataset only include small roads (2-lane and 4-lane). Users can edit the basic road assets (e.g., by copying, translating, or rotating)

to create wider lanes, such as 6-lane or 8-lane roads, to support the generation of more complex 3D scenes.

## 3.3. Controllable Mask-to-Scene Generation

The trimasks in the established assets offer valuable geometric information, including position, orientation, and scale. However, this is not enough for effective mask-to-semantic scene generation. We also need to extract detailed semantic information to ensure accurate object category generation. To tackle this, we propose a Geometric-Semantic Fusion Module (GSFM), as shown in Fig. 2(b), which consists of two branches: a geometric branch and a semantic branch.

**Geometric Branch.** The geometric branch encodes the trimask into mask embedding using an multi-layer perception (MLP), consisting of two linear layers and one activation layer. For simplicity, we first concatenate the trimask into a 2D feature maps $\mathcal{M}' \in \mathbb{R}^{N \times (X_m + Z_m) \times (Y_m + Z_m)}$, where N is the number of semantic classes, $X_m = \frac{X}{d}$, $Y_m = \frac{Y}{d}$ and $Z_m = \frac{Z}{d_z}$. The mask embedding $E_m \in \mathbb{R}^{N \times C_{emb}}$ can be obtained by

$$\text{MLP}(x) = \text{Linear}(\text{GeLU}(\text{Linear}(x) \tag{2}$$
$$E_m = \text{MLP}(\mathcal{M}') \tag{3}$$

The extracted mask embeddings currently operate independently and lack geometric information from other categories. To resolve this, we employ self-attention to capture the geometric relationships between masks of different categories through Eq. 4. This allows the model to detect targets of varying scales within the same category and identify overlapping regions between different category masks.

$$E'_m = E_m + \text{LayerNorm}(\text{SelfAttn}(E_m)). \tag{4}$$

**Semantic Branch.** In the semantic branch, we begin with an embedding layer to convert semantic labels into label embeddings $E_{label} \in \mathbb{R}^{N \times C_{emb}}$. However, the label embeddings offer only coarse-grained semantic information, which is inadequate for large-scale scene generation. The voxels generated within the mask regions may still contain a number of incorrect categories. To address this, we introduce a finer-grained semantic token $\mathbf{T}_{sem} \in \mathbb{R}^{N \times C_{emb}}$, which is defined as:

$$\mathbf{T}^i_{sem} = \text{Spatial Pooling}(\mathcal{M}_i \cdot \mathcal{T}) \tag{5}$$

where $i$ indicates the $i$-th semantic class and spatial pooling represents average pooling along the spatial dimension. As a result, the semantic embeddings $E_{sem} \in \mathbb{R}^{N \times C_{emb}}$ can be formulated as

$$E_{sem} = \text{MLP}(E_{label} + \mathbf{T}_{sem}) \tag{6}$$

4

| Model | FID ↓ | KID ↓ | IS ↑ | Prec. ↑ | Rec. ↑ |
|---|---|---|---|---|---|
| SemanticKITTI [2] | | | | | |
| SSD [15] | 117.46 | 0.12 | 2.15 ± 0.13 | 0.01 | 0.08 |
| SemCity [16] | 61.20 | 0.04 | 2.43 ± 0.11 | 0.19 | 0.12 |
| SSEditor (ours) | **47.93** | **0.03** | **2.55 ± 0.14** | **0.31** | **0.51** |
| CarlaSC [37] | | | | | |
| SSD [15] | 148.14 | 0.18 | 2.23 ± 0.10 | 0.15 | 0.01 |
| SemCity [16] | 137.94 | 0.13 | **3.03 ± 0.17** | 0.20 | 0.02 |
| SSEditor (ours) | **50.98** | **0.03** | 2.28 ± 0.08 | **0.37** | **0.18** |

Table 1. Quantitative results on SemanticKITTI and CarlaSC. The metrics are measured between the rendered image of the generated scene and the real scene. Prec. and Rec. indicates precision and recall, respectively.

**Geometric-Semantic Fusion Module.** The GSFM integrates mask embeddings $E_m$ and semantic embeddings $E_{sem}$ through cross-attention. We use the geometric embeddings as the query $Q \in \mathbb{R}^{N \times C_{emb}}$ and concatenate the geometric and semantic embeddings to form the key $K$ and value $V \in \mathbb{R}^{2N \times C_{emb}}$. The fused embeddings $E_{fused}$ can then be represented as:

$$E_{fused} = E'_m + \text{LayerNorm}(\text{CrossAttn}(Q, K, V)) \quad (7)$$

**Mask Conditional Diffusion Model.** Following LDM [28], we conduct diffusion and denoising process on the triplane features $\mathcal{T}$ to learn our mask conditional diffusion model $D_\theta$. We add t steps of Gaussian noise to a clean triplane features $\mathcal{T}_0$ and obtain a noised triplane $\mathcal{T}_t \sim q(\mathcal{T}_t|\mathcal{T}_0) = \mathcal{N}(\sqrt{\overline{\alpha}_t}\mathcal{T}, (1 - \overline{\alpha}_t)\mathbf{I})$, where $\mathcal{N}$ is the Gaussian distribution, $\overline{\alpha}_t = \prod_{i=1}^{t} \alpha_i$ and $\alpha_t = 1 - \beta_t$ with a variance schedule $\beta_t$. Then the diffusion model $D_\theta$ can be trained with the mean-squared error loss:

$$\mathcal{L} = \mathbb{E}_{t \sim [1,T]} \|\mathcal{T}_0 - D_\theta(\text{Concat}(\mathcal{T}_t, \mathcal{M}), t)\|_2 \quad (8)$$

To support mask conditional generation, we inject the fused embedding $E_{fused}$ obtained from Eq. 7 into the cross attention of the diffusion model. In addition, we concatenate the trimask with $\mathcal{T}_t$ to further enhance the guidance of the mask. Following classifier-free guidance [10], we randomly set the trimask to zero during training to simulate the effect of not using the trimask.

### 3.4. Downstream Applications

Unlike unconditional scene generation [16], our SSEditor can flexibly handle various downstream tasks based on the created assets, such as controllable scene inpainting and controllable scene outpainting. Note that our method does not require a resampling strategy [20].
**Controllable Scene Inpainting** can facilitate basic scene editing, such as adding or removing objects. Based on this,

| Model | Input | IoU ↑ | mIoU ↑ |
|---|---|---|---|
| Symphonies [12] | RGB | 41.92 | 14.89 |
| SCPNet [38] | Point Cloud | 50.24 | 37.55 |
| SSEditor (ours) | 3D Mask | **57.85** | **43.09** |

Table 2. Quantitative results on SemanticKITTI validation set. IoU and mIoU indicate how effectively our method handles geometric information and comprehends semantic information during generation, respectively.

SSEditor can simulate corner cases in autonomous driving scenarios, such as vehicle congestion at intersections, bicycles haphazardly parked on the roadside, and pedestrians crossing the street. Furthermore, the accumulation of multiple LiDAR frames causes trailing artifacts in dynamic objects within the SemanticKITTI dataset [2]. Our SSEditor effectively resolves this issue. In addition, by editing background assets such as roads and sidewalks, SSEditor can also widen roads to simulate scenarios with greater traffic.
**Controllable Scene Outpainting** can assist in scene extension. By selecting appropriate background assets and combining them, such as stitching together continuous roads, we can controllably extend the scene.
**Novel Urban Scene Generation** enables the rapid construction of 3D occupancy datasets. Imagine that we want to build a 3D semantic scene for a new city: we can create different assets based on LiDAR point clouds, and then generate a novel urban scene based on these assets.
**Removing trailing artifacts.** SemanticKITTI [2] aggregates multiple LiDAR frames to create dense 3D occupancy scenes, but this introduces trailing artifacts for moving objects in the ground truth, as shown in Fig. 1(b). Our method can effectively remove these artifacts and utilizes existing object assets to generate new objects.

## 4. Experiments

### 4.1. Datasets

We conduct our experiments on the SemanticKITTI [2] and CarlaSC [37] datasets. SemanticKITTI dataset is a large-scale real-world benchmark for semantic scene understanding in autonomous driving. It contains 20 semantic classes. Each scene is represented by a 256×256×32 voxel grid with a voxel resolution of 0.2m. CarlaSC dataset is a synthetic dataset with labels for 11 semantic classes, generated using the CARLA simulator. Each scene has a resolution of 128×128×8, covering an area of 25.6 meters around the vehicle, with a height of 3 meters. Additionally, we validated the cross-dataset transferability of SSEditor on Occ3D-Waymo [33]. We only included the occupancy labels from Occ3D-Waymo [33] as trimasks in our asset library and then simulated the generation of unknown
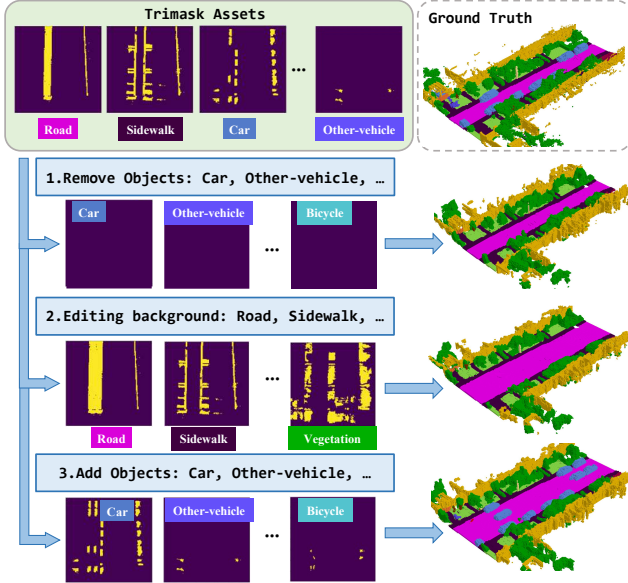
Figure 4. The details of editing 3D scenes with SSEditor: 1. When the mask of an object is set to 0, the corresponding object can be completely removed. 2. The background can be edited, such as widening roads to simulate heavier traffic. 3. Objects can be added to the edited scene.

urban scenes. Note that we disregard the Occ3D-Waymo categories not present in SemanticKITTI.

## 4.2. Implementation Details

All experiments are conducted on a single NVIDIA RTX 3090-24G GPU. For the 3D scene autoencoder, the batch size is set to 4, while for the controllable mask-to-scene generation, the batch size is set to 1. The downsampling factors are configured as $d = 2$ and $d_z = 1$. The loss weight $\alpha$ in the Eq. 1 is set to 1, the latent channel of triplane features $\mathcal{T}$ equals 16 and the embedding channel $C_{emb} = 64$. The learning rate for the autoencoder is 1e-3, while the learning rate for the diffusion model is 1e-4. Following the settings of [15, 16], the sampling time steps is set to 100 during both training and testing of the diffusion model. We utilize DDPM sampling strategy [11] for downstream tasks, omitting the need for the resampling strategy in RePaint [20].

## 4.3. Evaluation Metrics

We adopt evaluation metrics from prior works [16, 32, 40] rendering 3D scenes into 2D images and use traditional 2D evaluation metrics to assess the quality and diversity of generated scenes:

**Fréchet Inception Distance (FID)** [9] measures the similarity between the real and generated data distributions by comparing their feature statistics in the latent space of the ImageNet-pretrained Inception network.

**Inception Score (IS)** [30] evaluates both the quality and diversity of generated samples by computing a statistical score from the Inception network.

**Kernel Inception Distance (KID)** [4] computes the squared Maximum Mean Discrepancy (MMD) between the real and generated data distributions using features extracted from the Inception network.

**Precision** measures the proportion of generated samples that fall within the support of the real data distribution, while **Recall** measures the proportion of the real data distribution covered by the generated samples.

In addition, we use the intersection over union (**IoU**) and mean IOU (**mIoU**) metrics to evaluate the overall scene reconstruction quality and the reconstruction quality for each class, respectively.

## 4.4. Quantitative Results

**Generation.** Table 1 provides quantitative results on SmeanticKITTI and CarlaSC comparing with SSD [15] and SemCity [16]. In overall generation quality and diversity, our SSEditor outperforms the previous methods [15, 16] on SemanticKITTI [2], particularly in FID and recall, where we achieve improvements of 21.68% and 39%, respectively, compared to SemCity. On CarlaSC [37], SSEditor leads in all metrics except for IS, with FID improving by 63.04% over SemCity. Note that SemCity do not disclose which image sets are used for evaluation, making the results non-reproducible. To ensure a fair comparison, we train on the training set and generate scenes on the validation set to obtain the evaluation results.

**Semantic Scene Completion.** We assess the controllability and scene reconstruction capabilities of our method through semantic scene completion. Table 2 demonstrates that SSEditor performs well on the SemanticKITTI validation set. We only reference two state-of-the-art methods from different modalities, as other unconditional diffusion models [15, 16] lack the ability to reconstruct 3D semantic scenes. The IoU metric indicates that our method provides strong control over the position and size of objects during scene generation, while the mIoU score reflects a robust understanding of the semantics of the generated objects.

## 4.5. Qualitative Results

**Generation.** Fig. 5 showcases the qualitative results of the proposed SSEditor and SemCity [16] on the SemanticKITTI [2] and CarlaSC [37] datasets. While SemCity [16] effectively generates a variety of scenes using triplane representations, it lacks sufficient control, making scene customization challenging. In contrast, SSEditor allows for precise generation of 3D scenes guided by masks, offering enhanced controllability. In Fig. 5, we create trimasks based on ground truth to verify our method's controllability. The results demonstrate that SSEditor excels in
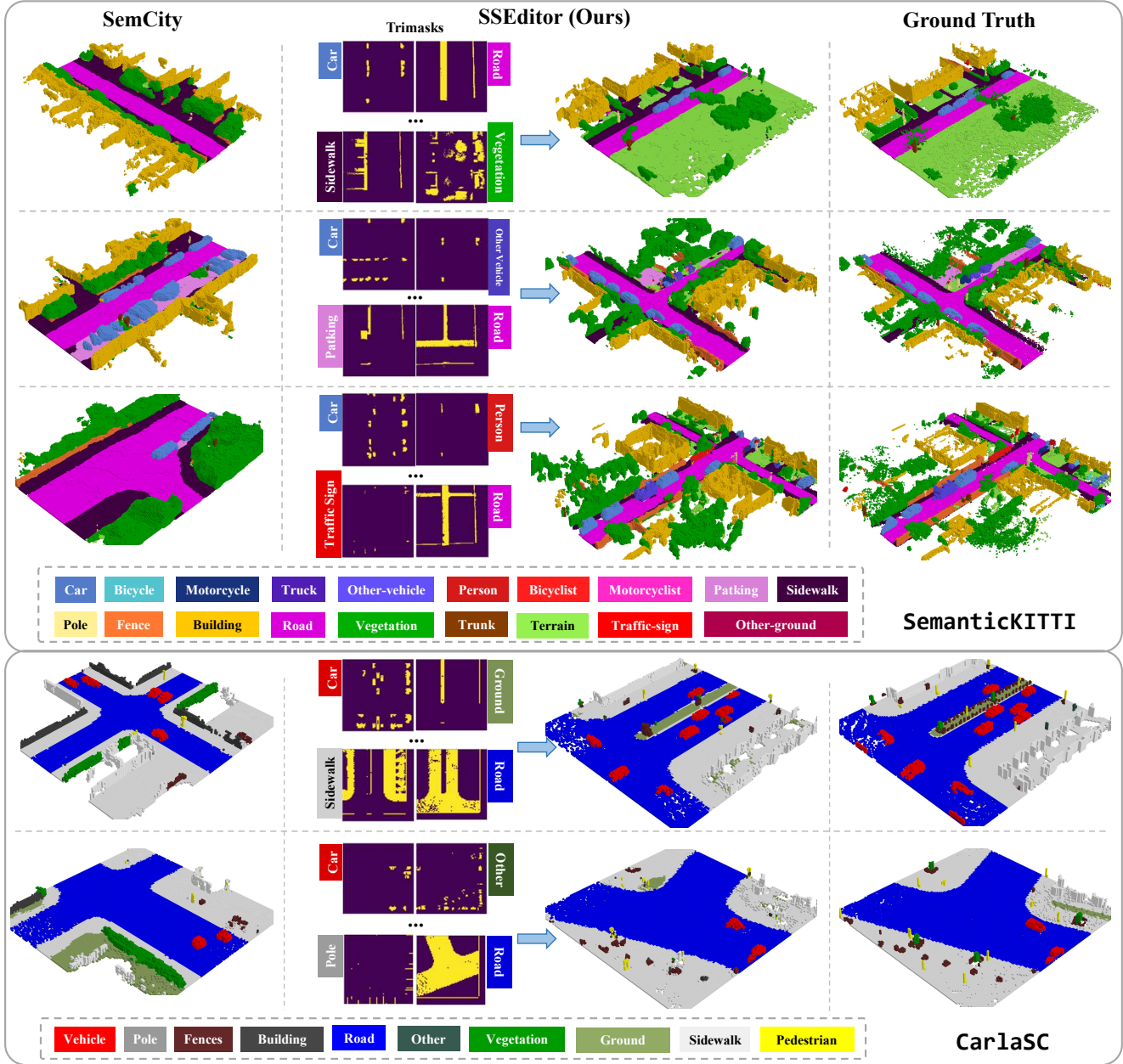
Figure 5. Visualization of semantic scene generation comparing with SemCity [16] on SemanticKITTI [2] and CalarSC [37]. Under the guidance of the trimask as a condition, SSEditor demonstrates its strong controllability.

controlling both the overall background (e.g., road, vegetation) and specific objects (e.g., vehicles, pedestrians).

**Scene Editing.** Fig. 4 highlights the details of scene editing with SSEditor. By setting the trimask of a target object or background to zero, we can effectively remove it from the scene. We can also edit background assets for more realistic scenarios, like creating four-lane or eight-lane assets. Once the background is adjusted, we can add objects, like increasing the number of cars to simulate higher traffic volumes, to create more dynamic scenarios.

**Novel Scene Generation.** To further validate the controllability of SSEditor in generating new scenes, we apply the trained model to the Occ-3D Waymo dataset [33]. We adjust the trimasks from Occ-3D Waymo through interpolation to align with the standard size of trimasks in our asset library, due to the different resolutions of the datasets. Note that we only create trimasks for categories that appear in SemanticKITTI [2]. The generated results in Fig. 6 demon-
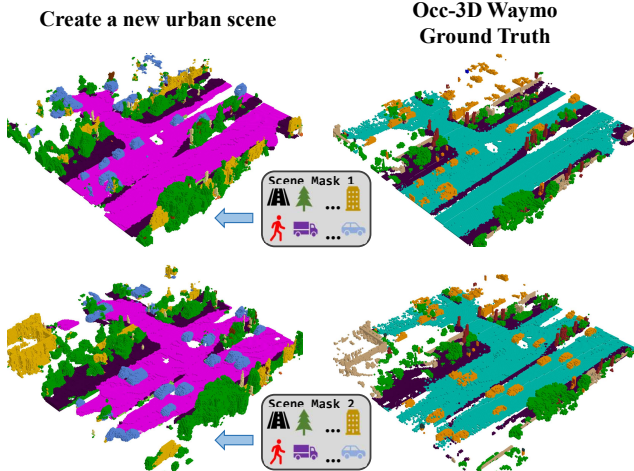
**Create a new urban scene**     **Occ-3D Waymo Ground Truth**

Figure 6. Create a novel urban scene from masks. The novel scene generation is tested on the unseen Occ-3D Waymo dataset [33].

| Method | FID ↓ | KID ↓ | IS ↑ | Prec. ↑ | Rec. ↑ |
|---|---|---|---|---|---|
| w/o geometric branch | 60.32 | 0.05 | 2.45 ± 0.15 | 0.24 | 0.28 |
| w/o semantic branch | 54.96 | 0.05 | 2.49 ± 0.13 | 0.27 | 0.37 |
| w/o semantic tokens | 53.67 | 0.04 | 2.49 ± 0.12 | 0.27 | 0.38 |
| w/o mask concat | 54.08 | 0.04 | 2.43 ± 0.17 | 0.23 | 0.19 |
| SSEditor (ours) | **47.93** | **0.03** | **2.55 ± 0.14** | **0.31** | **0.51** |

Table 3. Ablation studies on scene generation. We validated the effectiveness of the geometric branch, semantic branch, and the concatenated input of the trimask on SemanticKITTI [2].

| Method | Steps | Sampling | Inference Time |
|---|---|---|---|
| SSEditor | 100 | Resampling [20] | 56.44s |
| | | DDPM [11] | 13.40s |
| | 20 | Resampling [20] | 13.89s |
| | | DDPM [11] | 3.66s |
| | 10 | Resampling [20] | 6.91s |
| | | DDPM [11] | 2.08s |

Table 4. Ablation studies on sampling strategy. The inference time is reported based on 100 sample runs.

strate that SSEditor can effectively adapt to new scene generation, enabling the rapid creation of urban environments.

### 4.6. Ablation Studies

We conduct ablation experiments on the SemanticKITTI [2] validation set to assess the contribution of each component of SSEditor, as shown in Table 3.

First, we evaluate the effectiveness of the geometric branch by retaining the semantic branch and concatenating the trimask with the noised triplane $\mathcal{T}_t$ as input. Next, we remove the semantic branch, followed by the semantic tokens within the branch, to examine their individual impact. Finally, we input only the noised triplane $\mathcal{T}_t$ to assess the role of concatenating the trimask. In all ablation experiments, removing any component results in a performance drop, highlighting the necessity of each component for optimal performance.

Additionally, as shown in Table 4, we compared two sampling strategies: DDPM [11] and the resampling technique from RePaint [20]. While resampling improves object integration with the environment during generation, it greatly increases inference time for 3D scene generation. In contrast, our method employs traditional DDPM sampling, which maintains high quality and controllability in both scene inpainting and outpainting, while reducing inference time.

### 5. Limitations

Although SSEditor demonstrates strong capabilities for controllable scene generation, it still faces challenges with generating small objects, such as bicyclists and pedestrians. The generated areas sometimes contain incorrectly classified voxels, and the model's performance is highly sensitive to surrounding objects, which can lead to inaccuracies. These issues negatively affect the performance of downstream tasks that rely on high-quality scene generation. Predicting small objects in semantic scene completion is inherently challenging due to their low visibility and the complex interactions they have with the environment, resulting in lower mIoU performance. Future work could focus on addressing the long-tail distribution of data by incorporating more robust methods for representing and detecting small objects, as well as developing more fine-grained representation techniques that can improve the handling of these challenging cases.

### 6. Conclusion

In this paper, we propose SSEditor, a two-stage controllable scene generation framework based on the diffusion model. In the first stage, we leverage a 3D scene autoencoder to learn triplane representations. We then create a trimask asset library as a preparatory step for the second phase of training. In the second stage, we train a mask-conditional diffusion model for mask-to-scene generation, incorporating a geometric-semantic fusion module to enhance the extraction of geometric and semantic information. Experimental results on SemanticKITTI, CarlaSC, and Occ-3D Waymo demonstrate that our method outperforms existing unconditional diffusion approaches, offering superior controllability and high-quality scene generation. Moreover, SSEditor supports a wide range of applications, including the generation of novel 3D urban scenes (such as cross-dataset generation and road widening), controllable generation of dynamic objects, and scene outpainting.

# References

[1] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18208–18218, 2022. 2

[2] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Jurgen Gall. Semantickitti: A dataset for semantic scene understanding of lidar sequences. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9297–9307, 2019. 1, 5, 6, 7, 8

[3] Maxim Berman, Amal Rannen Triki, and Matthew B Blaschko. The lovász-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4413–4421, 2018. 4

[4] Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. *arXiv preprint arXiv:1801.01401*, 2018. 6

[5] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16123–16133, 2022. 2

[6] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VIII 14*, pages 628–644. Springer, 2016. 2

[7] Abdelrahman Eldesokey and Peter Wonka. Build-a-scene: Interactive 3d layout control for diffusion-based image generation. *arXiv preprint arXiv:2408.14819*, 2024. 2, 3

[8] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022. 2

[9] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6

[10] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 5

[11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2, 3, 6, 8

[12] Haoyi Jiang, Tianheng Cheng, Naiyu Gao, Haoyang Zhang, Tianwei Lin, Wenyu Liu, and Xinggang Wang. Symphonize 3d semantic scene completion with contextual instance queries. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20258–20267, 2024. 5

[13] Xiaoliang Ju, Zhaoyang Huang, Yijin Li, Guofeng Zhang, Yu Qiao, and Hongsheng Li. Diffindscene: Diffusion-based high-quality 3d indoor scene generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4526–4535, 2024. 1

[14] Animesh Karnewar, Andrea Vedaldi, David Novotny, and Niloy J Mitra. Holodiffusion: Training a 3d diffusion model using 2d images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18423–18433, 2023. 1

[15] Jumin Lee, Woobin Im, Sebin Lee, and Sung-Eui Yoon. Diffusion probabilistic models for scene-scale 3d categorical data. *arXiv preprint arXiv:2301.00527*, 2023. 1, 2, 5, 6

[16] Jumin Lee, Sebin Lee, Changho Jo, Woobin Im, Juhyeong Seon, and Sung-Eui Yoon. Semcity: Semantic scene generation with triplane diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28337–28347, 2024. 1, 2, 3, 5, 6, 7

[17] Muheng Li, Yueqi Duan, Jie Zhou, and Jiwen Lu. Diffusion-sdf: Text-to-shape via voxelized diffusion. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12642–12651, 2023. 2

[18] Chenguo Lin and Yadong Mu. Instructscene: Instruction-driven 3d indoor scene synthesis with semantic graph prior. *arXiv preprint arXiv:2402.04717*, 2024. 2, 3

[19] Yuheng Liu, Xinke Li, Xueting Li, Lu Qi, Chongshou Li, and Ming-Hsuan Yang. Pyramid diffusion for fine 3d large scene generation. *arXiv preprint arXiv:2311.12085*, 2023. 1, 2, 3

[20] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11461–11471, 2022. 2, 3, 5, 6, 8

[21] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13492–13502, 2022. 2

[22] Paritosh Mittal, Yen-Chi Cheng, Maneesh Singh, and Shubham Tulsiani. Autosdf: Shape priors for 3d completion, reconstruction and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 306–315, 2022. 2

[23] George Kiyohiro Nakayama, Mikaela Angelina Uy, Jiahui Huang, Shi-Min Hu, Ke Li, and Leonidas Guibas. Difffacto: Controllable part-based 3d point cloud generation with cross diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14257–14267, 2023. 2

[24] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 2

[25] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image gener-

ation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 2

[26] Xuanchi Ren, Jiahui Huang, Xiaohui Zeng, Ken Museth, Sanja Fidler, and Francis Williams. Xcube: Large-scale 3d generative modeling using sparse voxel hierarchies. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4209–4219, 2024. 1, 2

[27] Luis Roldao, Raoul de Charette, and Anne Verroust-Blondet. Lmscnet: Lightweight multiscale 3d semantic completion. In *2020 International Conference on 3D Vision (3DV)*, pages 111–119. IEEE, 2020. 4

[28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 5

[29] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 conference proceedings*, pages 1–10, 2022. 2

[30] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016. 6

[31] J Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d neural field generation using triplane diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20875–20886, 2023. 1

[32] Jiapeng Tang, Yinyu Nie, Lev Markhasin, Angela Dai, Justus Thies, and Matthias Nießner. Diffuscene: Denoising diffusion models for generative indoor scene synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20507–20518, 2024. 1, 2, 3, 6

[33] Xiaoyu Tian, Tao Jiang, Longfei Yun, Yucheng Mao, Huitong Yang, Yue Wang, Yilun Wang, and Hang Zhao. Occ3d: A large-scale 3d occupancy prediction benchmark for autonomous driving. *Advances in Neural Information Processing Systems*, 36, 2024. 5, 7, 8

[34] Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, Karsten Kreis, et al. Lion: Latent point diffusion models for 3d shape generation. *Advances in Neural Information Processing Systems*, 35:10021–10039, 2022. 2

[35] Lening Wang, Wenzhao Zheng, Yilong Ren, Han Jiang, Zhiyong Cui, Haiyang Yu, and Jiwen Lu. Occsora: 4d occupancy generation models as world simulators for autonomous driving. *arXiv preprint arXiv:2405.20337*, 2024. 1, 2

[36] Tengfei Wang, Ting Zhang, Bo Zhang, Hao Ouyang, Dong Chen, Qifeng Chen, and Fang Wen. Pretraining is all you need for image-to-image translation. *arXiv preprint arXiv:2205.12952*, 2022. 2

[37] Joey Wilson, Jingyu Song, Yuewei Fu, Arthur Zhang, Andrew Capodieci, Paramsothy Jayakumar, Kira Barton, and Maani Ghaffari. Motionsc: Data set and network for real-time semantic mapping in dynamic environments. *IEEE Robotics and Automation Letters*, 7(3):8439–8446, 2022. 5, 6, 7

[38] Zhaoyang Xia, Youquan Liu, Xin Li, Xinge Zhu, Yuexin Ma, Yikang Li, Yuenan Hou, and Yu Qiao. Scpnet: Semantic scene completion on point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 17642–17651, 2023. 5

[39] Qiangeng Xu, Weiyue Wang, Duygu Ceylan, Radomir Mech, and Ulrich Neumann. Disn: Deep implicit surface network for high-quality single-view 3d reconstruction. *Advances in neural information processing systems*, 32, 2019. 2

[40] Guangyao Zhai, Evin Pınar Örnek, Shun-Cheng Wu, Yan Di, Federico Tombari, Nassir Navab, and Benjamin Busam. Commonscenes: Generating commonsense 3d indoor scenes with scene graphs. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 2, 6

[41] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 2

[42] Guangcong Zheng, Xianpan Zhou, Xuewei Li, Zhongang Qi, Ying Shan, and Xi Li. Layoutdiffusion: Controllable diffusion model for layout-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22490–22499, 2023. 2

[43] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5826–5835, 2021. 1