# OccGen: Generative Multi-modal 3D Occupancy Prediction for Autonomous Driving

Guoqing Wang[1], Zhongdao Wang[2], Pin Tang[1], Jilai Zheng[1], Xiangxuan Ren[1], Bailan Feng[2], and Chao Ma[1]⋆

[1]MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University,
[2]Huawei Noah's Ark Lab
{guoqing.wang,pin.tang,zhengjilai,bunny_renxiangxuan,chaoma}@sjtu.edu.cn
{wangzhongdao,fengbailan}@huawei.com
Project page: https://occgen-ad.github.io/

**Abstract.** Existing solutions for 3D semantic occupancy prediction typically treat the task as a one-shot 3D voxel-wise segmentation perception problem. These discriminative methods focus on learning the mapping between the inputs and occupancy map in a single step, lacking the ability to gradually refine the occupancy map and the reasonable scene imaginative capacity to complete the local regions somewhere. In this paper, we introduce OccGen, a simple yet powerful generative perception model for the task of 3D semantic occupancy prediction. OccGen adopts a "noise-to-occupancy" generative paradigm, progressively inferring and refining the occupancy map by predicting and eliminating noise originating from a random Gaussian distribution. OccGen consists of two main components: a conditional encoder that is capable of processing multi-modal inputs, and a progressive refinement decoder that applies diffusion denoising using the multi-modal features as conditions. A key insight of this generative pipeline is that the diffusion denoising process is naturally able to model the coarse-to-fine refinement of the dense 3D occupancy map, therefore producing more detailed predictions. Extensive experiments on several occupancy benchmarks demonstrate the effectiveness of the proposed method compared to the state-of-the-art methods. For instance, OccGen relatively enhances the mIoU by 9.5%, 6.3%, and 13.3% on nuScenes-Occupancy dataset under the muli-modal, LiDAR-only, and camera-only settings, respectively. Moreover, as a generative perception model, OccGen exhibits desirable properties that discriminative models cannot achieve, such as providing uncertainty estimates alongside its multiple-step predictions.

**Keywords:** Occupancy · Generative Model · Diffusion · Multi-modal

## 1 Introduction

The precise 3D perception of the surrounding environment constitutes the cornerstone of modern autonomous driving systems, as it directly affects down-
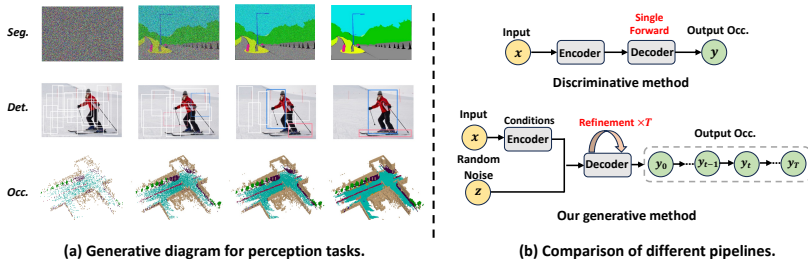
---

⋆ Corresponding author

**Fig. 1:** (a) The generative diagram of semantic segmentation (seg.), object detection (det.), and 3D semantic occupancy prediction (occ.). (b) Compared to previous discriminative methods with a single forward evaluation scheme, our OccGen is a generative model that can generate occupancy map in a coarse-to-fine manner.

stream tasks such as planning and vehicle control [17, 22]. In recent years, advancements in 3D object detection and segmentation [18,27,30,31,33,36,38,39, 55,56,64,66,67] have propelled the field of 3D perception. However, these methods require either rigid bounding boxes, which oversimplify the object shapes, or Bird's-Eye View (BEV) predictions that involve compromises in projecting 3D scenes onto 2D ground planes. Such methods can significantly impede the ability to accurately perceive structural information along the vertical axis, particularly when dealing with irregular objects.

To address this limitation, 3D semantic occupancy prediction [19,50,52,53,58] has been proposed to assign semantic labels to every spatially occupied region within the perceptive range. Most previous methods for 3D semantic occupancy prediction can be roughly divided into three categories: LiDAR-based [8,28,42, 61], vision-based [5,19,29,51,65], and multi-modal based [57] methods. These methods typically formulate the 3D occupancy prediction as a one-shot voxel-wise segmentation problem with a single forward evaluation scheme. While these works achieve promising results, this perception pipeline faces two critical issues: 1) Discriminative methods primarily focus on learning the mapping between the input-output pairs in a single forward step and neglect the modeling of the underlying occupancy map distribution. 2) Inferring only once is not enough for the model to complete the fine-grained scene well, just like humans need continuous observation to fully perceive the entire scene.

On the other hand, the diffusion model [16,49] has demonstrated its powerful generation capability and has also led to the successful application in numerous discriminative tasks, such as depth estimation [21,46], object detection [6], and segmentation [1,59,60]. We observe that the diffusion denoising process is naturally able to model the coarse-to-fine refinement of the dense 3D occupancy map, therefore producing more detailed predictions. Motivated by this, we propose OccGen, a simple yet powerful generative perception model for 3D semantic occupancy. As shown in Fig. 1, OccGen adopts a "noise-to-occupancy" generative paradigm, progressively inferring and eliminating noise originating from a random 3D Gaussian distribution. The proposed OccGen consists of two main components: a conditional encoder and a progressive refinement decoder. The

conditional encoder only needs to run once, while the decoder runs multiple times to fulfill progressive refinement. Since the encoder only runs once during the entire inference process, running the decoder step-by-step for diffusion denoising does not introduce significant computational overhead, thereby achieving comparable latency to single forward methods. During the training phase, we obtain a 3D noise map by gradually adding Gaussian noise to the ground truth occupancy. Subsequently, this noise map is fed into the progressive refinement decoder, which utilizes the multi-scale fusion features from the conditional encoder as conditions to generate noise-free predictions. In the inference phase, OccGen progressively generates the occupancy in a coarse-to-fine refinement manner, which is implemented by gradually denoising a 3D Gaussian noise map given the multi-modal condition inputs.

As a generative perception model, OccGen exhibits desirable properties that are not achievable by discriminative models: (1) progressive inference supports trading compute for prediction quality; (2) uncertainty estimation can be readily made alongside the predictions. We evaluate the effectiveness of OccGen on several benchmarks and show promising results compared with the state-of-the-art methods. Notably, OccGen has exhibited performance gains of 9.5%, 6.3%, and 13.3% on mIoU compared with the state-of-the-art method under the mulimodal, LiDAR-only, and camera-only settings on nuScenes-Occupancy.

Our contributions are summarized as follows:

- We introduce OccGen, a simple yet powerful generative framework following the "noise-to-occupancy" paradigm.
- OccGen adapts an efficient design that the encoder only runs once during the entire inference process, and the decoder runs step-by-step for progressive refinement, achieving a comparable latency to single forward methods.
- We extensively validate the proposed OccGen on multiple occupancy benchmarks, demonstrating its remarkable performance and desirable properties compared to previous discriminative methods.

## 2  Related Work

### 2.1  3D Semantic Occupancy Prediction.

The majority of popular 3D perception methods, whether the input is LiDAR sweeps, multi-view images, or multi-modal data, construct BEV feature representations and subsequently perform various downstream tasks in the BEV space [18, 30, 31, 33, 36]. However, these BEV-based methods typically project the 3D scene onto the ground plane, leading to the potential loss of information in the vertical dimension. Compared with BEV representation, 3D semantic occupancy provides a more detailed representation of the environment by explicitly modeling the occupancy status of each voxel in a 3D grid.

SSCNet [50] has first introduced the task of semantic scene completion, integrating both geometry and semantics. Subsequent works commonly utilized geometric inputs with explicit depth information [8, 28, 42, 61]. MonoScene [5] has

proposed the first monocular approach for semantic scene completion, employing a 3D UNet [44] to process voxel features generated through sight projection. TPVFormer [19] proposed a tri-perspective view representation for describing 3D scenes in semantic occupancy prediction. VoxFormer [29] introduced a two-stage transformer-based semantic scene completion framework that can output complete 3D volumetric semantics from only 2D images. OccFormer [65] introduced a dual-path transformer network for effective processing of 3D volumes in semantic occupancy prediction, achieving long-range, dynamic, and efficient encoding of camera-generated 3D voxel features. Furthermore, many concurrent works are dedicated to proposing surrounding-view benchmarks for 3D semantic occupancy prediction, contributing to the flourishing of the occupancy community [52, 53, 57, 58]. OpenOccupancy [57] constructed the first 3D multi-modal occupancy prediction benchmark and proposed an effective CONet to alleviate the computational burden of high-resolution occupancy predictions. In this paper, we propose OccGen, a simple yet powerful generative perception framework for 3D multi-modal semantic occupancy that can progressively refine the occupancy in a coarse-to-fine manner.

## 2.2   Diffusion Model

Diffusion models [16, 48] have been extensively researched due to their powerful generation capability. Denoising diffusion probabilistic models (DDPM) [16] proposed a diffusion model where the forward and reverse processes exhibit the Markovian property. Denoising diffusion implicit models (DDIM) [49] accelerated DDPM [16] by replacing the original diffusion process with non-Markovian chains to enhance sampling speed. On the other hand, conditional diffusion models have also been actively studied. Text-to-image generation models [43] and image-to-image translation models [45] achieved surprising results. Recently, diffusion models for visual perception have attracted widespread attention. Several pioneering works [1, 7, 47, 59, 60] attempted to apply the diffusion model to visual perception tasks, e.g. image segmentation or depth estimation tasks. For all the diffusion models listed above, one or two parameter-heavy convolutional U-Nets [44] are adopted, leading to low efficiency, slow convergence, and sub-optimal performance. DiffusionDet [6] proposed a denoising diffusion process from noisy boxes to object boxes. DDP [21] followed the "noise-to-map" generative paradigm for prediction by progressively removing noise from a random Gaussian distribution, guided by the image. In this work, as illustrated in Fig. 2, we extend the generative diffusion process into the occupancy perception pipeline while maintaining accuracy and efficiency.

## 3   Method

In this section, we first introduce the preliminaries on 3D semantic occupancy perception and conditional diffusion model. Then, we present the pipeline of the "noise-to-occupancy" and the overall architecture of OccGen. Finally, we show the details of the training and inference process.
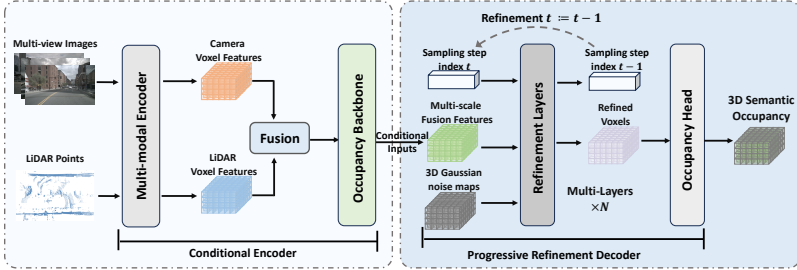
**Fig. 2:** The proposed OccGen framework. It has an encoder-decoder structure. The conditional encoder extracts the features from the inputs as the condition. The progressive refinement decoder consists of a stack of refinement layers and an occupancy head, which takes the 3D noise map, sampling step, and conditional multi-scale fusion features as inputs and progressively generates the occupancy prediction.

## 3.1  Preliminaries

**3D semantic occupancy perception.** The objective of 3D semantic occupancy perception is to predict a complete 3D representation of volumetric occupancy and semantic labels for a scene in the surround-view driving scenarios given inputs, such as images and LiDAR points. We utilize LiDAR point cloud $X_p \in \mathbb{R}^{N_L \times (3+d)}$ and multi-view camera images $X_c \in \mathbb{R}^{N_C \times H_C \times W_C \times 3}$ as multimodal inputs, denoted by $X = \{X_p, X_c\}$. Subsequently, we train a neural network $f_\theta$ to generate an occupancy voxel map $Y \in \{c_0, c_1, ..., c_N\}^{H \times W \times Z}$, where each voxel is assigned either an empty label $c_0$ or occupied by a specific semantic class from $\{c_1, c_2, ..., c_N\}$. Here, $N$ represents the total number of interested classes, and $\{H, W, D\}$ indicates the volumetric dimensions of the entire scene.

**Diffusion model.** The diffusion model is a type of generative model that demonstrates greater potential in the generative domain compared to Generative Adversarial Network (GAN) [12]. It can be divided into two categories: unconditional diffusion models learn an explicit approximation of the data distribution $P(z)$, while conditional diffusion models learn the distribution given a certain condition $k$, denoted as $p(z|k)$. In the conditional diffusion model, the data distribution is learned by recovering a data sample from Gaussian noise through an iterative denoising process. The forward diffusion process gradually adds noise to the data sample $z_0$, denoted as:

$$z_t = \sqrt{\alpha_t} z_0 + \sqrt{1 - \alpha_t} Z, \quad Z \sim \mathcal{N}(0, I), \tag{1}$$

which transforms the $z_0$ to a latent noisy sample $z_t$ for $t \in \{0, 1, \ldots, T\}$. The constant $\alpha_t = \prod_{i=1}^{t}(1 - \beta_i)$ and $\beta_s$ represents the noise schedule. In the training process, the conditional diffusion model $f_\theta(z_t, t | k)$ is trained to predict $z_0$ from $z_t$ under the guidance of condition $k$ by minimizing the training objective function (*i.e.*, $l_2$ loss). In the inference process, the predicted sample $z_0$ is reconstructed from a random noise $z_T$ with the model $f_\theta$ and conditional input $k$ following the denoising process of DDPM [16] or DDIM [49].

(a) The architecture of multi-modal encoder.      (b) The architecture of refinement layer.
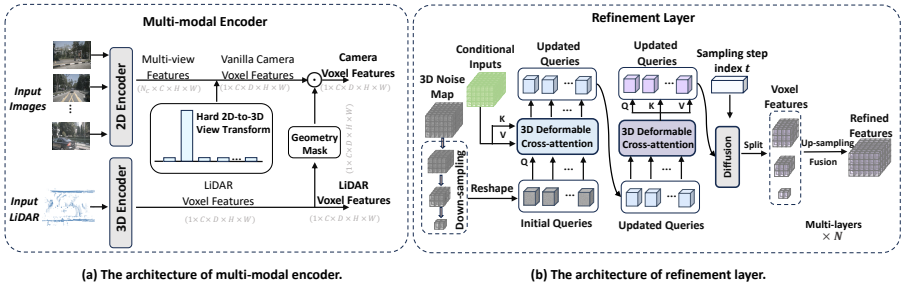
**Fig. 3:** The detailed architectures of (a) multi-modal encoder and (b) refinement layer. The multi-modal encoder is a two-stream structure, comprising LiDAR and camera streams. The refinement layer consists of three main components, i.e., 3D deformable cross-attention, self-attention, and time diffusion modules.

### 3.2  OccGen Framework

We first depict the proposed "noise-to-occupancy" generative paradigm and then introduce the overall architecture. As shown in Fig. 2, OccGen consists of a conditional encoder and a progressive refinement decoder.

**Noise-to-occupancy generative paradigm.** We regard the 3D semantic occupancy prediction as a generative process, which progressively generates the surrounding 3D environment with detailed geometry and semantics from single or multi-modal inputs. The goal of noise-to-occupancy is to learn an occupancy perception model $f_\theta$ which can model the coarse-to-fine refinement of the dense 3D occupancy map through a total of $T$ diffusion steps:

$$Y_T \xrightarrow{f_\theta} Y_{T-1} \xrightarrow{f_\theta} \dots \xrightarrow{f_\theta} Y_0, \tag{2}$$

where the diffusion step $T \to 0$ represents the coarse-to-fine refinement process from a 3D Gaussian voxel map to the refined occupancy. Thus, the generative occupancy prediction paradigm can be formulated as:

$$\Delta Y_t = f_\theta\left(x, t, Y_{t+1}\right), \quad Y_t = Y_{t+1} \oplus \Delta Y_t, \tag{3}$$

where the model $f_\theta$ refines the current prediction occupancy by giving the diffusion step index $t$ and the previous-step prediction occupancy $Y_{t+1}$, and $\oplus$ is element-wise summation.

**Conditional encoder.** The conditional encoder consists of three main components: a multi-modal encoder, a fusion module, and an occupancy backbone. As shown in Fig. 3 (a), the multi-modal encoder is a two-stream structure, comprising of LiDAR and camera streams. For the LiDAR stream, we follow VoxelNet [66] and 3D sparse convolutions [62] to transform raw LiDAR points to LiDAR voxel features. In the camera stream, we utilize the pre-trained 2D backbones [9, 15, 35, 54] and Feature Pyramid Network (FPN) [34] to extract multi-view image features given multi-view images. We obtain the vanilla camera voxel features through the 2D-to-3D view transformation.

Different from the previous 2D-to-3D view transformation [30, 33, 36, 41] methods that estimate the probabilistic of a set of discrete depths, OccGen proposes a hard 2D-to-3D view transformation to guarantee more accurate depth. We opt for predicting a one-hot vector for depth, as opposed to utilizing softmax on discrete depth values when lifting each image individually into a frustum of features for each camera. However, obtaining one-hot encoding directly through *argmax* operation is non-differentiable. To address this issue, we propose using Gumbel-Softmax [20] to convert the predicted depth into one-hot encoding.

The previous multi-modal methods for 3D occupancy prediction do not pay much attention to the interaction between multi-modal features. Therefore, we propose a straightforward solution to fully exploit the geometry-aware correspondence between camera and LiDAR modalities. We directly generate a geometry mask by leveraging LiDAR voxel features and then applying it to the vanilla camera voxel features to get the camera voxel features. This feature aggregation strategy effectively bridges the gap between the camera voxel features and the true spatial distribution in the real-world scene. We follow [57] and fuse the camera and LiDAR voxel features using the adaptive fusion module:

$$
\begin{aligned}
W &= \mathcal{G}_{\text{C}}\left(\left[\mathcal{G}_{\text{C}}\left(F_p\right), \mathcal{G}_{\text{C}}\left(F_c\right)\right]\right), \\
F_m &= \sigma(W) \odot F_p + (1 - \sigma(W)) \odot F_c,
\end{aligned}
\tag{4}
$$

where $\mathcal{G}_{\text{C}}$ is the 3D convolution, $[\cdot, \cdot]$ is the concatenation along feature channel. $\sigma$ and $\odot$ denote the Sigmoid function and element-wise product, respectively. Finally, we fed the multi-modal voxel features $F_m$ into the occupancy backbone to get the multi-scale fusion features for the following progressive refinement decoder. Additional details and experiments of hard 2D-to-3D view transformation and geometry mask are presented in the supplementary materials.

**Progressive Refinement Decoder.** The progressive refinement decoder of OccGen consists of a stack of refinement layers and an occupancy head. As illustrated in Fig. 3 (b), the refinement layer takes as input the random noise map or the predicted noise map $Y_{t+1}$ from the last step, the current sampling step $t$, and the multi-scale fusion features $F_m$. The refinement layer utilizes efficient 3D deformable cross-attention and self-attention to refine the 3D Gaussian noise map. Compared with traditional deformable attention [68] in 2D vision, 3D deformable attention samples the interest points around the reference point in the 3D pixel coordinate system to compute the attention results. Mathematically, 3D deformable attention can be represented by the following general equation:

$$
\text{DA}_{3\text{D}}(q, p, F) = \sum_{k=1}^{N} A_k W_k F(p + \Delta p_k),
\tag{5}
$$

where $q$ and $p$ denote the 3D query and 3D reference point, $F$ represents the flattened 3D voxel features, and $k$ indexes the sampled point from a total of $N$ points around the reference point $p$. $W_k$ represents the learnable weights for value generation, while $A_k$ corresponds to the learnable attention weight. $\Delta p_k$ denotes the predicted offset to the reference point $p$, and $F(p + \Delta p_k)$ signifies

the feature at the location $p + \Delta p_k$ extracted through bilinear interpolation. For brevity, we present the formulation for single-head attention only.

Directly operating on the original 3D Gaussian noise map $Y_t$ with high resolution is computationally intensive. Therefore, we first downsample it to obtain smaller multi-scale noise maps $Y_t^i \in \mathbb{R}^{\frac{D}{2^i} \times \frac{H}{2^i} \times \frac{W}{2^i} \times C_i} (i = 1, 2, 3)$. Then, we reshape these downsampled multi-scale noise maps to obtain initial queries. For each initial query $q$ in the multi-scale noise maps $Y_t^i$, we get the corresponding reference points $p$ on conditional inputs based on their corresponding spatial and level positions. We get the updated queries using 3D deformable cross-attention $(\mathrm{DCA_{3D}})$ by

$$\mathrm{DCA_{3D}}\left(Y_t^i, F_m\right) = \sum_{n \in F_m} \mathrm{DA_{3D}}\left(q, proj(q, n), F_m\right) \tag{6}$$

where $n$ denotes the hit multi-scale features. For each query $q$, we use $proj$ operation to obtain the reference point on multi-scale fusion features.

After one round of 3D deformable cross-attention, the initial queries gather knowledge from the condition inputs. To further enhance self-completion capability, we utilize the 3D deformable self-attention to update the queries,

$$\mathrm{DSA_{3D}}\left(Y_t^i, Y_t^i\right) = \sum_{n \in Y_t^i} \mathrm{DA_{3D}}\left(q, p, \mathbf{Y}_t^i\right). \tag{7}$$

Then, we split the learned queries into the down-sampled voxel sizes. We further apply a diffusion denoising step on the down-sampled multi-scale noise maps by

$$Y_t^i := \mathrm{Diff}(Y_t^i, \mathrm{ToEmbed}(t)), \tag{8}$$

where $\mathrm{ToEmbed}(\cdot)$ denote the embedding network that transfors a step index $t$ from scalar into a feature vector. $\mathrm{Diff}(\cdot)$ represents the diffusion module that applies the scale and shift operation along the time embedding. Furthermore, we upsample and project the downsampled voxels to the size of the original 3D noise map and obtain the refined voxel features. Finally, we obtain the 3D semantic occupancy by feeding the refined voxel features to the occupancy head. This process can be performed multiple times to progressively infer and refine the occupancy map by predicting and eliminating noise originating from a random Gaussian distribution.

### 3.3   Training

During training, we first construct a denoising diffusion process from the ground truth $Y_0$ to the 3D Gaussian noise map $Y_T$ and then train the progressive refinement decoder to reverse this process. Detailed information on the training procedure for OccGen is available in the supplementary materials.

**Occupancy Corruption.** We add Gaussian noise to corrupt the encoded ground truth, obtaining the 3D Gaussian noise map $Y_T$. As shown in Eq. 1, the intensity of corruption noise is controlled by $\alpha_t$, which follows a monotonically decreasing

schedule across different time steps $t \in [0,1]$. Different noise scheduling strate-
gies, including cosine schedule [40] and linear schedule [16], are compared and
discussed in the supplementary materials. We found that the cosine schedule
generally yields the best results in 3D semantic occupancy prediction.

**Loss Function.** The cross-entropy loss $\mathcal{L}_{ce}$ and lovasz-softmax loss $\mathcal{L}_{ls}$ [3] are
widely used to optimize the networks for semantic segmentation tasks. Following
[5,57], we also utilize affinity loss $\mathcal{L}_{scal}^{geo}$ and $\mathcal{L}_{scal}^{sem}$ to optimize the scene-wise and
class-wise metrics (*i.e.*, geometric IoU and semantic mIoU). Additionally, the
depth loss $\mathcal{L}_{d}$ [30] is used to optimize the predicted depth. Therefore, the overall
loss function can be derived as:

$$\mathcal{L}_{total} = \mathcal{L}_{ce} + \mathcal{L}_{ls} + \mathcal{L}_{scal}^{geo} + \mathcal{L}_{scal}^{sem} + \mathcal{L}_{d}. \tag{9}$$

### 3.4   Inference

Given multi-scale fusion features as conditional inputs, OccGen samples a ran-
dom noise map from a 3D Gaussian distribution and produces the occupancy
in a coarse-to-fine manner. The inference procedure for OccGen is provided in
supplementary materials.

**Sampling Rule.** Following [21], we choose the DDIM strategy [49] for the sam-
pling. In each sampling step $t$, the random noise map or the predicted noise
map from the last step and the conditional multi-scale fusion features are sent
to the progressive refinement decoder for occupancy prediction. After obtain-
ing the predicted result of the current step, we compute the refined noise map
for the next step using the reparameterization trick. Following [6,21], we use
the asymmetric time intervals (controlled by a hyper-parameter $td$) during the
inference stage. We empirically set $td = 1$ in our method.

**Progressively Inference.** According to the feature that the diffusion model
can generate the distribution step by step, we can perform progressive inference
to get fine-grained occupancy in a coarse-to-fine manner. Moreover, OccGen has
a natural awareness of the prediction uncertainty. As a comparison, previous
one-shot approaches for 3D semantic occupancy [29,57,58,65] can only output
a certain occupancy during the inference stage, and are unable to assess the
reliability and uncertainty of model predictions.

## 4   Experiments

### 4.1   Experimental Setup

**Dataset and Metrics.** We evaluate our proposed OccGen on two benchmarks,
i.e., nuScenes-Occupancy [57] and SemanticKITTI [2]. The nuScenes-Occupancy
extends the nuScenes [4] to provide dense annotations on keyframes for 3D multi-
modal semantic occupancy prediction. It covers 700 and 150 driving scenes in
the training and validation sets of nuScenes. SemanticKITTI [2] contains 22

**Table 1:** Semantic occupancy prediction results on nuScenes-Occupancy validation set. The $C, D, L, M$ denotes **camera, depth, LiDAR** and **multi-modal**. For **Surround**=✓, the method directly predicts surrounding semantic occupancy with 360-degree inputs. Best camera-only, LiDAR-only, and multi-modal results are marked red, blue, and **black**, respectively.

| Method | Input | Surround | IoU | mIoU | barrier | bicycle | bus | car | const. veh. | motorcycle | pedestrian | traffic cone | trailer | truck | drive. suf. | other flat | sidewalk | terrain | manmade | vegetation |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MonoScene [5] | C | ✗ | 18.4 | 6.9 | 7.1 | 3.9 | 9.3 | 7.2 | 5.6 | 3.0 | 5.9 | 4.4 | 4.9 | 4.2 | 14.9 | 6.3 | 7.9 | 7.4 | 10.0 | 7.6 |
| TPVFormer [19] | C | ✓ | 15.3 | 7.8 | 9.3 | 4.1 | 11.3 | 10.1 | 5.2 | 4.3 | 5.9 | 5.3 | 6.8 | 6.5 | 13.6 | 9.0 | 8.3 | 8.0 | 9.2 | 8.2 |
| 3DSketch [8] | C&D | ✗ | 25.6 | 10.7 | 12.0 | 5.1 | 10.7 | 12.4 | 6.5 | 4.0 | 5.0 | 6.3 | 8.0 | 7.2 | 21.8 | 14.8 | 13.0 | 11.8 | 12.0 | 21.2 |
| AICNet [28] | C&D | ✗ | 23.8 | 10.6 | 11.5 | 4.0 | 11.8 | 12.3 | 5.1 | 3.8 | 6.2 | 6.0 | 8.2 | 7.5 | 24.1 | 13.0 | 12.8 | 11.5 | 11.6 | 20.2 |
| LMSCNet [42] | L | ✓ | 27.3 | 11.5 | 12.4 | 4.2 | 12.8 | 12.1 | 6.2 | 4.7 | 6.2 | 6.3 | 8.8 | 7.2 | 24.2 | 12.3 | 16.6 | 14.1 | 13.9 | 22.2 |
| JS3C-Net [61] | L | ✓ | 30.2 | 12.5 | 14.2 | 3.4 | 13.6 | 12.0 | 7.2 | 4.3 | 7.3 | 6.8 | 9.2 | 9.1 | 27.9 | 15.3 | 14.9 | 16.2 | 14.0 | 24.9 |
| C-OpenOccupancy [57] | C | ✓ | 19.3 | 10.3 | 9.9 | 6.8 | 11.2 | 11.5 | 6.3 | 8.4 | 8.6 | 4.3 | 4.2 | 9.9 | 22.0 | 15.8 | 14.1 | 13.5 | 7.3 | 10.2 |
| L-OpenOccupancy [57] | L | ✓ | 30.8 | 11.7 | 12.2 | 4.2 | 11.0 | 12.2 | 8.3 | 4.4 | 8.7 | 4.0 | 8.4 | 10.3 | 23.5 | 16.0 | 14.9 | 15.7 | 15.0 | 17.9 |
| OpenOccupancy [57] | M | ✓ | 29.1 | 15.1 | 14.3 | 12.0 | 15.2 | 14.9 | 13.7 | 15.0 | 13.1 | 9.0 | 10.0 | 14.5 | 23.2 | 17.5 | 16.1 | 17.2 | 15.3 | 19.5 |
| C-CONet [57] | C | ✓ | 20.1 | 12.8 | 13.2 | 8.1 | 15.4 | 17.2 | 6.3 | 11.2 | 10.0 | 8.3 | 4.7 | 12.1 | 31.4 | 18.8 | 18.7 | 16.3 | 4.8 | 8.2 |
| L-CONet [57] | L | ✓ | 30.9 | 15.8 | 17.5 | 5.2 | 13.3 | 18.1 | 7.8 | 5.4 | 9.6 | 5.6 | 13.2 | 13.6 | 34.9 | 21.5 | 22.4 | 21.7 | 19.2 | 23.5 |
| CONet [57] | M | ✓ | 29.5 | 20.1 | 23.3 | 13.3 | 21.2 | 24.3 | 15.3 | 15.9 | 18.0 | 13.3 | 15.3 | 20.7 | 33.2 | 21.0 | 22.5 | 21.5 | 19.6 | 23.2 |
| C-OccGen | C | ✓ | 23.4 | 14.5 | 15.5 | 9.1 | 15.3 | 19.2 | 7.3 | 11.3 | 11.8 | 8.9 | 5.9 | 13.7 | 34.8 | 22.0 | 21.8 | 19.5 | 6.0 | 9.9 |
| L-OccGen | L | ✓ | 31.6 | 16.8 | 18.8 | 5.1 | 14.8 | 19.6 | 7.0 | 7.7 | 11.5 | 6.7 | 13.9 | 14.6 | 36.4 | 22.1 | 22.8 | 22.3 | 20.6 | 24.5 |
| OccGen | M | ✓ | 30.3 | 22.0 | 24.9 | 16.4 | 22.5 | 26.1 | 14.0 | 20.1 | 21.6 | 14.6 | 17.4 | 21.9 | 35.8 | 24.5 | 24.7 | 24.0 | 20.5 | 23.5 |

sequences including monocular images, LiDAR points, point cloud segentation labels and semantic scene completion annotations. We follow previous works [29, 57, 65] to report the Intersection of Union (IoU) as the geometric metric and the mean Intersection over Union (mIoU) of each class as the semantic metric. Besides, the noise class is excluded from the evaluation.

**Implementation Details.** We follow the same experiment settings of [57,65] to make a fair comparison with previous methods [5,29,57,65] on both nuScenes-Occupancy and SemanticKITTI. We simply stack six refinement layers with 3D deformable attention for the progressive refinement decoder. For training, we leverage the AdamW [24] optimizer with a weight decay of 0.01 and an initial learning rate of 2e-4. We adopt the cosine learning rate scheduler with linear warming up in the first 500 iterations, and a similar augmentation strategy as BEVDet [18]. The models are trained for 24 epochs with a batch size of 8 on 8 V100 GPUs. The implementation details on nuScenes-Occupancy and SemanticKITTI are listed in supplementary materials.

### 4.2   Comparison with the state-of-the-art

**Results on nuScenes-Occupancy.** As shown in Tab. 1, we report the quantitative comparison of existing LiDAR-based, camera-based, and multi-modal methods on nuScenes-Occupancy. Observations show that OccGen outperforms all existing competitors, regardless of whether the camera-only, LiDAR-only, or multi-modal methods. Compared with the current SOTA method CO-Net [57], OccGen achieves a remarkable boost of 1.7%, 0.4%, and 1.9% mIoU for camera-only, LiDAR-only, and multi-modal benchmarks, respectively. This demonstrates the effectiveness of OccGen for semantic occupancy prediction. We also note that

**Table 2: Semantic Scene Completion results on SemanticKITTI [2] validation set.** † denotes the results provided by MonoScene [5].

| Method | IoU | road. (%) | sidewalk. (%) | parking. (%) | otherground. (%) | building. (%) | car. (%) | truck. (%) | bicycle. (%) | motorcycle. (%) | othervehicle. (%) | vegetation. (%) | trunk. (%) | terrain. (%) | person. (%) | bicyclist. (%) | motorcyclist. (%) | fence. (%) | pole. (%) | trafficsign. (%) | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LMSCNet† [42] | 28.61 | 40.68 | 18.22 | 4.38 | 0.00 | 10.31 | 18.33 | 0.00 | 0.00 | 0.00 | 0.00 | 13.66 | 0.02 | 20.54 | 0.00 | 0.00 | 0.00 | 1.21 | 0.00 | 0.00 | 6.70 |
| AICNet† [28] | 29.59 | 43.55 | 20.55 | 11.97 | 0.07 | 12.94 | 14.71 | 4.53 | 0.00 | 0.00 | 0.00 | 15.37 | 2.90 | 28.71 | 0.00 | 0.00 | 0.00 | 2.52 | 0.06 | 0.00 | 8.31 |
| JS3C-Net† [61] | 38.98 | 50.49 | 23.74 | 11.94 | 0.07 | 15.03 | 24.65 | 4.41 | 0.00 | 0.00 | 6.15 | 18.11 | 4.33 | 26.86 | 0.67 | 0.27 | 0.00 | 3.94 | 3.77 | 1.45 | 10.31 |
| MonoScene [5] | 37.12 | 57.47 | 27.05 | 15.72 | **0.87** | 14.24 | 23.55 | 7.83 | 0.20 | 0.77 | 3.59 | 18.12 | 2.57 | 30.76 | 1.79 | 1.03 | 0.00 | 6.39 | 4.11 | 2.48 | 11.50 |
| TPVFormer [19] | 35.61 | 56.50 | 25.87 | 20.60 | 0.85 | 13.88 | 23.81 | 8.08 | 0.36 | 0.05 | 4.35 | 16.92 | 2.26 | 30.38 | 0.51 | 0.89 | 0.00 | 5.94 | 3.14 | 1.52 | 11.36 |
| VoxFormer [29] | 44.02 | 54.76 | 26.35 | 15.50 | 0.70 | 17.65 | 25.79 | 5.63 | 0.59 | 0.51 | 3.77 | 24.39 | **5.08** | 29.96 | 1.78 | 3.32 | 0.00 | **7.64** | 7.11 | 4.18 | 12.35 |
| OccFormer [65] | 36.50 | 58.85 | 26.88 | 19.61 | 0.31 | 14.40 | 25.09 | **25.53** | 0.81 | 1.19 | 8.52 | 19.63 | 3.93 | **32.62** | 2.78 | 2.82 | 0.00 | 5.61 | 4.26 | 2.86 | 13.46 |
| Symphonize [23] | 41.44 | 55.78 | 26.77 | 14.57 | 0.19 | **18.76** | **27.23** | 15.99 | **1.44** | 2.28 | 9.52 | **24.50** | 4.32 | 28.49 | 3.19 | **8.09** | 0.00 | 6.18 | **8.99** | **5.39** | 13.44 |
| **OccGen** (ours) | 36.87 | **61.28** | **28.30** | 20.42 | 0.43 | 14.49 | 26.83 | 15.49 | **1.60** | **2.53** | **12.83** | 20.04 | 3.94 | 32.44 | **3.20** | 3.37 | 0.00 | 6.94 | 4.11 | 2.77 | **13.74** |

that OccGen consistently delivers the best IoU results across almost all categories, which indicates that our method can better complete the scenes due to our coarse-to-fine generation property. It is also worth noting that OccGen with multi-modal inputs can improve camera-only and LiDAR-only by 7.5% and 5.8% mIoU, which demonstrates the effectiveness of the camera modality in capturing small objects (e.g., bicycle, pedestrian, motorcycle, traffic cone) and LiDAR modality on large objects structured regions (e.g., drivable surface, sidewalk, vegetation). This lays a solid foundation for us to further explore how to improve the role of images during fusion.

**Results on SemanticKITTI.** We also compare the proposed OccGen with the state-of-the-art vision-based works [23, 29, 65] on SemanticKITTI. For a fair comparison, we removed the LiDAR stream and fusion module from the conditional encoder. As shown in Tab. 2, we can see that OccGen achieves the highest mIoU compared with all existing competitors. Compared with the state-of-the-art OccFormer [65], our proposed method has an improvement of 0.3% mIoU, demonstrating the effectiveness of OccGen for semantic scene completion. We also notice that the transformer-based methods [10, 23, 29, 65] achieve higher performance than other previous methods. This reveals the superior capability of transformer-based structure in representation learning.

### 4.3 Ablation Study

**Overall architecture.** The ablation results on the conditional encoder and progressive refinement decoder are shown in Tab. 3. It is obvious that both the conditional encoder and progressive refinement decoder can achieve performance improvement. We also notice that "with proposed decoder" has a higher performance than "with proposed encoder", demonstrating the effectiveness of our generative pipeline.

**Progressive refinement decoder.** We conduct the ablations on the detailed components of the progressive refinement decoder. From Tab. 4 (a) and (b), it is

**Table 3:** Ablations on the conditional encoder and progressive refinement decoder on nuScenes-Occupancy under the multi-modal setting. (a), (b) and (c) denote our baseline, baseline "with proposed encoder" and "with proposed decoder", respectively.

|  | Encoder | Decoder | IoU | mIoU |
|---|---|---|---|---|
| (a) | - | - | 28.1 | 20.4 |
| (b) | ✓ | - | 28.6 | 20.7 |
| (c) | - | ✓ | 30.1 | 21.6 |
| (d) | ✓ | ✓ | 30.3 | 22.0 |

**Table 4:** Ablations on the progressive refinement decoder on nuScenes-Occupancy under the multi-modal setting. "**DCA**" and "**DSA**" denote the 3D deformable cross- and self-attention.

|  | Method | IoU | mIoU |
|---|---|---|---|
| (a) | w/o DSA | 30.1 | 21.4 |
|  | w/o CSA | 29.7 | 21.2 |
|  | w/o CSA and DSA | 29.1 | 20.7 |
| (b) | DSA + DCA | 29.4 | 21.6 |
|  | CSA + DSA | 30.3 | 22.0 |
| (c) | w/o Diffusion | 29.3 | 21.7 |
|  | OccGen | 30.3 | 22.0 |

evident that both 3D deformable cross- and self-attention lead to noticeable improvements in results. Compared to self-attention, cross-attention has a greater impact on performance, which is intuitive: learning knowledge from conditional inputs is always more comprehensive. Additionally, we also observed that the order of DCA and DSA in the decoder has a certain impact on the results. We also see that removing the temporal diffusion process leads to a decrease in results from Tab. 4 (c).

**Conditional encoder.** We also conduct the ablations on the detailed components of the conditional encoder. From Tab. 5, We also observe that the two solutions in the conditional encoder have both achieved promising performance. The reason is that the accurate depth estimation and geometry guidance can keep the fine-grained spatial structures. This effectively limits the impact of disruptive information from the images, leading to notable performance enhancements.

### 4.4 Further Discussion

The desirable properties of OccGen compared with the previous discriminative occupancy methods in a single-forward process are shown in Fig. 4 and 6. OccGen provides the flexibility to balance computational cost against prediction quality in a coarse-to-fine manner. Additionally, the stochastic sampling process enables the computation of voxel-wise uncertainty in the prediction.

**Progressive Refinement.** We evaluate OccGen with 1, 3, and 6 refinement layers by increasing their sampling steps from 1 to 10. The results are presented in Fig 5. It can be seen that OccGen can continuously improve its performance by using more sampling steps. For instance, OccGen with 6 refinement layers shows an increase from 21.7% mIoU (1 step) to 22.0% mIoU (3 steps), and we visualize the inference results of different steps in Fig. 4. In comparison to the previous single-step discriminative method, OccGen has the flexibility to balance computational cost against accuracy. This means our method can be adapted to different trade-offs between speed and accuracy under various scenarios without the need to retrain the network.

**Efficiency vs. Accuracy.** We report the results of IoU and mIoU to represent the accuracy of different methods and latency(ms) to represent the efficiency
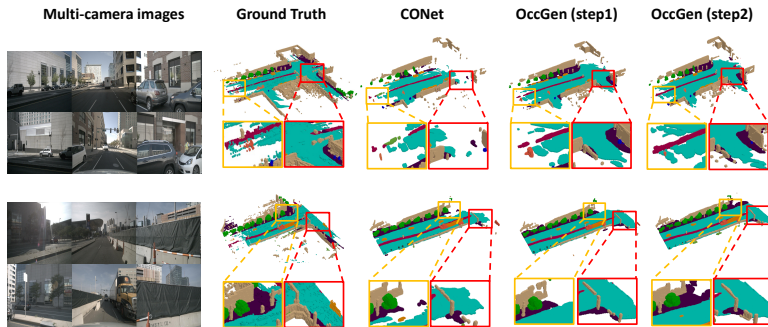
**Fig. 4:** Qualitative results of the 3D semantic occupancy predictions on nuScenes-Occupancy. The leftmost column shows the input surrounding images, and the following four columns visualize the 3D semantic occupancy results from the ground truth, CONet [57], OccGen(step1), and OccGen(step2). The regions highlighted by rectangles indicate that these areas have obvious differences (better viewed when zoomed in).

**Table 5:** Ablations on the multi-modal encoder on nuScenes-Occupancy under the multi-modal setting.'**Hard LSS**" and "**Geo. Mask**" denote hard 2D-to-3D view transformation and Geometry mask.

| Hard LSS | Geo. Mask | IoU | mIoU |
|:---:|:---:|:---:|:---:|
| - | - | 29.8 | 21.4 |
| ✓ | - | 30.2 | 21.5 |
| - | ✓ | 30.3 | 21.6 |
| ✓ | ✓ | 30.3 | 22.0 |

**Table 6:** The latency, and performance on nuScenes-Occupancy under camera-only and multi-modal settings.

| Models | Latency(ms) | IoU | mIoU |
|:---|:---:|:---:|:---:|
| C-Baseline [57] | 172.4 | 19.3 | 10.3 |
| C-CONet [57] | 285.7 | 20.1 | 12.8 |
| C-OccGen(step1) | 294.1 | 23.0 | 14.2 |
| C-OccGen(step2) | 312.5 | 23.3 | 14.4 |
| Baseline [57] | 243.9 | 29.1 | 15.1 |
| CONet [57] | 344.8 | 29.5 | 20.1 |
| OccGen(step1) | 357.1 | 29.3 | 21.7 |
| OccGen(step2) | 400.0 | 29.7 | 21.8 |

of the models. The results are shown in Tab. 6. Compared with the representative discriminative methods, OccGen achieves better results than state-of-the-art CONet [57] when using only one sampling step, with comparable latency on the camera-only, and multi-modal settings. When adopting two sampling steps, the performance is further boosted to 21.8% and 14.4% on the multi-modal and camera-only benchmarks, at a loss of 20 ∼ 50 ms. These results show that OccGen can progressively refine the output occupancy multiple times with reasonable time cost.

**Uncertainty Awareness.** In addition to the performance gains, the proposed OccGen can naturally provide uncertainty estimates. In the multi-step sampling process, we can simply count the voxels where the predicted result of each step differs from the result of the previous step, thereby obtaining an uncertainty occupancy result. We can see from Fig. 6 that the areas with high uncertainty in the uncertainty map often align with those in the error map, which indicates incorrect prediction regions. In comparison, OccGen offers a straightforward and
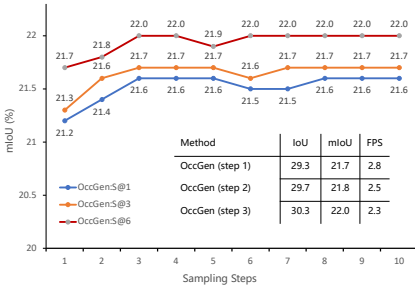
**Fig. 5:** The results of multiple inferences on nuScenese-Occupancy under multi-modal setting.
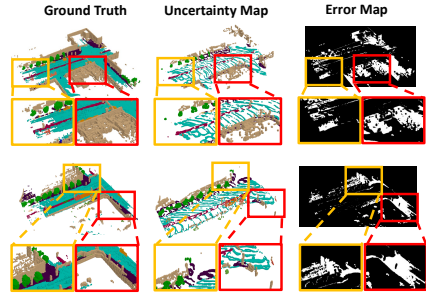


**Fig. 6:** The visualization of uncertainty map and error map on nuScenese-Occupancy under multi-modal setting.

inherently capable approach, whereas previous methods [14, 37] require complicated modeling such as Bayesian networks.

### 4.5 Qualitative Results

In Fig. 4, we visualize the predicted results of 3D semantic occupancy on nuScenes-Occupancy from CONet [57] and our proposed OccGen. Compared with CONet, our method can better understand the scene-level semantic layout and perform local region completion. It is obvious that the regions of "drivable surface" and "sidewalk" predicted by our OccGen have higher continuity and integrity, and can effectively reduce a large number of hole areas compared with CONet. One more interesting observation is that due to the ground truth being initially constructed based on sparse LiDAR data, the shape of voxels in space is not very well-defined, especially in the drivable area. However, both CONet [57] and OccGen yields smoother predictions for these occupancy results.

## 5    Conclusion

In this paper, we propose OccGen, a simple yet powerful generative perception model for 3D semantic occupancy prediction. OccGen adapts a "noise-to-occupancy" generative paradigm, progressively inferring and refining the occupancy map from a random Gaussian distribution. OccGen consists of two main components: a conditional encoder that processes the multi-modal inputs as condition inputs and a progressive refinement decoder that produces fine-grained occupancy in a coarse-to-fine manner. OccGen has achieved state-of-the-art performance for 3D semantic occupancy prediction on nuScenes-Occupancy and SemanticKITTI. In addition, the proposed OccGen has shown desirable properties that discriminative models cannot achieve, such as progressive inference and uncertainty estimates. Currently, the latency of our OccGen is comparable to the previous state-of-the-art methods and has not achieved a significant speed advantage. Next, we will explore a more lightweight generative architecture for 3D semantic occupancy prediction.

# References

1. Amit, T., Nachmani, E., Shaharbany, T., Wolf, L.: Segdiff: Image segmentation with diffusion probabilistic models. arXiv preprint arXiv:2112.00390 (2021) 2, 4
2. Behley, J., Garbade, M., Milioto, A., Quenzel, J., Behnke, S., Stachniss, C., Gall, J.: Semantickitti: A dataset for semantic scene understanding of lidar sequences. In: ICCV. pp. 9297–9307 (2019) 9, 11, 22, 23, 27, 28
3. Berman, M., Triki, A.R., Blaschko, M.B.: The lovász-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In: CVPR. pp. 4413–4421 (2018) 9
4. Caesar, H., Bankiti, V., Lang, A.H., Vora, S., Liong, V.E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuscenes: A multimodal dataset for autonomous driving. In: CVPR. pp. 11621–11631 (2020) 9, 22
5. Cao, A.Q., de Charette, R.: Monoscene: Monocular 3d semantic scene completion. In: CVPR. pp. 3991–4001 (2022) 2, 3, 9, 10, 11, 23
6. Chen, S., Sun, P., Song, Y., Luo, P.: Diffusiondet: Diffusion model for object detection. In: ICCV. pp. 19830–19843 (2023) 2, 4, 9
7. Chen, T., Li, L., Saxena, S., Hinton, G., Fleet, D.J.: A generalist framework for panoptic segmentation of images and videos. arXiv preprint arXiv:2210.06366 (2022) 4
8. Chen, X., Lin, K.Y., Qian, C., Zeng, G., Li, H.: 3d sketch-aware semantic scene completion via semi-supervised structure prior. In: CVPR. pp. 4193–4202 (2020) 2, 3, 10
9. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020) 6
10. Everingham, M., Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. IJCV (2009) 11
11. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: CVPR. pp. 3354–3361 (2012) 22
12. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. NeurIPS **27** (2014) 5, 21
13. Graham, B., Engelcke, M., Van Der Maaten, L.: 3d semantic segmentation with submanifold sparse convolutional networks. In: CVPR. pp. 9224–9232 (2018) 20
14. Harakeh, A., Smart, M., Waslander, S.L.: Bayesod: A bayesian approach for uncertainty estimation in deep object detectors. In: ICRA. pp. 87–93. IEEE (2020) 14
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016) 6, 23
16. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. NeurIPS **33**, 6840–6851 (2020) 2, 4, 5, 9, 21, 22, 26
17. Hu, Y., Yang, J., Chen, L., Li, K., Sima, C., Zhu, X., Chai, S., Du, S., Lin, T., Wang, W., et al.: Planning-oriented autonomous driving. In: CVPR. pp. 17853–17862 (2023) 2
18. Huang, J., Huang, G., Zhu, Z., Du, D.: Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. arXiv preprint arXiv:2112.11790 (2021) 2, 3, 10
19. Huang, Y., Zheng, W., Zhang, Y., Zhou, J., Lu, J.: Tri-perspective view for vision-based 3d semantic occupancy prediction. In: CVPR. pp. 9223–9232 (2023) 2, 4, 10, 11

20. Jang, E., Gu, S., Poole, B.: Categorical reparameterization with gumbel-softmax. arXiv preprint arXiv:1611.01144 (2016) 7, 19
21. Ji, Y., Chen, Z., Xie, E., Hong, L., Liu, X., Liu, Z., Lu, T., Li, Z., Luo, P.: Ddp: Diffusion model for dense visual prediction. In: ICCV. pp. 21741–21752 (2023) 2, 4, 9
22. Jia, X., Gao, Y., Chen, L., Yan, J., Liu, P.L., Li, H.: Driveadapter: Breaking the coupling barrier of perception and planning in end-to-end autonomous driving. In: CVPR. pp. 7953–7963 (2023) 2
23. Jiang, H., Cheng, T., Gao, N., Zhang, H., Liu, W., Wang, X.: Symphonize 3d semantic scene completion with contextual instance queries. CVPR (2024) 11
24. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014) 10
25. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2015) 24
26. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114 (2013) 21
27. Lang, A.H., Vora, S., Caesar, H., Zhou, L., Yang, J., Beijbom, O.: Pointpillars: Fast encoders for object detection from point clouds. In: CVPR. pp. 12697–12705 (2019) 2, 27
28. Li, J., Han, K., Wang, P., Liu, Y., Yuan, X.: Anisotropic convolutional networks for 3d semantic scene completion. In: CVPR. pp. 3351–3359 (2020) 2, 3, 10, 11
29. Li, Y., Yu, Z., Choy, C., Xiao, C., Alvarez, J.M., Fidler, S., Feng, C., Anandkumar, A.: Voxformer: Sparse voxel transformer for camera-based 3d semantic scene completion. In: CVPR. pp. 9087–9098 (2023) 2, 4, 9, 10, 11, 23, 24, 25
30. Li, Y., Ge, Z., Yu, G., Yang, J., Wang, Z., Shi, Y., Sun, J., Li, Z.: Bevdepth: Acquisition of reliable depth for multi-view 3d object detection. arXiv preprint arXiv:2206.10092 (2022) 2, 3, 7, 9, 19, 26, 27
31. Li, Z., Wang, W., Li, H., Xie, E., Sima, C., Lu, T., Qiao, Y., Dai, J.: Bevformer: Learning bird's-eye-view representation from multi-camera images via spatiotemporal transformers. In: ECCV. pp. 1–18. Springer (2022) 2, 3
32. Li, Z., Yu, Z., Austin, D., Fang, M., Lan, S., Kautz, J., Alvarez, J.M.: Fb-occ: 3d occupancy prediction based on forward-backward view transformation. arXiv preprint arXiv:2307.01492 (2023) 24, 25
33. Liang, T., Xie, H., Yu, K., Xia, Z., Lin, Z., Wang, Y., Tang, T., Wang, B., Tang, Z.: Bevfusion: A simple and robust lidar-camera fusion framework. In: NeurIPS (2022) 2, 3, 7, 19
34. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: CVPR. pp. 2117–2125 (2017) 6, 23
35. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: CVPR. pp. 10012–10022 (2021) 6
36. Liu, Z., Tang, H., Amini, A., Yang, X., Mao, H., Rus, D., Han, S.: Bevfusion: Multitask multi-sensor fusion with unified bird's-eye view representation. In: ICRA (2023) 2, 3, 7, 19
37. Loquercio, A., Segu, M., Scaramuzza, D.: A general framework for uncertainty estimation in deep learning. IEEE Robotics and Automation Letters 5(2), 3153–3160 (2020) 14
38. Lu, H., Tang, J., Xu, X., Cao, X., Zhang, Y., Wang, G., Du, D., Chen, H., Chen, Y.: Scaling multi-camera 3d object detection through weak-to-strong eliciting. arXiv preprint arXiv:2404.06700 (2024) 2

39. Lu, H., Zhang, Y., Lian, Q., Du, D., Chen, Y.: Towards generalizable multi-camera 3d object detection via perspective debiasing. arXiv preprint arXiv:2310.11346 (2023) 2

40. Nichol, A.Q., Dhariwal, P.: Improved denoising diffusion probabilistic models. In: ICML. pp. 8162–8171. PMLR (2021) 9, 26

41. Philion, J., Fidler, S.: Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In: ECCV. pp. 194–210. Springer (2020) 7, 19

42. Roldao, L., de Charette, R., Verroust-Blondet, A.: Lmscnet: Lightweight multiscale 3d semantic completion. In: 3DV (2020) 2, 3, 10, 11

43. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR. pp. 10684–10695 (2022) 4

44. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: MICCAI. pp. 234–241. Springer (2015) 4

45. Saharia, C., Chan, W., Chang, H., Lee, C., Ho, J., Salimans, T., Fleet, D., Norouzi, M.: Palette: Image-to-image diffusion models. In: SIGGRAPH. pp. 1–10 (2022) 4

46. Saxena, S., Kar, A., Norouzi, M., Fleet, D.J.: Monocular depth estimation using diffusion models. arXiv preprint arXiv:2302.14816 (2023) 2

47. Saxena, S., Kar, A., Norouzi, M., Fleet, D.J.: Monocular depth estimation using diffusion models. arXiv preprint arXiv:2302.14816 (2023) 4

48. Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., Ganguli, S.: Deep unsupervised learning using nonequilibrium thermodynamics. In: ICML. pp. 2256–2265. PMLR (2015) 4

49. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 (2020) 2, 4, 5, 9, 21, 22, 26

50. Song, S., Yu, F., Zeng, A., Chang, A.X., Savva, M., Funkhouser, T.: Semantic scene completion from a single depth image. In: CVPR. pp. 1746–1754 (2017) 2, 3, 22

51. Tang, P., Wang, Z., Wang, G., Zheng, J., Ren, X., Feng, B., Ma, C.: Sparseocc: Rethinking sparse latent representation for vision-based semantic occupancy prediction. arXiv preprint arXiv:2404.09502 (2024) 2

52. Tian, X., Jiang, T., Yun, L., Wang, Y., Wang, Y., Zhao, H.: Occ3d: A large-scale 3d occupancy prediction benchmark for autonomous driving. arXiv preprint arXiv:2304.14365 (2023) 2, 4, 22, 23, 24

53. Tong, W., Sima, C., Wang, T., Chen, L., Wu, S., Deng, H., Gu, Y., Lu, L., Luo, P., Lin, D., et al.: Scene as occupancy. In: ICCV. pp. 8406–8415 (2023) 2, 4

54. Vaswani, A., Shazeer, N.M., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. NeurIPS (2017) 6

55. Vora, S., Lang, A.H., Helou, B., Beijbom, O.: Pointpainting: Sequential fusion for 3d object detection. In: CVPR. pp. 4604–4612 (2020) 2

56. Wang, C., Ma, C., Zhu, M., Yang, X.: Pointaugmenting: Cross-modal augmentation for 3d object detection. In: CVPR. pp. 11794–11803 (2021) 2

57. Wang, X., Zhu, Z., Xu, W., Zhang, Y., Wei, Y., Chi, X., Ye, Y., Du, D., Lu, J., Wang, X.: Openoccupancy: A large scale benchmark for surrounding semantic occupancy perception. In: ICCV. pp. 17850–17859 (2023) 2, 4, 7, 9, 10, 13, 14, 22, 23, 24, 26, 27, 28, 29

58. Wei, Y., Zhao, L., Zheng, W., Zhu, Z., Zhou, J., Lu, J.: Surroundocc: Multi-camera 3d occupancy prediction for autonomous driving. In: ICCV. pp. 21729–21740 (2023) 2, 4, 9, 24, 25

59. Wolleb, J., Sandkühler, R., Bieder, F., Valmaggia, P., Cattin, P.C.: Diffusion models for implicit image segmentation ensembles. In: MIDL. pp. 1336–1348 (2022) 2, 4
60. Wu, J., Fang, H., Zhang, Y., Yang, Y., Xu, Y.: Medsegdiff: Medical image segmentation with diffusion probabilistic model. arXiv preprint arXiv:2211.00611 (2022) 2, 4
61. Yan, X., Gao, J., Li, J., Zhang, R., Li, Z., Huang, R., Cui, S.: Sparse single sweep lidar point cloud segmentation via learning contextual shape priors from scene completion. In: AAAI. vol. 35, pp. 3101–3109 (2021) 2, 3, 10, 11
62. Yan, Y., Mao, Y., Li, B.: SECOND: Sparsely embedded convolutional detection. Sensors (2018) 6, 20, 27
63. Yin, T., Zhou, X., Krahenbuhl, P.: Center-based 3d object detection and tracking. In: CVPR. pp. 11784–11793 (2021) 27
64. Zhang, Y., Zhou, Z., David, P., Yue, X., Xi, Z., Gong, B., Foroosh, H.: Polarnet: An improved grid representation for online lidar point clouds semantic segmentation. In: CVPR. pp. 9601–9610 (2020) 2
65. Zhang, Y., Zhu, Z., Du, D.: Occformer: Dual-path transformer for vision-based 3d semantic occupancy prediction. In: ICCV. pp. 9433–9443 (2023) 2, 4, 9, 10, 11, 23, 24, 25, 27, 28
66. Zhou, Y., Tuzel, O.: Voxelnet: End-to-end learning for point cloud based 3d object detection. In: CVPR. pp. 4490–4499 (2018) 2, 6, 23
67. Zhu, X., Zhou, H., Wang, T., Hong, F., Ma, Y., Li, W., Li, H., Lin, D.: Cylindrical and asymmetrical 3d convolution networks for lidar segmentation. In: CVPR. pp. 9939–9948 (2021) 2
68. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable detr: Deformable transformers for end-to-end object detection. In: ICLR (2021), https://openreview.net/forum?id=gZ9hCDWe6ke 7

# Supplementary Materials for
# OccGen: Generative Multi-modal 3D Occupancy
# Prediction for Autonomous Driving

In the supplementary material, we first present the methodology details of our proposed OcccGen, including hard 2D-to-3D view transformation, geometry mask, discriminative vs. generative modeling, and DDPM vs. DDIM. Then, we provide the details of datasets and implementation. Furthermore, we present additional experimental results to demonstrate the effectiveness of OccGen. Finally, we discuss the broader impact statement and limitations.

## 6   Methodology Details

### 6.1   Hard 2D-to-3D View Transformation

The previous LSS-based methods [33, 36, 41] associate a set of discrete depths for every pixel, covering the full range of potential depth values. These methods typically choose the softmax operation, which is a smooth approximation of $argmax$, which normalizes the vector, enabling gradient computation while the values can also represent probabilities. Intuitively, this soft approach to depth prediction allows the network to learn depth information that is more conductive to feature optimization, rather than more precise depth information. As per the findings in [30], this soft depth prediction fails to obtain precise depth information, consequently leading to the presence of ambiguous grids in the camera voxel. In light of this, we propose a hard 2D-to-3D view transformation method that utilizes hard Gumbel-Softmax [20] to obtain a deterministic discrete output vector. The key point lies in the fact that the $argmax$ operation for obtaining a one-hot vector is non-differentiable, making it impossible to calculate gradients, and consequently, network updates cannot be performed. Hard Gumbel-Softmax is a deterministic version of Gumbel-Softmax, where instead of sampling from the Gumbel-Softmax distribution. The formula for hard Gumbel-Softmax can be expressed as follows:

$$\mathbf{hard\_gumbel\_softmax}(z) = \mathbf{one\_hot}(\mathbf{argmax}((z + g)/\tau)) \qquad (10)$$

where $z$ represents the logits. $g$ is sampled from the Gumbel distribution, typically calculated as $-log(-log(u))$ with u being a uniform random variable. $\tau$ is the temperature parameter controlling the softness of the distribution. When $\tau$ approaches zero, the hard Gumbel-Softmax approaches the one-hot encoding. The introduced hard 2D-to-3D view transformation enables gradient backpropagation during training and yields a definitive assignment of discrete depth during inference.

### 6.2   Geometry Mask

The camera voxel features $F_c$ obtained through the hard 2D-to-3D view transformation module contain some voxels with misleading information, leading to a

blurred feature distribution. This issue arises because all points along a camera ray in 3D space are projected to the same location on the 2D image plane, resulting in some voxels sharing the same image features. This inaccurate 3D spatial structure hinders subsequent feature fusion and adversely impacts the final detection results. To overcome this challenge, we propose a straightforward method to fully exploit the geometry-aware correspondence between camera and LiDAR modalities. We generate a geometry mask by leveraging the LiDAR voxel features and then applying it to the camera voxel features. This process effectively bridges the gap between the image voxel representation and the true spatial distribution, resulting in improved feature representation and better alignment with the real-world scene.

Due to the limitations of LiDAR sensors, the raw point cloud data only covers a portion of the real scene, resulting in incomplete object shapes. When the point cloud is regularized into voxel grids, only a small portion of these voxels contain non-empty information. As a result, the initial sparse voxels fail to adequately represent the 3D geometry-aware correspondence between LiDAR and camera features. Commonly used 3D sparse convolutions [13, 62] include two types, regular and submanifold sparse convolutions. The submanifold sparse convolution performs convolution operations only on the sparse features located at the center of its receptive field. Consequently, the output features of submanifold convolution maintain consistency with the positions of the input sparse features. On the other hand, regular sparse convolution performs convolution operations as long as there are features within its receptive field, resulting in the propagation of sparse features to their surrounding positions. This leads to a substantial increase in the density of the output sparse features. Through performing several 3D sparse convolutional operations, non-empty voxels are significantly increased, effectively covering more 3D space. Compared with the original LiDAR features, most of the 3D voxels are completed, and the generated 3D geometry-aware constraints can better reflect the real scene distribution.

Specifically, we can generate the geometry mask as follows,

$$L_{mask} = f_{\text{dense}}(f_{\text{reg}}(f_{\text{sp}}(F_{\text{p}}))) \tag{11}$$

where $f_{sp}$ denotes a 3D sparse convolution block that mixes both regular and submanifold convolution, $f_{\text{reg}}$ denotes regular 3D sparse convolution, and $f_{\text{dense}}$ indicates that sparse camera voxel features are densified by padding zeros in empty positions. However, there is no guarantee that all foreground objects can be adequately represented even after the expansion of sparse features. Directly utilizing the aforementioned constraints on the camera voxel grids can obtain numerous all-zero features, potentially losing meaningful features. Therefore, we employ a softmax operation along the height dimension to maintain the density of camera voxel features. Subsequently, we apply the 3D constraint to the image voxel feature, which can be represented as,

$$F_c = \textbf{Softmax}(L_{mask}) \cdot F_c. \tag{12}$$

This weight assignment reduces the influence of misleading or ambiguous features during the transformation, effectively guaranteeing the robustness of the spatial information.

### 6.3   Discriminative vs. Generative Modeling

**Discriminative Modeling.** Discriminative methods for 3D occupancy semantic prediction aim to predict the occupancy and semantic labels of voxels in a 3D scene. These methods typically focus on learning the conditional distribution $P(Y|X)$, where $Y$ represents the occupancy and semantic labels of voxels, and $X$ represents the observed 3D scene data (e.g., point clouds or multi-view images). However, only learning these mapping between inputs and outputs may result in a limited understanding of the overall scene context. This can lead to incomplete or inaccurate scene completions, especially in complex scenes with intricate spatial relationships between objects. From the perspective of uncertainty, discriminative methods often do not explicitly model uncertainty in predictions, which can be crucial for 3D occupancy prediction where the inputs may be noisy or incomplete. This can lead to overconfident predictions in uncertain regions of the scene. Furthermore, these methods may struggle to incorporate prior knowledge or constraints into the learning process, which can be important for ensuring the coherence and consistency of the completed scene.

**Generative Modeling.** Generative modeling for 3D occupancy semantic prediction aim to generate complete 3D scenes from partial or incomplete input data. These methods often leverage generative models, such as GAN [12], VAE [26] and diffusion model [16, 49], to predict the semantic occupancy. Inspired by the great success of diffusion models, we adopt the diffusion model as our pipeline. In contrast to the previous discriminative works, which typically optimize the posterior probability $P(Y|X)$, our method establishes the joint probability $P(X,Y) = P(Y|X) * P(X)$. This optimization process can be denoted as

$$\operatorname*{argmax}_{f_\theta} p(X,Y) = \operatorname*{argmax}_{f_\theta}\{\underbrace{\log p(Y|X)}_{\text{data term}} + \underbrace{\log p(X)}_{\text{prior term}}\}. \tag{13}$$

The data term can be formally defined as the process of learning the mapping between multi-modal inputs and occupancy map, while the prior term can be denoted as a denoising diffusion process for the multi-modal inputs. This holistic procedure can be seamlessly integrated into the framework of the conditional diffusion model. Generative models inherently capture uncertainty in the predictions. This can be useful in scenarios where the inputs are noisy or ambiguous, as the model can provide a measure of confidence in its predictions. In addition, generative models can generate more coherent and contextually relevant completions thanks to the prior knowledge about the scene incorporated in the training process.

## 6.4   DDPM and DDIM

A significant limitation of denoising diffusion probabilistic models (DDPM) [16] is their requirement for numerous iterations to generate high-quality samples. This is due to the generative process, which transforms noise into data, approximating the reverse of the forward diffusion process. The forward diffusion process may involve thousands of steps, necessitating iteration over all these steps to produce a single sample. This process is much slower compared to GANs, which only require one pass through a network. To address this efficiency gap between DDPM and GAN, denoising diffusion implicit model (DDIM) [49]. DDIM enables significantly faster sampling without compromising on the training objective, making generative models using this architecture competitive with GAN of the same model size and sample quality. This is achieved by estimating the cumulative effect of multiple Markov chain steps and incorporating them simultaneously. Since each Markov jump is modeled as a Gaussian distribution, they approximate the combined effect of multiple jumps by using a higher-variance Gaussian distribution with the same mean. It is worth noting that the sum of two Gaussians remains Gaussian. In this paper, we utilize DDPM and DDIM to corrupt the ground truth occupancy and progressively refine the noise map to obtain the final results. The results on DDIM and DDPM are listed in the experimental part.

## 7   Dataset and Implementation

### 7.1   Details of Dataset.

We provide results on nuScenes-Occupancy [57], Occ3D-nuScenes [52] and SemanticKITTI [2]. nuScenes-Occupancy and Occ3D-nuScenes are extended from the large-scale nuScenes [4] dataset with dense semantic occupancy annotation. SemanticKITTI [2] provides dense semantic annotations for each LiDAR sweep from the KITTI Odometry Benchmark [11]. We will introduce nuScenes-Occupancy [57], Occ3D-nuScnes [52] and SemanticKITTI [2] in sequence.

**nuScenes-Occupancy.** In the pursuit of establishing a large-scale surrounding occupancy perception dataset, wang [57] et al. introduced the nuScenes-Occupancy based on nuScenes [4]. Notably, the nuScenes-Occupancy dataset boasts approximately 40 times more annotated scenes and about 5 times more annotated frames compared to the work presented in [50]. To efficiently achieve this extensive annotation and densification of occupancy labels, they introduced an Augmenting And Purifying (AAP) pipeline. The pipeline initiates annotation through the superimposition of multi-frame LiDAR points. Acknowledging the sparsity inherent in the initial annotation attributed to occlusion or limitations in LiDAR channels, they employed a strategy to augment it with pseudo-occupancy labels. These pseudo labels are constructed using a pre-trained baseline. To further enhance the quality of the annotations by reducing noise and artifacts, human efforts are enlisted in the purification process.

**Occ3D-nuScenes.** Occ3D-nuScenes [52] is a comprehensive autonomous driving dataset comprising 700 training scenes and 150 validation scenes. Each frame in this dataset features a 32-beam LiDAR point cloud and six RGB images captured by six cameras positioned at different angles around the LiDAR. These frames are densely annotated with voxel-wise semantic occupancy labels. The dataset's occupancy scope spans from $-40m$ to $40m$ along the $X$ and $Y$ axes, and from $-1m$ to $5.4m$ along the $Z$ axis in the ego coordinate system. The voxel size for the occupancy labels is $0.4m \times 0.4m \times 0.4m$. Semantic labels in the dataset encompass 17 categories, which include 16 known object classes and an additional "empty" class.

**SemanticKITTI.** The SemanticKITTI dataset [2] is focused on semantic scene understanding using LiDAR points and front cameras. OccGen is evaluated for semantic scene completion using the monocular left camera as input, following the approach of MonoScene [5] and OccFormer [65]. In this evaluation, the ground truth semantic occupancy is represented as $256 \times 256 \times 32$ voxel grids. Each voxel is $0.2m \times 0.2m \times 0.2m$ in size and is annotated with one of 21 semantic classes (19 semantics, 1 free, 1 unknown). Similar to previous work [5, 29, 65], the dataset's 22 sequences are split into 10/1/11 for training/validation/testing.

## 7.2 Implementation Details

In the camera stream, we adopt the ResNet-50 [15] model as our image backbone and employ the FPN [34] for multi-scale camera feature fusion, generating the image feature maps of size $6 \times 56 \times 100$, with 512 channels. Then, we utilize the proposed hard 2D-to-3D image view transformation to generate the camera voxel feature of size $128 \times 128 \times 10$, with 80 channels. In the LiDAR stream, the point cloud is constrained within the range of $[-51.2m, 51.2m]$ for $X$ and $Y$ axis, and $[-5m, 3m]$ for the $Z$ axis. Voxelization is performed with a voxel size of $(0.1m, 0.1m, 0.1m)$. We utilize VoxelNet [66] as the backbone and employ FPN-3D [34] to produce the LiDAR voxel features of size $128 \times 128 \times 10$, with 80 channels. In order to fully exploit the implicit geometry-aware cues between camera and LiDAR modalities, we utilize the structure knowledge in LiDAR modality to guide the camera modality to learn geometry mask $80 \times 128 \times 128 \times 10$, which can improve the generalization capability of the fused voxel features significantly. Subsequently, we follow [57] and utilize ResNet3D and FPN-3D to generate multi-scale voxel features as condition input for the progressive refinement module. The progressive refinement consists of six refinement layers with 3D deformable attention. The refinement layer takes as input the random noise map or the predicted noise map from the last step, the current sampling step, and the multi-scale fusion features. We downsample the random noise map three times to obtain smaller multi-scale noise maps to avoid the high-resolution 3D Gaussian noise map. Then, we reshape these downsampled multi-scale noise maps to obtain initial queries. After the learning process of several refinement layers, we upsample and project the downsampled voxels to the size of the original 3D noise map and obtain the refined voxel features. Finally, we obtain the

3D semantic occupancy by feeding the refined voxel features to the occupancy head [57] for full-scale evaluation.

---

**Algorithm 1** Training algorithm

---

**Input:** Multi-modal inputs: $\{X_p, X_c\}$; GT occupancy: $Y$;
**Output:** Training loss
 1: Extract multi-modal features $F_p$ and $F_c$. $F_p, F_c \leftarrow$ **Extractor**$(X_p, X_c)$
 2: Aggregate the camera features with a geometry mask. $F_c \leftarrow$ **Aggregate**$(F_c, M_p)$
 3: Obtain the multi-modal fusion features $F_m$. $F_m \leftarrow$ **Fuser**$(F_p, F_c)$;
 4: Encoding the ground truth occupancy. $Y_0 \leftarrow$ **Encoding**$(Y)$
 5: Construct noise signal and choose step index. $t \leftarrow$ **Randint**$(0, T)$, $\epsilon \leftarrow$ **Randn**(mean=0, std=1)
 6: Signal scaling. $Y_0 \leftarrow$ **Norm**$(Y_0)$
 7: Corrupt the occupancy input. $Y_t \leftarrow$ **Schedule**$(t) \times Y_0 + (1 -$ **Schedule**$(t)) \times \epsilon$
 8: Obtain the downsampled multi-scale noise maps. $Y_t^i \leftarrow$ **Downsample**$(Y_t)$
 9: Obtain the refined noise map. $Y_t \leftarrow$ **Refine**$(F_m, Y_{t+1}, t)$
10: Predict the occupancy results. $\hat{Y}_t \leftarrow$ **Voxel2Occ**$(Y_t)$
11: Calculate the training loss $\mathcal{L}_{\text{total}}$ (Eq. 9).

---

For the input, we follow the setting in [57] to take the image size as $900 \times 1600$, and utilize 10 sweeps to densify the LiDAR point cloud. During training, we adopt similar data augmentation strategies in [57] for both the image and LiDAR data. In our experiments, we utilize the AdamW [25] optimizer with a weight decay of 0.01 and an initial learning rate of $2e^{-4}$. We also utilize the cosine learning rate scheduler with linear warming up in the first 500 iterations. During training, we first construct the diffusion process from ground truth to noisy occupancy and then train the model to reverse this process. Algorithm 1 provides the pseudo-code of OccGen training procedure. The inference procedure of OccGen is a denoising sampling process from noise to 3D semantic occupancy. Starting from 3D voxel grids sampled in Gaussian distribution, the OccGen progressively refines its predictions, as shown in Algorithm 2. All models are trained and inferenced with a batch size of 8 on 8 V100 GPUs.

## 8    Additional Experiments

**Results on Occ3D-nuScenes.** We also compare our OccGen with the state-of-the-art vision-based 3D occupancy prediction methods [29, 32, 58, 65] on Occ3D-nuScenes [52]. For a fair comparison, we removed the LiDAR stream and fusion module from the conditional encoder and followed the same backbone and image size of FB-Occ [32]. As shown in Tab. 7, we can see that OccGen achieves the highest mIoU compared with all existing SOTA methods, demonstrating the effectiveness of OccGen for semantic scene completion.

**Additional Results on nuScenes-Occupancy.** We report the results of IoU and mIoU to represent the accuracy of different methods, and Params and FPS

**Algorithm 2** Inference algorithm

**Input:** Multi-modal inputs: $\{X_p, X_c\}$; Generative steps: $T$;
**Output:** Prediction occupancy $\hat{Y}$
1: Extract multi-modal features. $F_p, F_c \leftarrow$ **Extractor**$(X_p, X_c)$
2: Aggregate the camera features. $F_c \leftarrow$ **Aggregate**$(F_c, M_p)$
3: Obtain the multi-modal fusion features. $F_m \leftarrow$ **Fuser**$(F_p, F_c)$;
4: Initialize the noise map. $Y_T \leftarrow$ **Randn**(mean=0, std=1)
5: **for** $i = 1, 2, ..., T$ **do**
6:     **if** $i > 1$ **then**
7:         Update the current 3D noise map. $Y_i \leftarrow Y_{i-1}$
8:     **else**
9:         Obtain the refined noise map. $Y_i \leftarrow$ **Refine**$(F_m, Y_i, t)$
10:         Predict the occupancy results. $\hat{Y}_i \leftarrow$ **Voxel2Occ**$(Y_i)$
11:         Obtain the noise map for the next evaluation. $Y_i \leftarrow$ **DDIM**$(Y_i, i)$
12:     **end if**
13: **end for**

**Table 7:** Semantic occupancy prediction results on Occ3D-nuScenes validation set.

| Method | Backbone | Image Size | mIoU |
|---|---|---|---|
| OccFormer [65] | ResNet-50 | $900 \times 1600$ | 36.5 |
| SurroundOcc [58] | InternImage-B | $900 \times 1600$ | 40.7 |
| VoxFormer [29] | ResNet-101 | $900 \times 1600$ | 40.7 |
| FB-Occ [32] | ResNet-50 | $900 \times 1600$ | 41.8 |
| Ours | ResNet-50 | $900 \times 1600$ | **42.6** |

to represent the efficiency of the models. The results are shown in Tab. 8. Compared with the representative discriminative methods, OccGen achieves better results when using only one sampling step, with fewer parameters and comparable FPS on the camera-only, LiDAR-only, or multi-modal methods. When adopting three sampling steps, the performance is further boosted to 22.0%, 16.8%, and 14.5% on the multi-modal, camera-only, and LiDAR-only benchmarks, at a loss of $0.3 \sim 0.5$ FPS. These results show that OccGen can progressively refine the output occupancy multiple times with reasonable time cost. We note that OccGen consistently delivers the best IoU results across almost all categories in the third step, which indicates that our method can better complete the scenes due to our coarse-to-fine generation property. We also observe that camera-only methods are more time-consuming compared to LiDAR-only methods due to the 2D-to-3D view transformation. This indicates that a more efficient LSS method is urgent.

**Scaling factor.** The performance of different scaling factors is shown in Tab. 9a. As can be seen, we found the lower scaling factor 0.001 has achieved a bit lower performance than 0.01. A larger scaling factor means more noise is added to the estimate, typically resulting in faster convergence but potentially reducing the quality of sampling. Conversely, a smaller scaling factor reduces the noise level but may require more iteration steps to converge, thereby increasing sampling time.

**Table 8:** Performance on nuScenes-Occupancy (validation set). We report the geometric metric IoU, semantic metric mIoU, IoU, parameters and FPS for each semantic class. The $C, L, M$ denotes *camera, LiDAR* and *multi-modal*. The best results are in boldface (Best camera-only, LiDAR-only, and multi-modal results are marked <span style="color:red">red</span>, <span style="color:blue">blue</span>, and **black**, respectively.

| Method | Input | IoU | mIoU | Params | FPS | barrier | bicycle | bus | car | const. veh. | motorcycle | pedestrian | traffic cone | trailer | truck | drive. suf. | other flat | sidewalk | terrain | manmade | vegetation |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| C-Baseline [57] | C | 19.3 | 10.3 | 93M | 5.8 | 9.9 | 6.8 | 11.2 | 11.5 | 6.3 | 8.4 | 8.6 | 4.3 | 4.2 | 9.9 | 22.0 | 15.8 | 14.1 | 13.5 | 7.3 | 10.2 |
| L-Baseline [57] | L | 30.8 | 11.7 | 63M | 6.9 | 12.2 | 4.2 | 11.0 | 12.2 | 8.3 | 4.4 | 8.7 | 4.0 | 8.4 | 10.3 | 23.5 | 16.0 | 14.9 | 15.7 | 15.0 | 17.9 |
| M-Baseline [57] | M | 29.1 | 15.1 | 117M | 4.1 | 14.3 | 12.0 | 15.2 | 14.9 | 13.7 | 15.0 | 13.1 | 9.0 | 10.0 | 14.5 | 23.2 | 17.5 | 16.1 | 17.2 | 15.3 | 19.5 |
| C-CONet [57] | C | 20.1 | 12.8 | 111M | 3.5 | 13.2 | 8.1 | 15.4 | 17.2 | 6.3 | 11.2 | 10.0 | 8.3 | 4.7 | 12.1 | 31.4 | 18.8 | 18.7 | 16.3 | 4.8 | 8.2 |
| L-CONet [57] | L | 30.9 | 15.8 | 63M | 4.0 | 17.5 | 5.2 | 13.3 | 18.1 | 7.8 | 5.4 | 9.6 | 5.6 | 13.2 | 13.6 | 34.9 | 21.5 | 22.4 | 21.7 | 19.2 | 23.5 |
| M-CONet [57] | M | 29.5 | 20.1 | 137M | 2.9 | 23.3 | 13.3 | 21.2 | 24.3 | 15.3 | 15.9 | 18.0 | 13.3 | 15.3 | 20.7 | 33.2 | 21.0 | 22.5 | 21.5 | 19.6 | 23.2 |
| C-OccGen (step1) | C | 23.0 | 14.2 | 110M | 3.4 | 15.5 | 9.1 | 15.0 | 18.9 | 6.6 | 11.6 | 11.4 | 8.8 | 5.4 | 13.1 | 34.4 | 21.4 | 21.6 | 18.8 | 5.6 | 9.6 |
| C-OccGen (step2) | C | 23.3 | 14.4 | 110M | 3.2 | 14.8 | 8.5 | 15.2 | 19.0 | 7.3 | 11.4 | 11.9 | 8.3 | 6.0 | 13.9 | 34.6 | 22.0 | 21.6 | 19.5 | 5.7 | 9.8 |
| C-OccGen (step3) | C | 23.4 | 14.5 | 110M | 3.0 | 15.5 | 9.1 | 15.3 | 19.2 | 7.3 | 11.3 | 11.8 | 8.9 | 5.9 | 13.7 | 34.8 | 22.0 | 21.8 | 19.5 | 6.0 | 9.9 |
| L-OccGen (step1) | L | 31.1 | 16.1 | 62M | 4.0 | 17.6 | 4.1 | 14.3 | 19.1 | 6.6 | 7.1 | 11.0 | 6.2 | 13.2 | 14.3 | 35.8 | 21.3 | 22.2 | 20.9 | 20.1 | 24.2 |
| L-OccGen (step2) | L | 31.4 | 16.6 | 62M | 3.9 | 18.7 | 5.1 | 15.0 | 19.3 | 7.3 | 7.8 | 11.2 | 6.3 | 13.7 | 14.3 | 36.3 | 21.9 | 22.7 | 21.9 | 20.2 | 24.1 |
| L-OccGen (step3) | L | 31.6 | 16.8 | 62M | 3.7 | 18.8 | 5.1 | 14.8 | 19.6 | 7.0 | 7.7 | 11.5 | 6.7 | 13.9 | 14.6 | 36.4 | 22.1 | 22.8 | 22.3 | 20.6 | 24.5 |
| OccGen (step1) | M | 29.3 | 21.7 | 117M | 2.8 | 25.4 | 16.6 | 22.2 | 26.0 | 13.4 | 19.9 | 21.8 | 14.6 | 17.3 | 22.1 | 35.4 | 24.1 | 24.1 | 22.8 | 19.5 | 22.3 |
| OccGen (step2) | M | 29.7 | 21.8 | 117M | 2.5 | 24.8 | 16.8 | 22.4 | 25.9 | 13.8 | 20.3 | 21.7 | 14.6 | 17.5 | 21.9 | 35.2 | 24.5 | 24.3 | 23.5 | 19.5 | 22.5 |
| OccGen (step3) | M | 30.3 | 22.0 | 137M | 2.3 | 24.9 | 16.4 | 22.5 | 26.1 | 14.0 | 20.1 | 21.6 | 14.6 | 17.4 | 21.9 | 35.8 | 24.5 | 24.7 | 24.0 | 20.5 | 23.5 |

**Table 9:** The diffusion settings of progressive refinement layer on nuScenes-Occupancy. We report Iou and mIoU. Default settings are marked in ▢ gray .

**(a) Scaling factor.** The best scaling factor is 0.01.

| scale | IoU | mIoU |
|---|---|---|
| 0.001 | 30.0 | 21.8 |
| 0.01 | 30.3 | 22.0 |

**(b) Noise schedule.** Cosine works best.

| type | IoU | mIoU |
|---|---|---|
| cosine | 30.3 | 22.0 |
| linear | 29.9 | 21.4 |

**(c) Sampling strategy.** Using DDIM works best.

| type | IoU | mIoU |
|---|---|---|
| DDIM | 30.3 | 22.0 |
| DDPM | 29.2 | 21.7 |

**Noise schedule.** As shown in Tab. 9b, we compare the effectiveness of the cosine schedule [40] and linear schedule [16] in OccGen for occupancy prediction. We observe that the model using a cosine schedule achieves better performance (22.0% vs. 21.4%). The possible reason is that the cosine schedule allows for a smooth reduction in noise, promoting more stable learning dynamics and the linear schedule may sometimes exhibit a more abrupt transition, and its impact on model convergence and sample quality can differ from the cosine schedule.

**Sampling strategy.** As shown in Tab. 9c, we compare the effectiveness of the DDIM [49] and DDPM [16] sampling strategies in OccGen, and find that the model using DDIM is better than DDPM. DDIM uses a non-Markovian diffusion process to accelerate sampling and DDPM is defined as the reverse of a Markovian diffusion process.

**The effectiveness of hard 2D-to-3D view transformation.** We also conduct experiments to fully exploit the effectiveness of hard 2D-to-3D view transformation under the multi-modal setting. The results are shown in Tab. 10 We observe that "Depth supervision" proposed in BEVDepth [30] can boost the

**Table 10:** Ablations on hard 2D-to-3D view transformation in OccGen under the multi-modal setting. "*Hard LSS*" and "*Depth Supervision*" denote hard 2D-to-3D view transformation and the generated depth ground truth following [30, 57] , respectively.

|     | Hard LSS | Depth Supervision | IoU | mIoU |
|-----|----------|-------------------|------|------|
| (a) | -        | -                 | 25.1 | 19.4 |
| (b) | ✓        | -                 | 28.6 | 20.6 |
| (c) | -        | ✓                 | 29.5 | 20.1 |
| (d) | ✓        | ✓                 | 29.4 | 20.8 |

**Table 11:** Ablation study of backbone selection, input size of different modality, and number of denoising layers. *C,L* denotes camera and LiDAR.

|     | Method    | 2D Backbone | Input Size                | Layers | IoU  | mIoU |
|-----|-----------|-------------|---------------------------|--------|------|------|
|     | C-OccGen  | R-50        | $704 \times 256$          | six    | 21.8 | 13.0 |
| (a) | C-OccGen  | R-50        | $1600 \times 900$         | six    | 23.4 | 14.5 |
|     | C-OccGen  | R-101       | $1600 \times 900$         | six    | 23.3 | 15.0 |
| (b) | L-OccGen  | -           | 1 sweep                   | six    | 30.4 | 15.9 |
|     | L-OccGen  | -           | 10 sweeps                 | six    | 31.6 | 16.2 |
|     | OccGen    | R-50        | $1600 \times 900$ 10 sweeps | one  | 29.4 | 21.6 |
| (c) | OccGen    | R-50        | $1600 \times 900$ 10 sweeps | three | 29.9 | 21.7 |
|     | OccGen    | R-50        | $1600 \times 900$ 10 sweeps | six  | 30.4 | 22.0 |

performance of occupancy prediction significantly. This indicates that the accurately predicted depth can lead to more complete occupancy. We also note that our proposed hard 2D-to-3D view transformation can achieve comparable results without adopting depth supervision, which demonstrates the effectiveness of the hard Gumbel softmax on depth prediction.

**Different Experiment Setting.** In this subsection, we ablate the different experiment settings (*e.g.*, input size, backbone selection, number of refinement layers) in Tab. 11. For the camera-based OccGen, using a larger input size ($1600 \times 900$) relatively improves IoU and mIoU by 7.3% and 11.5%. Besides, replacing ResNet-50 with ResNet-101 can further improve the performance of mIoU. For the LiDAR-based OccGen, it is observed that utilizing multi-sweeps as input (following [27, 62, 63], 10 sweeps are used) perform well the single-sweep counterpart on IoU and mIoU. For the multi-modal OccGen, we observe that the number of refinement layers has a discernible impact on the performance. The performance tends to increase with a greater number of layers.

**More visualization.** We visualize the predicted results of OccFormer [65] and Our OccGen on SemanticKITTI [2] in Fig 7. We can observe that OccGen produces more reasonable results than OccFormer [65]. In Fig. 8, we visualize the predicted results of 3D semantic occupancy on nuScenes-Occupancy from CONet [57] and our proposed OccGen. It is evident that the "drivable surface" and "sidewalk" regions predicted by our OccGen exhibit superior continuity and integrity. This results in a significant reduction in the number of void areas com-

pared to the previous SOTA CONet [57]. In Fig. 9, we visualize the predicted results of 3D occupancy of different sampling steps. We observe that the results of the third step have more complete geometric structure and semantic information compared with the generated results of the first step. In Fig. 10, we note that the uncertainty maps of different steps clearly show that the proposed OccGen can iteratively refine the occupancy in a coarse-to-fine manner.
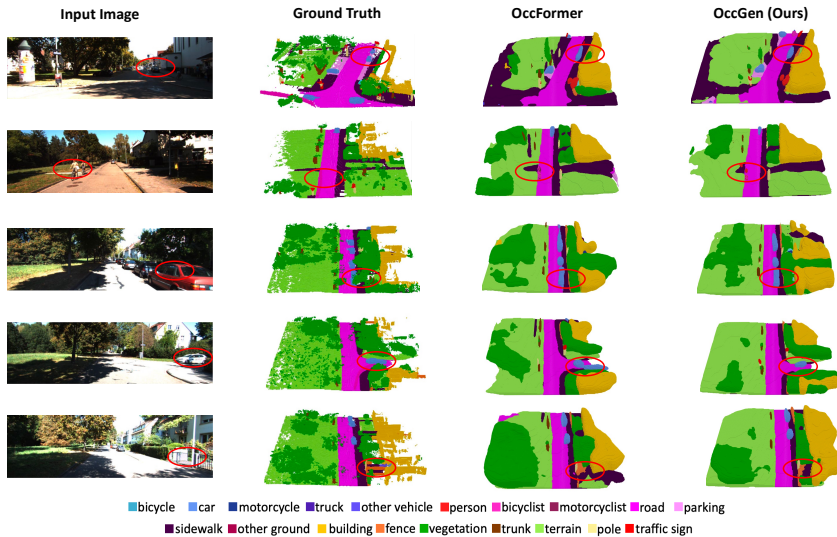


**Fig. 7:** Qualitative results of semantic Scene Completion on SemanticKITTI [2] validation set. The leftmost column shows the input image, the following three columns visualize the results from the ground truth, OccFormer [65], and Our OccGen.

# 9    Broader Impact Statement and Limitations

This paper studies a generative model for 3D occupancy semantic prediction and does not see potential privacy-related issues. Nevertheless, the deployment of a model that is biased toward the training data may introduce significant safety concerns and potential risks when utilized in real-world applications. This research is simple yet effective, which may inspire the community to produce follow-up generative studies for 3D occupancy.
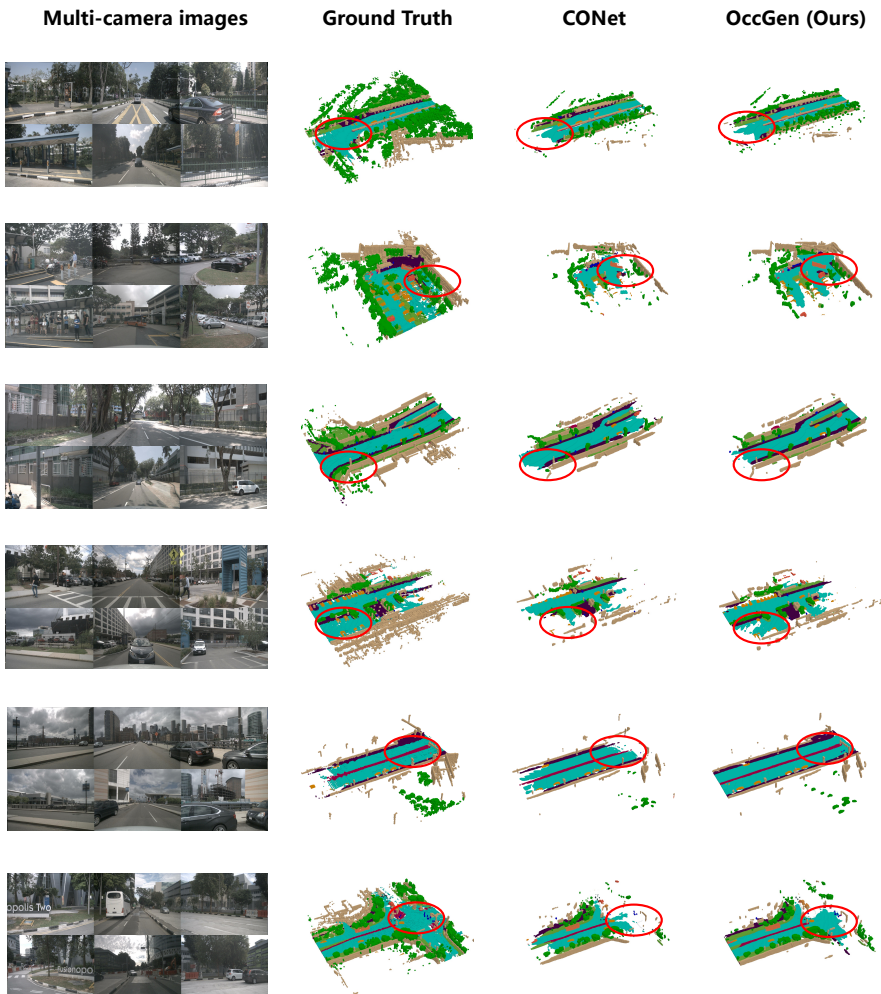
**Fig. 8:** Qualitative results of the 3D semantic occupancy predictions on nuScenes-Occupancy. The leftmost column shows the input surrounding images, the following three columns visualize the 3D semantic occupancy results from the ground truth, CONet [57], and Our OccGen. The regions highlighted by red circles indicate that these areas have obvious differences (better viewed when zoomed in).
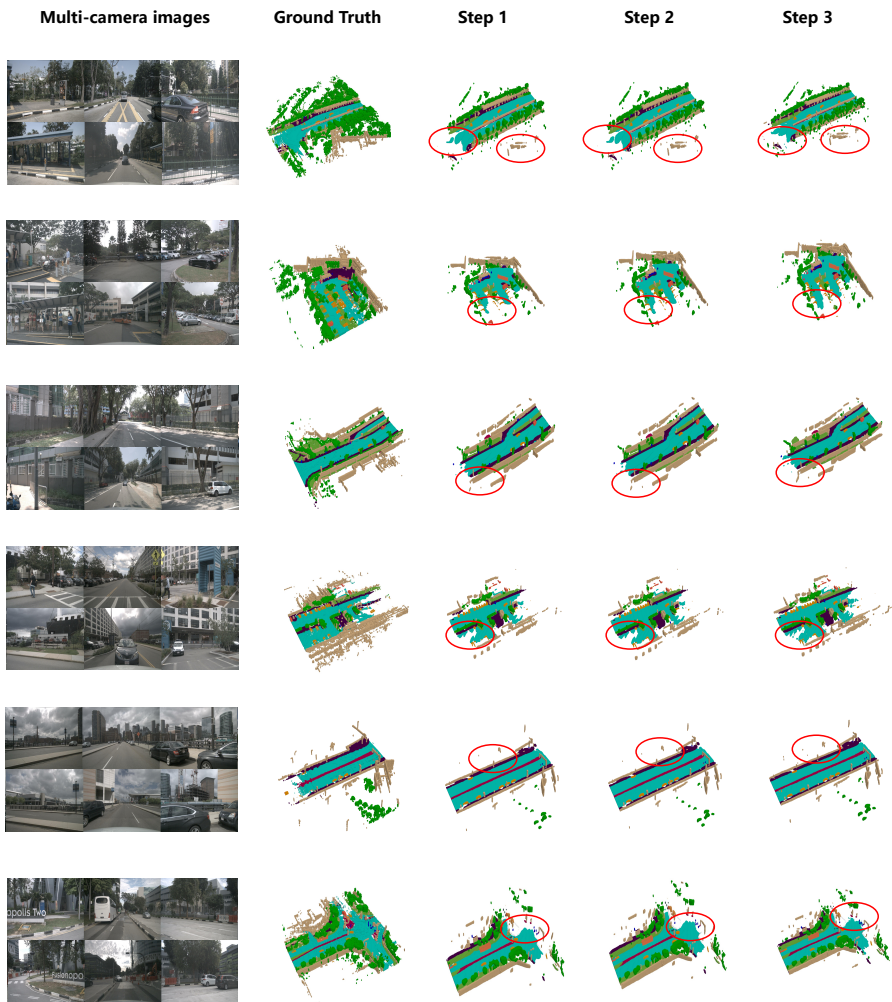
**Fig. 9:** The predicted occupancy results on the different steps of OccGen on nuScenes-Occupancy. The leftmost column shows the input surrounding images, the following four columns visualize the 3D semantic occupancy results from the ground truth, step 1, step 2, and step 3. The regions highlighted by red circles indicate that these areas have obvious differences (better viewed when zoomed in).
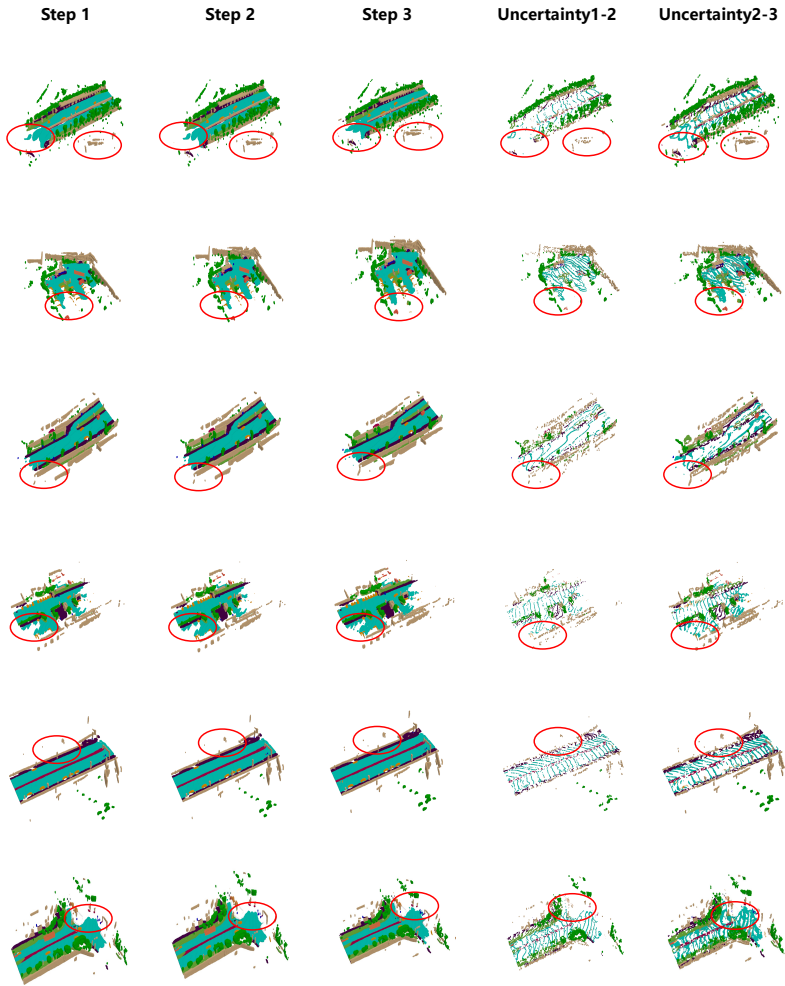
**Fig. 10:** The uncertainty estimates between different steps of OccGen on nuScenes-Occupancy. The left three columns show the predicted 3D semantic occupancy results from step 1, step 2, and step 3. The "Uncertainty 1-2" and "Uncertainty 2-3" represent the high estimated uncertainty voxels from step one to step two and from step two to step three, respectively. The regions highlighted by red circles indicate that these areas have obvious differences (better viewed when zoomed in).