
SyntheOcc: Synthesize Occupancy-Controlled Street View Images through 3D Semantic MPIs

Anonymous Author(s)

Affiliation

Address

email

Abstract

The advancement of autonomous driving is increasingly reliant on high-quality annotated datasets, especially in the task of 3D occupancy prediction, where the occupancy labels require dense 3D annotation with significant human effort. In this paper, we propose **SyntheOcc**, which denotes a diffusion model that Synthesize photorealistic and geometric-controlled images by conditioning Occupancy labels in driving scenarios. This yields an unlimited amount of diverse, annotated, and controllable datasets for applications like training perception models and simulation. SyntheOcc addresses the critical challenge of how to efficiently encode 3D geometric information as conditional input to a 2D diffusion model. Our approach innovatively incorporates 3D semantic multi-plane images (MPIs) to provide comprehensive and spatially aligned 3D scene descriptions for conditioning. As a result, SyntheOcc can generate photorealistic multi-view images and videos that faithfully align with the given geometric labels (semantics in 3D voxel space). Extensive qualitative and quantitative evaluations of SyntheOcc on the nuScenes dataset prove its effectiveness in generating controllable occupancy datasets that serve as an effective data augmentation to perception models.

1 Introduction

With the rapid development of generative models, they have shown realistic image synthesis and diverse controllability. This progress has opened up new avenues for dataset generation in autonomous driving [6, 16, 29, 37]. The task of dataset generation is usually modeled as controllable image generation, where the ground truth (*e.g.* 3D Box) is employed to control the generation of new datasets in downstream tasks (*e.g.* 3D detection). This approach helps to mitigate the data collection and annotation effort as it can generate labeled data for free. However, a novel task of vital importance, occupancy prediction [30, 34], poses new challenges for dataset generation compared with 3D detection. It requires finer and more nuanced geometry controllability, which refers to use the occupancy state and semantics of voxels in the whole 3D space to control the image generation. We argue that solving this problem not only allows us to synthesize occupancy datasets, but also empowers valuable applications such as editing geometry to generate rare data for corner case evaluation, as shown in Fig. 1. In the following, we first illustrate why prior work struggles to achieve the above objective, and then demonstrate how we address these challenges.

In the area of diffusion models, several representative works have displayed high-quality image synthesis; however, they are constrained by limited 3D controllability: they are incapable of editing 3D voxels for precise control. For example, BEVGen [29] generates street view images by conditioning BEV layouts using diffusion models. MagicDrive [6] extend BEVGen and additionally converts the 3D box parameters into text embedding through Fourier mapping that is similar to NeRF [24], and uses cross-attention to learn conditional generation. Although these methods achieve satisfactory results in image generation, their 3D controllability is inherently limited. These approaches are

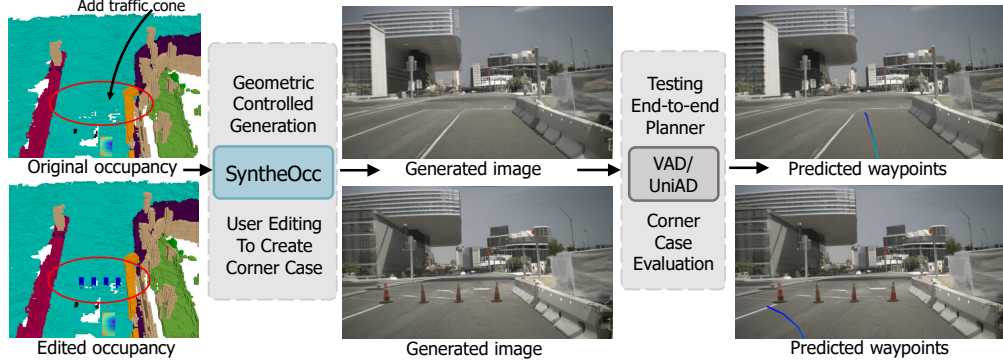


Figure 1: A showcase of application of **SyntheOcc**. We enable geometric-controlled generation that conveys the user editing in 3D voxel space to generate realistic street view images. In this case, we create a rare scene that traffic cones block the way. This advancement facilitates the evaluation of autonomous systems, such as the end-to-end planner VAD [12], in simulated corner case scenes.

restricted to manipulating the scene in types of 3D boxes and BEV layouts, and hardly adapt to finer geometry control such as editing the shape of objects and scenes. Meanwhile, they usually convert conditional input into 1D embedding that aligns with prompt embedding, which is less effective in 3D-aware generation due to lack of spatial alignment with the generated images. This limitation hinders their utility in downstream applications, such as occupancy prediction and editing scene geometry to create long-tailed scenes, where granular volumetric control is paramount in both tasks.

ControlNet [48] and GLIGEN [18] is another type of prominent method in the field of controllable image generation. These approaches exhibit several desirable attributes in terms of controllability. They leverage conditional images such as semantic masks for control, thereby offering a unified framework to manipulate both foreground and background. However, despite its precise spatial control, ControlNet does not align with our specific requirements. Their conditions of pixel-level images differ fundamentally from what we require in 3D contexts. Our experimental results also find that ControlNet struggles to handle overlapping objects with varying depths (see Fig. 7 (a)), as it only utilizes an ambiguous 2D semantic map as conditional input. As a result, it is non-trivial to extend the ControlNet framework and convey their desirable attributes for 3D conditioning.

To address the above challenges, we propose an innovative representation, 3D semantic multi-plane images (MPIs), which contribute to image generation with finer geometric control. In detail, we employ multi-plane images [50] to represent the occupancy, where each plane represents a slice of semantic label at a specific depth. Our 3D semantic MPIs not only preserve accurate and authentic 3D information, but also keep pixel-wise alignment with the generated images. We additionally introduce the MPI encoder to encode features, and the reweighing methods to ease the training with long-tailed cases. As a collection, our framework enables 3D geometry and semantic control for image generation and further facilitates corner case evaluation as depicted in Fig. 1. Finally, experimental results demonstrate that our synthetic data achieve better recognizability, and are effective in improving the perception model on occupancy prediction. In summary, our contributions include:

- We present **SyntheOcc**, an image and video generation framework provides finer and precise 3D geometric control, thereby unlocking a spectrum of applications such as 3D editing, dataset generation, and long-tailed scene generation.
- Incorporating the proposed 3D semantic MPI, MPI encoder, and reweighing strategy, we deliver a substantial advancement in image quality and recognizability over prior works.
- Our extensive experimental results demonstrate that our synthetic data yields an effective data augmentation in the realm of 3D occupancy prediction.

2 Related Work

2.1 3D Occupancy Prediction

The task of 3D occupancy prediction aims to predict the occupancy status of each voxel in 3D space, as well as its semantic label if occupied. Compared with previous perception methods like

3D object detection, occupancy prediction offers a more detailed and nuanced understanding of the environment, as it provides finer geometric details, is capable of handling general, out-of-vocabulary objects, and finally, enriches the planning stack with comprehensive 3D information. Recent methods perform vision-based 3D occupancy prediction [30, 31, 34, 36]. By predicting the geometric and semantic properties of both dynamic and static elements, 3D occupancy prediction offers a more comprehensive understanding of the surrounding environment.

2.2 Diffusion-based Image Generation

Recent advancements in diffusion models (DMs) have achieved remarkable progress in image generation. In particular, Stable Diffusion (SD) [27] employs DMs within the latent space of autoencoders, striking a balance between computational efficiency and high image quality. A noteworthy work is ControlNet [26, 48], which enhances controllability by using image control. We refer readers to recent survey [42] for more details.

2.3 Image Generation in Autonomous Driving

As training neural networks relies heavily on labeled data, numerous studies are delving into dataset generation to boost training. Lift3D [16] designs generative NeRF to synthesize labeled datasets for 3D detection for the first time. Several other works employ BEV layouts to synthesize image data, proving beneficial for perception models. For example, BEVGen [29] conditions BEV layouts to generate multi-view street images, while BEVControl [41] separately generates foregrounds and backgrounds from BEV layouts. MagicDrive [5, 6] generates images with 3D geometry controls by independently encoding objects and maps through a text encoder or map encoder. Compared with MagicDrive, our geometry control is characterized by a more detailed and lossless representation of 3D scenes for control, poses significant challenges than projected layout or box embedding.

Recently, DriveDreamer [33], DrivingDiffusion [17], Drive-WM [35], Panacea [37] and SimGen [51] use a ControlNet framework, which involves projecting primitives like bounding boxes and road maps onto 2D FoV images as a conditioning input. This approach has proven to be effective for geometric control. However, it is limited in that it only achieves alignment at the 2D-pixel level. Consequently, this method falls short in capturing the depth hierarchy and fails to account for the occlusion relationships present in the 3D real world. Besides, adding a depth channel like Panacea [37] may address the limitations of depth order, but it discards the occluded part and only contains partial observation. UrbanGiraffe [44] train a generative NeRF to perform image generation. Another line of research [7, 40] employs a next-frame prediction to achieve a world model [8] integration. WoVoGen [22] creates a future world volume feature using occupancy to guide the generation, but rely on object mask guidance. Recently, some works (Uniscene, Infinicube and Drivingsphere) [14, 23, 39] explore leveraging occupancy as an intermediate representation for generation and downstream application, which directly benefits from our framework.

As described above, most of the prior work is restricted to only modeling a projected primitive of 3D boxes and road maps as conditions. They suffer from ill-posed un-projection ambiguity. In contrast, we model 3D occupancy labels as conditions, as they provide finer geometric details and semantic information. However, designing an input representation of 3D occupancy labels into a 2D diffusion model is challenging. In this paper, we propose a novel representation: 3D semantic Multi-Plane Images (MPIs) as conditional inputs, which not only provide spatial alignment that improves visual consistency, but also encode comprehensive 3D geometric information including occluded parts.

3 Method

Overview The overview of our method is depicted in Fig. 2. Built upon the SD pipeline, we aim to perform geometry-controlled image generation by conditioning on 3D geometry labels with semantics (occupancy labels). One requirement is that the images should faithfully align with the given label. This task is more challenging than conditioned on 3D box due to the sparse and irregular nature of occupancy. We first discuss how to efficiently represent occupancy in Sec. 3.2, followed by our designed MPI encoder to enhance generation quality in Sec. 3.3, and reweighing strategy to handle the long-tailed depth and category in Sec. 3.5.

3.1 Representation of Condition: Local Control Aligns Better than Global Control

One of the key challenges is how to represent our conditional occupancy input. A straightforward method [3, 6] is to convert the 3D occupancy voxel to 1D global embedding that is similar to text

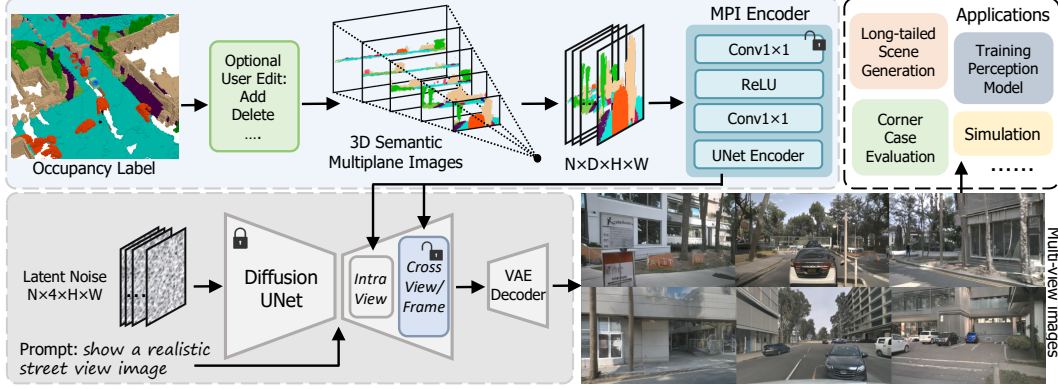


Figure 2: The overall architecture of **SyntheOcc**. We achieve 3D geometric control in image generation by utilizing our proposed 3D semantic multiplane images to encode scene occupancy. In our framework, we can edit the occupied state and semantics of every voxel in 3D space to control the image generation, thereby opening up a wide spectrum of applications as shown in the top right.

embedding, and then use cross-attention to learn controllable generation. However, these global methods can be less effective when dealing with dense or irregular data due to the following reasons: (i) They perform controllable generation through hard encoding the spatial relationship between 1D global embedding and 2D UNet features. (ii) Ignore the underlying geometry alignment between the conditional input and the generated image. In contrast, local methods like ControlNet, directly add spatial features to the UNet features, providing 2D local control with pixel-level spatial alignment. They are better than the global method (see Tab. 1), but suffer from 3D ambiguity (see Fig. 7 (a)). Consequently, this comparison motivates us to seek a more compact and efficient manner to encode and condition our 3D occupancy labels.

3.2 Represent Occupancy as 3D Semantic MPIs

It is non-trivial to design a 3D representation for conditioning. To efficiently store both the semantic and geometric information of the irregular occupancy input, we propose to use multiplane images (MPIs) [50] as representation. An MPI is composed of a series of fronto-parallel RGBA layers within the frustum of the source camera with a specific viewpoint. These planes are arranged at varying depths, from d_{min} to d_{max} , starting from the nearest to the farthest. Each layer of these images contains both an RGB image and an alpha map, which collectively capture the visual and geometric details of the scene at the respective depth. In our work, instead of storing RGB value and alpha map in the original MPI, we store our 3D semantic labels. Each layer of MPI represents the semantic index at the corresponding depth. We display the colored MPI in the top row of Fig. 2 for visual clarity, but we actually use the integer index for learning. We obtain our 3D semantic MPI by:

$$P_l = (u \times d_l, v \times d_l, d_l)^T, \quad (1)$$

$$d_l = d_{min} + (d_{max} - d_{min}) \times l/D, \quad (2)$$

$$\text{MPI}_{n,l} = \text{Interpolate}(\text{Occupancy}, \mathbf{T}_n \cdot \mathbf{K}_n^{-1} \cdot P_l), \quad (3)$$

$$\text{MPI} = \text{Concatenate}(\text{MPI}_{i,j}), \quad (4)$$

$$i \in (0, N), j \in (0, D), \quad (5)$$

where (u, v) is a pixel coordinate in image space, d_l is depth value of the l^{th} layer, n denotes the n^{th} camera view. This equation implies we first back project points P in camera frustum space (u, v, d) to Euclid space (x, y, z) by multiplying inverse intrinsic \mathbf{K}^{-1} . Then we use transformation matrix \mathbf{T} to map points from camera coordinates to occupancy coordinates. We then use the point coordinates to interpolate the nearest semantic index from the dense occupancy voxel to form a slice of MPI. Finally, we concatenate all slices to form $\text{MPI} \in \mathbb{R}^{N \times D \times H \times W}$, where D is the number of layers that is set at 256, N is the number of camera views in the case of batch size = 1.

By representing occupancy as 3D semantic MPI, every pixel in MPI contains geometry and semantic information with implicit depth, seamlessly integrating occluded elements, and ensuring a precise spatial alignment with the generated images.



Figure 3: Visualizations of geometric controlled generation. **Top row:** Fusion of 3D semantic MPI. **Bottom row:** our generation concatenated from neighboring views.

3.3 3D Semantic MPI Encoder

To enable local control with spatially aligned conditions, we develop a simple but effective MPI encoder that aligns the 3D multi-plane feature to the latent space of the diffusion model. The purpose of the MPI encoder is to obtain features from multi-plane images to perform 3D-aware image synthesis. Unlike the original ControlNet which downsampling conditional input through 3×3 convolutions with padding, we design a 1×1 convolutional encoder without downsampling to encode features. In detail, the 3D multiplane features which have the sample resolution with latent features, are transformed by a 1×1 convolution layer and ReLU activation [1] in the MPI encoder.

After obtaining the multi-scale feature after the MPI encoder, we add the feature to the decoder of diffusion UNet to provide spatial features. Experimental results in Tab. 3 will show that our 1×1 conv in MPI encoder is more effective than 3×3 conv, as the 1×1 conv with receptive field = 1 provides a spatial align feature to the latent feature in the diffusion UNet. In contrast, 3×3 conv is conducted in a camera frustum space rather than Euclid space, making an imprecise correspondence between 3D multiplane features and 2D image features. Moreover, using 3×3 conv to process 3D semantic MPI will introduce a large computational burden as the channel number increases from 3 channels of RGB to 256 planes. We display our 3D geometry and semantic control property in Fig. 3.

In summary, we chose MPIs as the representation because they (i) Incorporate lossless 3D information, including scene geometry rather than 2.5D depth. (ii) Provide spatially aligned conditional features that naturally extend the ControlNet framework from image level to 3D level. (iii) Capable of representing geometry and semantics including occluded elements.

3.4 Cross-View Attention

The sensor arrangement in a self-driving car usually requires a full surround view of cameras to capture the entire 360-degree environment. To effectively simulate the multi-view and subsequent multi-frame generation, zero-initialized [48] cross-view attention is integrated into the diffusion model to maintain consistency between views and frames. Following prior work [6, 35, 37, 38], each cross-view attention allows the target view to access information from its neighboring left and right views, thus training it using multi-view consistent images will enforce it to generate the same instance in the overlapping region of multi-view cameras.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V, \quad (6)$$

$$h_{out} = h_{in} + \sum_{i \in \{l, r\}} \text{Attention}(Q_{in}, K_i, V_i), \quad (7)$$

where l , and r is the camera view of left and right. Q_{in} and h_{in} denotes the query and the hidden state of input view. As for video generation, we use CogVideoX [46] as video diffusion backbone. Video results are provided in the supplementary.

3.5 Importance Reweighing

To deal with the extreme imbalance problem between foreground, background, and object categories, and also to ease the training, we propose three types of reweighting methods to improve the generation quality of foreground objects.

Progressive Foreground Enhancement To mitigate the complexity of the learning task, we propose a progressive reweighting method that incrementally enhances the loss associated with the foreground regions (based on semantics) as the training progresses. The detailed formulation is:

$$w(x, m, n) = \frac{(m-1)}{2} \cdot \left(1 + \cos\left(\frac{x}{n} \cdot \pi + \pi\right)\right) + 1, \quad (8)$$

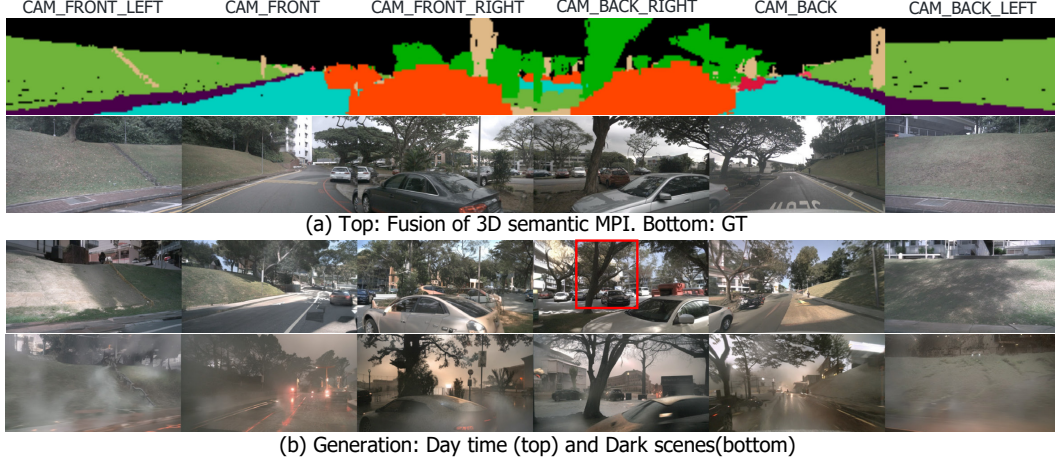


Figure 5: Visualizations of generated multi-view images. The generation conditions (occupancy labels) are from nuScenes validation set. We highlight that (i) Geometry alignment of trees in red rectangle in (b). (ii) Use text prompt to control high-level appearance.

where x is the current training step, m is the maximum value of weights that set at 2, and n is the total training steps. This approach is engineered to facilitate a learning trajectory that progresses from simplicity to complexity, thereby aiding in the convergence of the model. This curve can be interpreted as a cosine annealing but inverted to amplify the importance of the foreground region.

Depth-aware Foreground Reweighting In the meantime, we acknowledge the learning difficulty in different depth in 3D scenes. Following GeoDiffusion [3], we perform depth reweighing to foreground objects by adaptively assigning higher weights to farther foreground areas. This enables the model to focus more thoroughly on hard examples with depth-aware importance reweighing. Instead of using their exponential function to increase weights, we use our designed cosine function Eq. 8 for stability. Here x is the input depth value, and n is the maximum depth that set at 50.

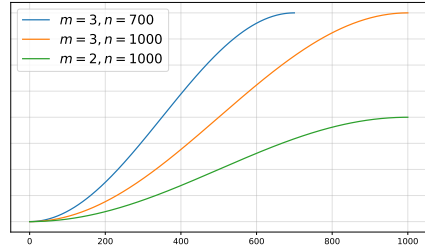


Figure 4: Visualizations of the reweighing function in Eq. 8.

CBGS Sampling To deal with the class imbalance problem in driving scenarios, where certain object categories appear infrequently, we employ the Class-Balanced Grouping and Sampling (CBGS) [52] to better handle the long-tailed classes. CBGS addresses the challenge of class imbalance by grouping and re-sampling training data to ensure each group has a balanced distribution of sample frequency across different object categories. This method reduces the bias towards more frequent classes and enables better generalization to rare scenarios.

3.6 Model Training

We separate the training of image generation and video generation. Our final objective function is formulated as a standard denoising objective with reweighing:

$$\mathcal{L} = \mathbb{E}_{\mathcal{E}(x), \epsilon, t} \|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|^2 \odot w, \quad (9)$$

where w is the multiplication of progressive reweighing and depth-aware reweighing.

4 Experiments

4.1 Dataset and Setups

We conduct our experiments on the nuScenes dataset [2], which is collected using 6 surrounded-view cameras that cover the full 360° field of view around the ego-vehicle. It contains 700 scenes for

Method	Train	Val	mIoU	barrier	bicycle	bus	car	cons. veh.	moto.	pedes.	traf. cone	trailer	truck	drive. suf.	other flat	sidewalk	terrain	manmade	vegetation
Oracle (FB-Occ)	Real	Real	39.3	45.4	28.2	44.1	49.4	25.9	28.8	28.0	27.7	32.4	37.3	80.4	42.2	49.9	55.2	42.0	37.7
SyntheOcc-Aug	Real+Gen	Real	40.3	45.4	27.2	46.6	49.5	26.4	27.8	28.4	29.4	34.0	37.2	81.3	46.0	52.4	56.5	43.3	38.9
MagicDrive	Real	Gen	13.4	0.7	0.0	11.8	32.4	0.0	6.6	2.8	0.3	2.6	19.6	60.1	12.1	26.2	23.4	15.5	12.8
ControlNet	Real	Gen	17.3	17.7	0.2	13.6	21.0	0.6	0.8	8.6	10.4	6.9	11.9	67.4	18.8	36.4	36.9	20.8	22.4
ControlNet+depth	Real	Gen	17.5	19.3	0.3	14.0	23.7	1.0	0.6	9.2	9.2	5.7	12.1	68.8	19.2	36.0	35.3	19.8	22.8
SyntheOcc-Gen	Real	Gen	25.5	32.6	13.8	27.7	33.4	7.5	6.5	15.7	16.5	16.5	25.6	74.3	24.5	39.4	40.5	28.6	28.8

Table 1: Downstream evaluation on the **nuScenes-Occupancy** validation set. Based on the used train and val data, two types of settings are reported. The first is to use generated training set to augment the real training set, and evaluate on the real validation set, denoted as Aug. The second is to use pretrained models trained on the real training datasets to test on the generated validation set, denoted as Gen.

226 training and 150 scenes for validation. We resize the original image from 1600×900 to 800×448 for
227 training. In our work, we use the occupancy label with a resolution of $0.2m$ from OpenOccupancy [34]
228 as condition input, while the benchmark of occupancy prediction uses a resolution of $0.4m$ from
229 Occ3D [30] dataset for its popularity.

230 **Networks** We use Stable Diffusion [27] v2.1 checkpoint as initialization and only train occupancy
231 encoder, cross-view attention. We adopt FB-Occ [19] as the target model for occupancy prediction
232 for its SOTA performance in this task. The pretrained checkpoint of FB-Occ is obtained from their
233 official repository. Since FB-Occ predicts occupancy using only single-frame images, we adopt
234 single-frame inference.

235 **Video generation** We implement the same architecture for CogVideoX [46] 2B for video generation.
236 Our model is capable of generating videos with 49 frames at a resolution of 480×256 pixels in 12hz.
237 Specifically, we set the first half of transformer blocks as our MPI encoder. The training process
238 is the same as the image generation. As only occupancy annotation of keyframes is available, we
239 perform linear interpolation of the MPI features to non-keyframes.

240 **Metrics** We use mIoU to measure the precision of occupancy prediction. We use FID [9] and
241 FVD [32] to measure the perceptual quality of our generation.

242 **Hyperparameters** We set $D = 256$, $d_{min} = 0$ and $d_{max} = 50$. The depth resolution of MPI is
243 thus higher than occupancy voxel. We train our model in 6 epochs with batch size = 8. The learning
244 rate is set at $2e^{-5}$. The training phase takes around 1 day using 8 NVIDIA A100 80G GPUs. We
245 use UniPC scheduler [49] with the classifier-free guidance (CFG) [10] that is set as 7.0. During
246 inference, we use 20 denoising steps for dataset generation.

247 **Baselines** We provide comparison in Tab. 1. ControlNet denotes we train a ControlNet using
248 an RGB semantic mask as the condition. ControlNet+depth denotes we add a depth channel after
249 the semantic mask to provide 2.5D depth information. The depth map rendered by occupancy is
250 normalized to [0-255] to accommodate the RGB value. The ControlNet+depth can be regarded as a
251 degradation of SyntheOcc which is reduced to a single plane. Then we evaluate MagicDrive since
252 it is the only open-sourced method in this area. MagicDrive separately encodes foreground and

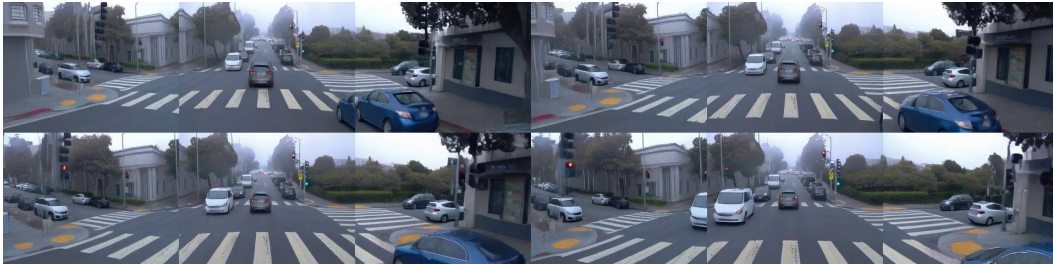


Figure 6: Video generation on Waymo dataset. More videos are provided in the supplementary.

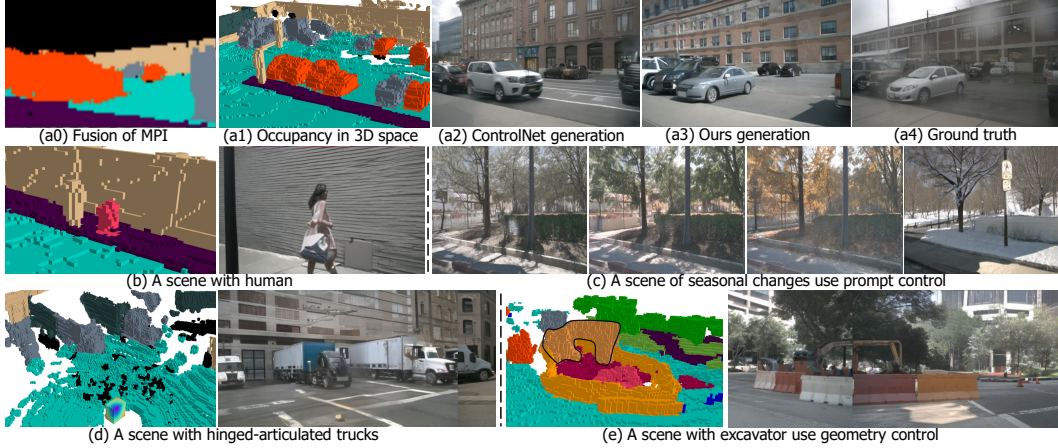


Figure 7: **Top row:** Comparison with ControlNet. We achieve a precise alignment between conditional labels and synthesized images, while ControlNet generates objects with incorrect pose due to ambiguous 2D condition. **Mid and Bottom row:** Visualizations of geometry-controlled image generation. We can faithfully generate objects with the desired topology in a specific 3D position.

background using prompt and BEV layout. Furthermore, we evaluate the image quality (FID [9]) of our method in Tab. 2. Compared with prior methods, we use a unified 3D representation that seamlessly handles foreground and background, surpassing them by a large margin.

4.2 Qualitative Results

High-level Control using Prompt In Fig. 5 (c,d) and Fig. 7 (c), we demonstrate the capability to employ user-defined prompts to generate images with specific weather conditions and high-level style. Although the nuScenes dataset doesn't contain rare weather images like snow and sandstorms, our method successfully conveys prior knowledge pretrained from stable diffusion to our scenes. Compared with visualization results in prior work like Fig. 8 of MagicDrive, our method shows better alignment with the text prompt, demonstrating the cross-domain generalization ability of our method.

3D Geometric Control Our flexible framework enables us to create novel scenes by manipulating voxels as displayed in Fig. 1 and Fig. 3. Basically, we can edit the occupied state and semantics of every voxel in our scenes for generation. We highlight that we can create a hinged-articulated truck and an excavator as shown in Fig. 7 (d,e). The generated excavator image exhibits a remarkable alignment with the input occupancy that is delineated by a black outline.

Long-tailed Scene Generation The flexibility of 3D semantic MPI has conferred significant advantages upon our approach. In the following, we create long-tail scenes that rarely occur in our real world for evaluation. In Fig. 1, we show that we manually add parallel traffic cones in front of the ego vehicle. This scene has never happened in the training dataset, but our geometric controllability provides us the capability to create such data. We then use the created scene to test autonomous driving systems such as end-to-end planner VAD [12] to validate its effectiveness. In this case, VAD successfully predicts correct waypoints with the high-level command 'turn left'. Moreover, in

Method	Condition Type	FID
BEVGen [29]	BEV map	25.54
BEVControl [41]	BEV map	24.85
DriveDreamer [33]	Box + FoV map	52.60
MagicDrive [6]	Box + BEV map	16.20
Panacea [37]	Box + FoV map	16.96
Ours	3D Semantic MPI	14.75

Table 2: Comparison of FID with previous methods on the nuScenes dataset.

MPI Encoder		Reweighting Method			Metric
Design		Progressive	Depth	CBGS	mIoU
3×3	-	-	-	-	21.96
1×1	-	-	-	-	23.05
1×1	✓	-	-	-	23.63
1×1	✓	✓	✓	-	24.40
1×1	✓	✓	✓	✓	25.50

Table 3: Ablation of different designs of the MPI encoder and reweighting methods.

275 supplementary, we generate long-tailed scenes with extreme weather such as snow and sandstorms,
276 and evaluate perception model on it to examine its generalizability of rare weather.

277 **Comparison with Baselines** In Fig. 7 (a), we visualize a comparison with ControlNet. We find
278 that ControlNet struggles to distinguish the overlapping instances in 2D-pixel space. This leads to the
279 two parked cars being merged into a single car with incorrect pose. In contrast, our 3D semantic MPIs
280 contain more than 2D semantic mask, but also account for complete scene geometry with occluded
281 parts. Together with our proposed MPI encoder and reweighing strategy, our framework yields a
282 realistic image generation with high-quality label alignment.

283 4.3 Quantitative Results

284 **Recognizability, Realism and Controllability Evaluation** To evaluate whether our generated
285 images aligned with given annotations, we provide Gen experiment in Tab. 1. Using the annotation of
286 val set, we synthesize a copy of val set’s images, then use perception model trained on real training set
287 to perform evaluation. The performance will be more effective as it is close to the oracle performance.
288 We find that local method (ControlNet) perform better than global method (MagicDrive).

289 **Data Augmentation for 3D Occupancy Prediction** Notably, we conduct experiments using our
290 synthesized dataset to enhance the real training set in Tab. 1. We first use the occupancy labels from
291 training set to create a synthetic training set. Then we modify the loading pipeline in perception model
292 to randomly sample images from real dataset or synthetic dataset and train network from scratch.
293 Therefore, our approach preserves the inherent training dynamics of the neural network by solely
294 modifying the training images, without any alteration to the number of training iterations or epochs.
295 As MagicDrive-Aug exhibits numerical overflow when training FB-Occ, which may attributed to
296 unsatisfactory recognizability, we have to omit it and only provide MagicDrive-Gen experiments.

297 As shown in Tab. 1, where SyntheOcc-Aug denotes the augmentation experiments using our generated
298 dataset, shows a satisfactory improvement over the prior state of the art. We emphasize that surpassing
299 the performance of the original dataset is not the primary objective of our work; rather, it is an
300 ancillary benefit that emerges from our framework for geometry-controlled generation.

301 **Ablations** In Tab. 3, we present ablation studies across several design spaces of our model, anal-
302 ogous to the Gen experiment in Tab. 1. We find that our designed MPI encoder of 1×1 conv has
303 significant improvement compared to the conventional 3×3 conv approach, and reweighing methods
304 demonstrate a consistent improvement. As a result, the improved image quality and label alignment
305 enable higher precision in downstream tasks. For more experiments, like Waymo dataset results, ab-
306 lation of plane numbers in MPIs, FVD evaluation, 3D detection evaluation, and view consistency
307 evaluation, we provide them in the supplementary.

308 **Occupancy-free Generation** To avoid relying on occupancy as input, we further employ a
309 generative model [13] to generate occupancy as condition for novel scene generation. This experiment
310 is presented in the supplementary.

311 5 Conclusion

312 In this paper, we propose **SyntheOcc**, an innovative image and video generation framework that
313 is empowered with geometry-controlled capabilities using occupancy. We introduce a novel 3D
314 representation, 3D semantic MPIs, to address the critical challenge of how to efficiently encode
315 occupancy. This representation not only preserves the authentic and complete 3D geometry details
316 with semantics, but also provides a spatial-align feature representation for 2D diffusion models. With
317 this property, our method enjoys photorealistic appearances and fine-grained 3D controllability, serves
318 as a generative data engine to enable a broad range of applications. Extensive experiments demonstrate
319 that our synthetic data facilitate the training for perception models on occupancy prediction, and
320 provide valuable corner case evaluation in a simulated world.

References

- [1] Abien Fred Agarap. Deep learning using rectified linear units (relu). [arXiv preprint arXiv:1803.08375](#), 2018.
- [2] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, 2020.
- [3] Kai Chen, Enze Xie, Zhe Chen, Lanqing Hong, Zhenguo Li, and Dit-Yan Yeung. Integrating geometric control into text-to-image diffusion models for high-quality detection data generation via text prompt. [arXiv preprint arXiv:2306.04607](#), 2023.
- [4] Jaeyoung Chung, Suyoung Lee, Hyeongjin Nam, Jaerin Lee, and Kyoung Mu Lee. Luciddreamer: Domain-free generation of 3d gaussian splatting scenes. [arXiv preprint arXiv:2311.13384](#), 2023.
- [5] Ruiyuan Gao, Kai Chen, Bo Xiao, Lanqing Hong, Zhenguo Li, and Qiang Xu. Magicdrivedit: High-resolution long video generation for autonomous driving with adaptive control. [arXiv preprint arXiv:2411.13807](#), 2024.
- [6] Ruiyuan Gao, Kai Chen, Enze Xie, Lanqing Hong, Zhenguo Li, Dit-Yan Yeung, and Qiang Xu. Magicdrive: Street view generation with diverse 3d geometry control. In *ICLR*, 2024.
- [7] Shenyan Gao, Jiazhi Yang, Li Chen, Kashyap Chitta, Yihang Qiu, Andreas Geiger, Jun Zhang, and Hongyang Li. Vista: A generalizable driving world model with high fidelity and versatile controllability. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [8] David Ha and Jürgen Schmidhuber. World models. [arXiv preprint arXiv:1803.10122](#), 2018.
- [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. *NeurIPS*, 2017.
- [10] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. [arXiv preprint:2207.12598](#), 2022.
- [11] Lukas Höllein, Ang Cao, Andrew Owens, Justin Johnson, and Matthias Nießner. Text2room: Extracting textured 3d meshes from 2d text-to-image models. In *ICCV*, 2023.
- [12] Bo Jiang, Shaoyu Chen, Qing Xu, Bencheng Liao, Jiajie Chen, Helong Zhou, Qian Zhang, Wenyu Liu, Chang Huang, and Xinggang Wang. Vad: Vectorized scene representation for efficient autonomous driving. In *ICCV*, 2023.
- [13] Jumin Lee, Sebin Lee, Changho Jo, Woobin Im, Juhyeon Seon, and Sung-Eui Yoon. Semcity: Semantic scene generation with triplane diffusion. [arXiv preprint arXiv:2403.07773](#), 2024.
- [14] Bohan Li, Jiazhe Guo, Hongsi Liu, Yingshuang Zou, Yikang Ding, Xiwu Chen, Hu Zhu, Feiyang Tan, Chi Zhang, Tiancai Wang, et al. Uniscene: Unified occupancy-centric driving scene generation. [arXiv preprint arXiv:2412.05435](#), 2024.
- [15] Kaican Li, Kai Chen, Haoyu Wang, Lanqing Hong, Chaoqiang Ye, Jianhua Han, Yukuai Chen, Wei Zhang, Chunjing Xu, Dit-Yan Yeung, et al. Coda: A real-world road corner case dataset for object detection in autonomous driving. In *ECCV*, 2022.
- [16] Leheng Li, Qing Lian, Luozhou Wang, Ningning Ma, and Ying-Cong Chen. Lift3d: Synthesize 3d training data by lifting 2d gan to 3d generative radiance field. In *CVPR*, 2023.
- [17] Xiaofan Li, Yifu Zhang, and Xiaoqing Ye. Drivingdiffusion: Layout-guided multi-view driving scene video generation with latent diffusion model. [arXiv preprint arXiv:2310.07771](#), 2023.
- [18] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *CVPR*, 2023.
- [19] Zhiqi Li, Zhiding Yu, David Austin, Mingsheng Fang, Shiyi Lan, Jan Kautz, and Jose M Alvarez. Fb-occ: 3d occupancy prediction based on forward-backward view transformation. [arXiv preprint arXiv:2307.01492](#), 2023.
- [20] Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela L Rus, and Song Han. Bevfusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation. In *2023 IEEE international conference on robotics and automation (ICRA)*, pages 2774–2781. IEEE, 2023.
- [21] William Ljungbergh, Adam Tonderski, Joakim Johnander, Holger Caesar, Kalle Åström, Michael Felsberg, and Christoffer Petersson. Neuroncap: Photorealistic closed-loop safety testing for autonomous driving. [arXiv preprint arXiv:2404.07762](#), 2024.
- [22] Jiachen Lu, Ze Huang, Jiahui Zhang, Zeyu Yang, and Li Zhang. Wovogen: World volume-aware diffusion for controllable multi-camera driving scene generation. [arXiv preprint arXiv:2312.02934](#), 2023.
- [23] Yifan Lu, Xuanchi Ren, Jiawei Yang, Tianchang Shen, Zhangjie Wu, Jun Gao, Yue Wang, Siheng Chen, Mike Chen, Sanja Fidler, et al. Infinicube: Unbounded and controllable dynamic 3d driving scene generation with world-guided video models. [arXiv preprint arXiv:2412.03934](#), 2024.

- [24] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- [25] Julian Ost, Fahim Mannan, Nils Thuerey, Julian Knodt, and Felix Heide. Neural scene graphs for dynamic scenes. In *CVPR*, 2021.
- [26] Can Qin, Shu Zhang, Ning Yu, Yihao Feng, Xinyi Yang, Yingbo Zhou, Huan Wang, Juan Carlos Niebles, Caiming Xiong, Silvio Savarese, et al. Unicontrol: A unified diffusion model for controllable visual generation in the wild. *arXiv preprint arXiv:2305.11147*, 2023.
- [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- [28] Liangchen Song, Liangliang Cao, Hongyu Xu, Kai Kang, Feng Tang, Junsong Yuan, and Yang Zhao. Roomdreamer: Text-driven 3d indoor scene synthesis with coherent geometry and texture. *arXiv preprint arXiv:2305.11337*, 2023.
- [29] Alexander Szwedlow, Runsheng Xu, and Bolei Zhou. Street-view image generation from a bird’s-eye view layout. *IEEE RAL*, 2024.
- [30] 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. *NeurIPS*, 2024.
- [31] Wenwen Tong, Chonghao Sima, Tai Wang, Li Chen, Silei Wu, Hanming Deng, Yi Gu, Lewei Lu, Ping Luo, Dahua Lin, et al. Scene as occupancy. In *ICCV*, 2023.
- [32] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Fvd: A new metric for video generation. *arXiv preprint*, 2019.
- [33] Xiaofeng Wang, Zheng Zhu, Guan Huang, Xinze Chen, and Jiwen Lu. Drivedreamer: Towards real-world-driven world models for autonomous driving. *arXiv preprint arXiv:2309.09777*, 2023.
- [34] Xiaofeng Wang, Zheng Zhu, Wenbo Xu, Yunpeng Zhang, Yi Wei, Xu Chi, Yun Ye, Dalong Du, Jiwen Lu, and Xingang Wang. Openoccupancy: A large scale benchmark for surrounding semantic occupancy perception. In *ICCV*, 2023.
- [35] Yuqi Wang, Jiawei He, Lue Fan, Hongxin Li, Yuntao Chen, and Zhaoxiang Zhang. Driving into the future: Multiview visual forecasting and planning with world model for autonomous driving. *arXiv preprint arXiv:2311.17918*, 2023.
- [36] Yi Wei, Linqing Zhao, Wenzhao Zheng, Zheng Zhu, Jie Zhou, and Jiwen Lu. Surroundocc: Multi-camera 3d occupancy prediction for autonomous driving. In *ICCV*, 2023.
- [37] Yuqing Wen, Yucheng Zhao, Yingfei Liu, Fan Jia, Yanhui Wang, Chong Luo, Chi Zhang, Tiancai Wang, Xiaoyan Sun, and Xiangyu Zhang. Panacea: Panoramic and controllable video generation for autonomous driving. *arXiv preprint arXiv:2311.16813*, 2023.
- [38] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *ICCV*, 2023.
- [39] Tianyi Yan, Dongming Wu, Wencheng Han, Junpeng Jiang, Xia Zhou, Kun Zhan, Cheng-zhong Xu, and Jianbing Shen. Drivingsphere: Building a high-fidelity 4d world for closed-loop simulation. *arXiv preprint arXiv:2411.11252*, 2024.
- [40] Jiazhi Yang, Shenyan Gao, Yihang Qiu, Li Chen, Tianyu Li, Bo Dai, Kashyap Chitta, Penghao Wu, Jia Zeng, Ping Luo, Jun Zhang, Andreas Geiger, Yu Qiao, and Hongyang Li. Generalized predictive model for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [41] Kairui Yang, Enhui Ma, Jibin Peng, Qing Guo, Di Lin, and Kaicheng Yu. Bevcontrol: Accurately controlling street-view elements with multi-perspective consistency via bev sketch layout. *arXiv preprint arXiv:2308.01661*, 2023.
- [42] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 2023.
- [43] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. *arXiv preprint arXiv:2310.06114*, 2023.
- [44] Yuanbo Yang, Yifei Yang, Hanlei Guo, Rong Xiong, Yue Wang, and Yiyi Liao. Urbangiraffe: Representing urban scenes as compositional generative neural feature fields. In *ICCV*, 2023.
- [45] Ze Yang, Yun Chen, Jingkan Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. Unisim: A neural closed-loop sensor simulator. In *CVPR*, 2023.

- 432 [46] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi
433 Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert
434 transformer. [arXiv preprint arXiv:2408.06072](#), 2024.
- 435 [47] Hong-Xing Yu, Haoyi Duan, Junhwa Hur, Kyle Sargent, Michael Rubinstein, William T Freeman, Forrester
436 Cole, Deqing Sun, Noah Snavely, Jiajun Wu, et al. Wonderjourney: Going from anywhere to everywhere.
437 [arXiv preprint arXiv:2312.03884](#), 2023.
- 438 [48] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion
439 models. In *ICCV*, 2023.
- 440 [49] Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-corrector
441 framework for fast sampling of diffusion models. *NeurIPS*, 2023.
- 442 [50] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification:
443 Learning view synthesis using multiplane images. [arXiv preprint arXiv:1805.09817](#), 2018.
- 444 [51] Yunsong Zhou, Michael Simon, Zhenghao Peng, Sicheng Mo, Hongzi Zhu, Minyi Guo, and Bolei Zhou.
445 Simgen: Simulator-conditioned driving scene generation. [arXiv preprint arXiv:2406.09386](#), 2024.
- 446 [52] Benjin Zhu, Zhengkai Jiang, Xiangxin Zhou, Zeming Li, and Gang Yu. Class-balanced grouping and
447 sampling for point cloud 3d object detection. [arXiv preprint arXiv:1908.09492](#), 2019.

Appendix

A Potential Discussion

To help a comprehensive understanding of our paper, we discuss intuitive questions that might be raised.

How to define geometric control? In our paper, we refer the geometric controllable generation as using a voxel grid in 3D space to control the image generation. Although the voxel is a quantized representation of the 3D world, when the resolution goes larger, it can already faithfully represent the geometry detail of scenes. Currently, we are limited by the precision of ground truth labels. The $0.2m$ occupancy grid is a tensor of $500 \times 500 \times 40$ that cover a space in x-axis spanning $[-50m, 50m]$, y-axis spanning $[-50m, 50m]$, z-axis spanning $[-5m, 3m]$. In the future, we plan to explore a higher resolution of geometric control to refine our generation.

Can 3D semantic MPI extend to other representations beyond occupancy? Except for occupancy, several other 3D representations can be expressed by 3D semantic MPI, such as mesh, dense point clouds, and even 3D boxes or HD maps. The underlying mechanism is to cast several slices of multi-plane images at different depths to retrieve geometric information. Our application scope is wide, and we left them for future work. As a result, our 3D semantic MPI can be regarded as a general 3D conditioning representation to benefit a wide spectrum of practical systems. These encompass but are not limited to 3D generation such as text2room [11], RoomDreamer [28], WonderJourney [47], and LucidDreamer [4], each of which stands to benefit from the rich geometric context provided by our approach.

Occupancy is complex. How to edit occupancy for controllable generation? We agree that occupancy is more complex than the 3D box, but it provides a more nuanced scene description. To ease the editing, we provide a strategy that disentangles the foreground control and background control in occupancy data. If we want to edit a car’s trajectory, we can keep the background occupancy unchanged and select the car’s first frame occupancy using the 3D box. During the following frames, we remove foreground occupancy and simply place our foreground target’s occupancy in a certain location using trajectory. By doing so, we only add minor steps by using occupancy but provide more precise 3D control, which makes it a favorable choice for conditioning. We provide a user example of adding traffic cones in our supplementary video.

B Long-Tailed Scene Evaluation

In this section, we explore to use SyntheOcc to create long-tailed scenes for downstream evaluation. This also stands for evaluating our model using several corner cases. Similar to the SytheOcc-Gen experiment, we generate a synthetic validation set but use prompt control to manipulate weather patterns or the intensity of illumination.

In Tab. 4, we observe that all kinds of extreme weather lead to a degradation in performance. This observation underscores the limitations of the perception model in terms of its generalizability to infrequent weather scenarios. Among them, we find that foggy, rainy, and day night exert the most severe impact, as they contribute to a large reduction in visibility. To improve the generalizability to handle various weather conditions, future work can leverage our generated data to cover the long-tailed scenes, or use adversarial search to find severe scenes based on our framework.

Furthermore, we perform long-tailed scene evaluation in Fig. 8. We display the failure of the downstream model VAD [12] in our synthetic long-tailed scene. In this case, we simulate a foggy environment that the dense fog obscures the majority of the ego view. Our experiment reveals that due to the lack of training images of foggy scenes, VAD erroneously predicts waypoints that would

Scenes	Sandstorm	Snow	Foggy	Rainy	Day night	Day time (raw data)
FB-Occ mIOU	22.88	18.25	10.29	9.71	9.95	25.50

Table 4: Experiments of downstream evaluation on long-tailed scenes with extreme weather.

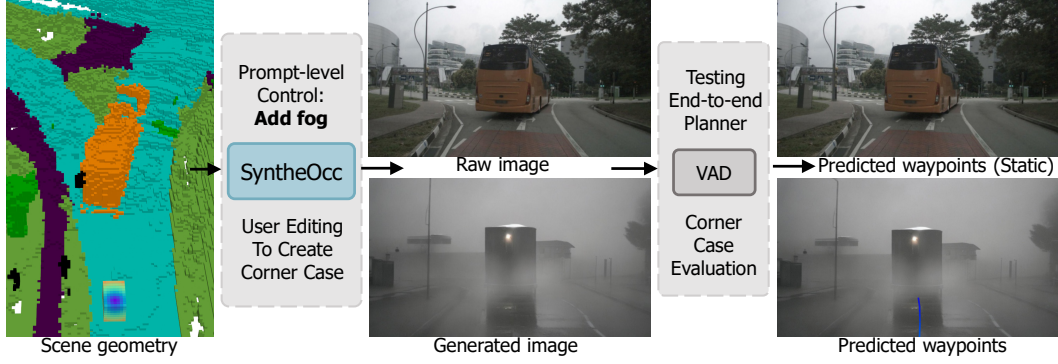


Figure 8: Use **SyntheOcc** to create long-tailed scenes for testing. **Top:** In the ordinary scene of a bus placed in front of the ego vehicle, the end-to-end planner VAD [12] predicts future waypoints without movement, thus not plotted in the image. **Bottom:** By harnessing the prompt-level control in our framework, we simulate a scene with the same layout but filled with fog. VAD predicts wrong waypoints that will collide with the bus.

492 result in a collision with the bus. This experiment elucidates the boundaries and failure cases of the
 493 VAD model [12]. It exposes the limitations of the system under certain conditions, thereby providing
 494 insights into scenarios where the model’s performance may be compromised.



Figure 9: Comparison with baselines.

495 C Ablation of plane number of MPIs

496 In our proposed 3D semantic MPIs, the number of planes is a hyperparameter that affects the precision
 497 of 3D representation. The plane number can be regarded as the 3D resolution in depth axis. The
 498 larger the plane number, the MPI will contain more details. We find that an increase in the number of
 499 planes is associated with improved accuracy in downstream tasks. This finding denotes that more
 500 condition information leads to better downstream task performance.

Number of Planes	96	128	256
FB-Occ mIOU	23.36	24.28	25.50

Table 5: Ablation of the number of multi-plane images.

D Qualitative Comparison with Baselines and SOTA

In Fig. 9, we conduct a qualitative comparison of our method against MagicDrive, ControlNet, and ControlNet+depth. We find that all the methods display a satisfactory image quality, as they build upon the foundation of the stable diffusion model. The generation of MagicDrive fails to synthesize barriers as shown in the bottom row. ControlNet struggles to generate objects with the correct pose solely from only 2D conditions as shown in the second row. ControlNet+depth, a degradation of our method, an enhancement over ControlNet in terms of alignment, nevertheless suffers from a loss of finer detail in scenes with heavy occlusion, as shown in the human of the third row. Our method, in contrast, aims to address these challenges and provide a more nuanced and accurate generation of complex scenes.

E Data Augmentation Experiments

Data Augmentation using ControlNet We provide experiments that use ControlNet and ControlNet+depth to enable data augmentation. This experiment is analogous to the Aug experiment. In the experiment of ControlNet and ControlNet+depth, due to the potential for input-generation ambiguity, the augmented data could lead to the propagation of inaccurate gradients, thereby affecting the training process. These experiments demonstrate that our approach outperforms the ControlNet baseline in terms of effectiveness.

Methods	No aug	ControlNet	ControlNet+depth	SyntheOcc
FB-Occ mIOU	39.3	39.0	39.1	40.3

Table 6: Experiments of evaluating the data augmentation effects using different generative model.

Evaluate on 3D Detection We assess the accuracy of 3D detection using the BEVFusion [20]. This experiment corresponds to the Generation experiment, with the distinction that BEVFusion is employed for evaluating 3D detection precision. Owing to the effective pixel-aligned alignment offered by our method, SyntheOcc yields superior detection accuracy compared to previous studies.

Methods	MagicDrive-mAP	MagicDrive-NDS	SyntheOcc-mAP	SyntheOcc-NDS
Results	20.8	30.2	22.3	31.3

Table 7: Experiments of evaluating the generation quality using 3D detection accuracy.

F Extend to Video Generation

As described in the main paper, we perform video generation based on a strong and open-source video generation backbone CogVideoX [46]. We use the same architecture of the MPI encoder as in image experiments in stable diffusion. Our generation results can be found in the supplementary video. In practice, we use the keyframe annotation of the occupancy label of the nuScenes dataset to train our video model. After that, we upsample the MPI feature to all frames using linear interpolation. We further evaluate the Fréchet Video Distance (FVD) score [32] to evaluate the video generation quality in Tab. 8. Attributed to our commendable controllable image generation quality, SyntheOcc achieves competitive performance that is on par with other models.

Methods	DriveGAN	DriveDreamer	DrivingDiffusion	Panacea	Ours
FVD	502	340	332	139	34

Table 8: Experiments of evaluating the quality of video generation on nuScenes dataset.

Given that our primary contribution does not lie in video generation, this experiment serves as a proof of concept, demonstrating the potential adaptability of our framework. Future research may extend our methodology to facilitate the generation of longer video sequences, thereby expanding the scope and applicability of our framework.

G The Influence of the Amount of Augmented Data

As SyntheOcc is capable of generating an infinite number of synthetic data, we investigate the influence of the amount of augmented data on downstream tasks in Tab. 9. We find that when our augmented data is expanded from one-fold to two-fold of the training dataset, the performance of perception model slightly decreases. This may indicate the generated data has an optimal ratio for downstream tasks. Due to limited computational resources, we only experiment with a limited amount of ratio. Future work can conduct more thorough experiments to find a universal theorem.

Amount of Augmented Data	0 (no augmentation)	1	2
FB-Occ mIOU	39.3	40.3	40.1

Table 9: Ablation of the amount of augmented data.

H Additional Experiments

Evaluation on multi-view consistency. We evaluate view consistency using View Consistency Score (VCS) from BEVGen [29]. VCS is calculated using the confidence of matching points between different views (large better).

Method	BEVGen	MagicDrive	Panacea	SyntheOcc
VCS	6.24	6.45	6.53	6.80

Table 10: Evaluation of view consistency on nuScenes dataset.

Occupancy-free image generation: avoid the dependence on occupancy data. To avoid relying on occupancy as input, we further train a generative model SemCity [13] (on nusenes) to generate occupancy as a condition for novel scene generation. We use a two-stage pipeline to mitigate the issue of obtaining occupancy. First, we train a diffusion model to generate novel occupancy in 3D space. Second, we leverage the generated occupancy as the condition and convert it to a multi-plane image. Finally, our SyntheOcc generates images or videos by the occupancy instruction, as shown in Fig. 10. We further evaluate image generation quality as FID in Tab. 11. We find that they achieve similar results as they share the same image generative model. In this manner, we show that our generation framework can achieve an unconditional manner that does not need to rely on existing occupancy labels that are hard to annotate.



Figure 10: Top: Projection of occupancy generated by SemCity. Bottom: Generated frame.

Source	SemCity Occupancy	GT Occupancy
FID	15.28	14.75

Table 11: Evaluation of FID across using different sources of occupancy.

Experiment with Waymo Open Dataset. To further showcase our effectiveness, we train our model on Waymo dataset as shown in Fig. 11. We use the 3 front-view cameras to train our video generation model. We achieve reasonable FID = 17.2 and FVD = 86.3, demonstrating our effectiveness and scalability.

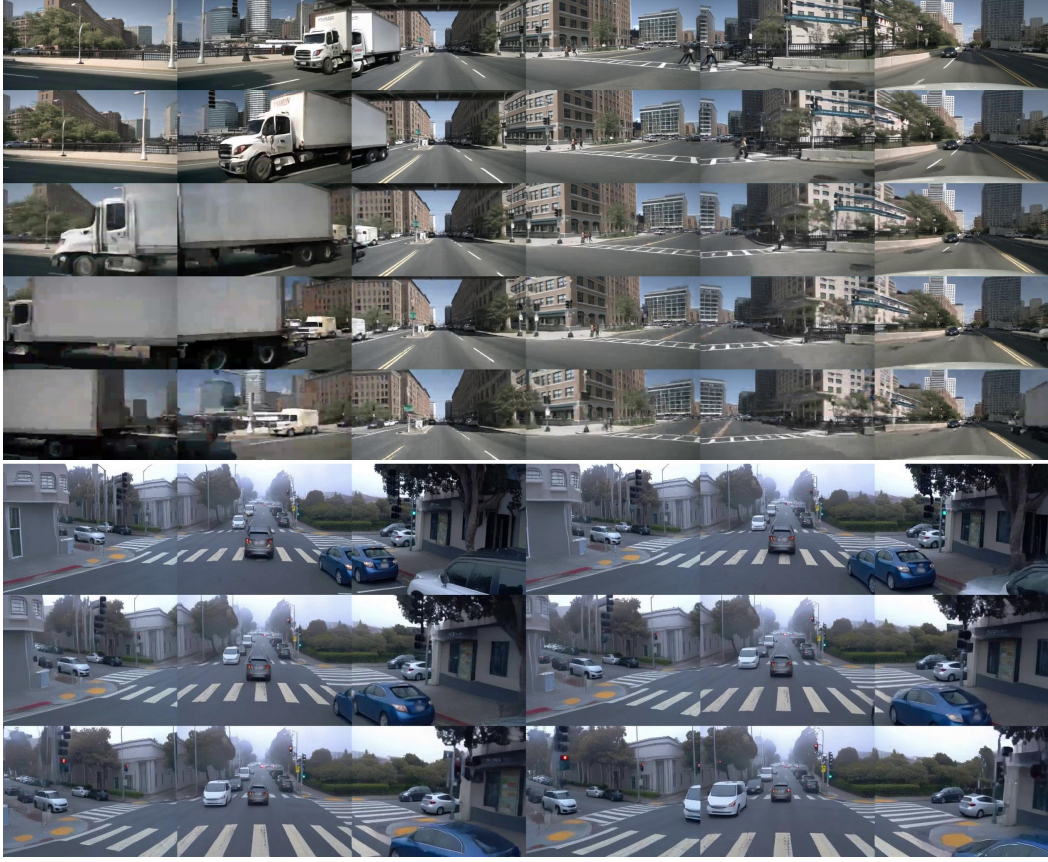


Figure 11: Video generation results. (see supplementary for videos.)

Comparison with Existing Work. Drive-WM utilizes a rough 3D box and layout as the conditional input, failing to capture objects' detailed shapes and irregular elements such as sidewalk and terrain. In contrast, our method enables finer and more precise 3D geometric control by using occupancy. WovoGen, on the other hand, performs a future prediction task with a distinct setting from ours. It leverages historical occupancy as input to predict future feature volumes for conditions. WovoGen's object guide can be regarded as a simple ControlNet implementation that uses a single plane without semantics. Moreover, WovoGen lags behind our method in terms of FID, FVD, and lacks critical experiments on downstream applications for 3D occupancy prediction, as it only provides 3D detection results.

More baselines. While our improvement may not be substantial in numerical, it aligns with the objective laws of the occupancy prediction task. For instance, SurroundOcc achieved only a marginal improvement of less than one point compared to previous sota (31.49 vs 30.86). Moreover, our method achieves a satisfactory improvement (40.3 vs 39.3) under a strong baseline (FB-Occ), which demonstrates our effectiveness and robustness. As suggested by the reviewer, we add experiments that apply data augmentation to CVT-Occ, which achieve reasonable improvement in mIOU (41.46 vs 40.34).

Method	No-aug	SyntheOcc
CVT-Occ mIOU	40.34	41.46

Table 12: Evaluation of data augmentation on occupancy prediction.

575 I Limitation and Broader Impacts

576 **Long-tailed Scene Generation** In this paper, we investigate a series of long-tailed scene generation
577 and corner case evaluations such as rare layout and extreme weather in Sec. B. Future work can
578 extend our framework to (i) Synthesize more samples for tail classes to boost performance. (ii)
579 Generate or replicate large-scale databases of corner cases [15] for robust perception.

580 **Closed-loop Simulation** Given the underlying diverse and controllable image generation of our
581 method, it would be advantageous and valuable to extend our work to a broader domain such as closed-
582 loop simulation [21, 45], to enable high-fidelity autonomous systems testing. This line of work can
583 be conducted by utilizing motion conditions to generate future frames as in world model [22, 35, 43],
584 or by explicitly modeling scene graph as in the case of UniSim [25, 45] and NeuroNCAP [21].

NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”.**
- **Keep the checklist subsection headings, questions/answers and guidelines below.**
- **Do not modify the questions and only use the provided macros for your answers.**

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [Yes]

Justification: Please find this part in Sec. 3.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Please find this part in Sec. I.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper does not include theoretical results.

Guidelines: Do not have theoretical results.

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Please find this part in Sec. 4.

Guidelines:

- The answer NA means that the paper does not include experiments.

- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Please find this part in Sec. 4.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: Please find this part in Sec. 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[Yes\]](#)

Justification: Please find this part in Sec. 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [\[Yes\]](#)

Justification: Please find this part in Sec. 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.

- 789 • The paper should disclose whether the full research project required more compute
790 than the experiments reported in the paper (e.g., preliminary or failed experiments that
791 didn't make it into the paper).

792 **9. Code Of Ethics**

793 Question: Does the research conducted in the paper conform, in every respect, with the
794 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

795 Answer: [Yes]

796 Justification: It should be fine.

797 Guidelines:

- 798 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
799 • If the authors answer No, they should explain the special circumstances that require a
800 deviation from the Code of Ethics.
801 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
802 eration due to laws or regulations in their jurisdiction).

803 **10. Broader Impacts**

804 Question: Does the paper discuss both potential positive societal impacts and negative
805 societal impacts of the work performed?

806 Answer: [Yes]

807 Justification: Please find this part in Sec. I.

808 Guidelines:

- 809 • The answer NA means that there is no societal impact of the work performed.
810 • If the authors answer NA or No, they should explain why their work has no societal
811 impact or why the paper does not address societal impact.
812 • Examples of negative societal impacts include potential malicious or unintended uses
813 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
814 (e.g., deployment of technologies that could make decisions that unfairly impact specific
815 groups), privacy considerations, and security considerations.
816 • The conference expects that many papers will be foundational research and not tied
817 to particular applications, let alone deployments. However, if there is a direct path to
818 any negative applications, the authors should point it out. For example, it is legitimate
819 to point out that an improvement in the quality of generative models could be used to
820 generate deepfakes for disinformation. On the other hand, it is not needed to point out
821 that a generic algorithm for optimizing neural networks could enable people to train
822 models that generate Deepfakes faster.
823 • The authors should consider possible harms that could arise when the technology is
824 being used as intended and functioning correctly, harms that could arise when the
825 technology is being used as intended but gives incorrect results, and harms following
826 from (intentional or unintentional) misuse of the technology.
827 • If there are negative societal impacts, the authors could also discuss possible mitigation
828 strategies (e.g., gated release of models, providing defenses in addition to attacks,
829 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
830 feedback over time, improving the efficiency and accessibility of ML).

831 **11. Safeguards**

832 Question: Does the paper describe safeguards that have been put in place for responsible
833 release of data or models that have a high risk for misuse (e.g., pretrained language models,
834 image generators, or scraped datasets)?

835 Answer: [NA]

836 Justification: Our paper poses no such risks.

837 Guidelines:

- 838 • The answer NA means that the paper poses no such risks.

- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [\[Yes\]](#)

Justification: Please find this part in Sec. 4.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[NA\]](#)

Justification: Our paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[NA\]](#)

Justification: Our paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: Our paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: LLMs are not core methods in our research.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.