

DrivingSphere: Building a High-fidelity 4D World for Closed-loop Simulation

Supplementary Material

818 6. Implement details

819 In this section, we provide detailed implementation settings
820 to facilitate the reproducibility of this work. Specifically,
821 we elaborate on the datasets used, evaluation metrics, the
822 architecture of our generation models.

823 6.1. Datasets

824 **Scene Generation.** Our experiments are primarily based
825 on the nuScenes [3] dataset. For the scene generation task,
826 we use the nuScenes-OpenOcc dataset as the data source.
827 Each scene provides complete occupancy annotations and
828 BEV maps, supporting the evaluation of both static and dy-
829 namic elements. We also use GPT4-V to obtain the scene
830 description as the text prompts.

831 **Video Generation.** Following prior works [13, 28, 44],
832 we use a standard split of 700 scenes for training and 150
833 scenes for validation. Each sequence is recorded at 12 Hz
834 and lasts approximately 20 seconds, with annotations pro-
835 vided at 2 Hz. To train higher-frequency models, we inter-
836 polate the sequences to generate 12 Hz annotations and train
837 models with both 2 Hz and 12 Hz versions. To achieve fine-
838 grained control over the generated scenes and traffic actors’
839 appearance, we utilize GPT-4v to generate detailed scene
840 captions and object captions. These captions provide high-
841 level semantic descriptions of the overall scene and detailed
842 attributes for each traffic actor, enabling precise guidance
843 during video generation. Additionally, we track each fore-
844 ground actor within a single sequence to assign a unique
845 ID, ensuring appearance consistency across frames in the
846 generated video sequence.

847 **Open-loop and Closed-loop Settings.** We support
848 two map environments, *singapore-onenorth* and *boston-*
849 *seaport*, both aligned with the DriveArena platform. A total
850 of 100 simulation sequences are defined as the validation
851 set to evaluate both open-loop and closed-loop generation
852 performance.

853 6.2. Evaluation Metrics

854 **Frechet Video Distance (FVD).** This metric evaluates the
855 visual quality and temporal consistency of generated video
856 clips, following prior methods [13, 45].

857 **Mean Average Precision (mAP) and NuScenes Detection
858 Score (NDS).** We adopt mAP and NDS to assess the detec-
859 tion accuracy on generated data to validate the fidelity.

860 We utilize metrics from [8, 50] for comprehensive eval-
861 uation in both open-loop and closed-loop setups: Pro-
862 gressive Driving Metric Suite (PDMS): Initially proposed
863 by NavSim, PDMS evaluates trajectory outputs at each

timestep based on the following criteria:

864 No Collisions (NC): Measures whether the agent avoids
865 collisions with road users. Drivable Area Compliance
866 (DAC): Assesses whether the agent remains within the driv-
867 able area. Ego Progress (EP): Quantifies how effectively
868 the agent progresses along its intended route. Time-to-
869 Collision (TTC): Evaluates the safety of the trajectory in
870 terms of time remaining before a collision occurs. Com-
871 fort (C): Ensures the smoothness of the driving trajectory,
872 minimizing abrupt accelerations or turns. Arena Driving
873 Score (ADS): ADS integrates trajectory-level performance
874 (PDMS) with route completion to provide a holistic metric.
875 The Route Completion (Rc) is defined as the percentage of
876 the total route distance completed by the agent, where
877 $Rc \in [0, 1]$. ADS is particularly suited for closed-loop
878 evaluation, as it considers both safety and consistency. For
879 instance, collisions or deviations from the road terminate
880 the simulation, making ADS an effective differentiator of
881 agent performance. 882

883 6.3. Model Details

884 **OccDreamer.** For the scene tokenizer \mathcal{F}_{VAE}^{occ} , we follow-
885 ing [41, 57] and train a 3D occupancy VAE, which takes
886 the occupancy data with the size of $192 \times 192 \times 16$. \mathcal{F}_{VAE}^{occ}
887 compresses the data \mathcal{S}_k to a latent space $\mathcal{Z}^{\mathcal{S}_k}$ with a dimen-
888 sion of $48 \times 48 \times 4$ and the channel is set to 8. For the
889 BEV map, a pre-trained encoder [16, 17] encodes the BEV
890 representation at the same resolution as the latent feature.

891 As for the denoiser ϵ_θ^s and the ControlNet branch ϵ_ϕ^s , we
892 use 3D U-Net [16, 17] as the backbone for the 3D input
893 data. For the basic scene generation model, we train ϵ_θ^s and
894 ϵ_ϕ^s for 60k iterations on 8 NVIDIA A800 GPUs. For the
895 scene extension version, we freeze ϵ_θ^s and fine-tune ϵ_ϕ^s with
896 extra channels to take the partial scene as the condition. In
897 the inference stage, we adopt DDIM [16, 17] with 100 steps
898 sampling. Additionally, we set the classifier-free guidance
899 scale as 7 for the condition.

900 **Videodreamer.** Our implementation is based on the Open-
901 Sora codebase [59], initialized with pre-trained weights.
902 The training process is carried out on 8 NVIDIA A800
903 GPUs, comprising 30k iterations. As for the 4D occupancy
904 encoder $\mathcal{F}_{VAE}^{4Docc}$, we also borrow the network architecture
905 from [41, 57]. $\mathcal{F}_{VAE}^{4Docc}$ takes 4D occupancy data as the in-
906 put and extract the embedding. The number of DiT blocks
907 is set to 26 and $N = 13$. For inference, we utilize recti-
908 fied flow [59] with a classifier-free guidance scale of 7.0,
909 performing 30 sampling steps to generate videos at vari-
910 ous resolutions from 480p to 1080p. For the open-loop and
911 closed-loop evaluation, we generate the video of 4 frames

Methods	Downsampling Scale	IoU _(↑)	mIoU _(↑)
OccWorld [57]	$H/4 \times W/4 \times T$	62.29	66.38
OccSora [41]	$H/8 \times W/8 \times T/8$	27.4	37
<i>DrivingSphere</i> _{4D}	$H/4 \times W/4 \times T$	93.1	73.89
Semcity [22]	-	95.8	76.9
<i>DrivingSphere</i> _{3D}	$H/4 \times W/4$	97.2	86.81

Table 5. **Quantitative results of Occupancy Tokenizer for Occupancy Reconstruction.** *DrivingSphere*_{4D} indicates $\mathcal{F}_{\text{VAE}}^{\text{4Docc}}$ in Sec. 3.2 while *DrivingSphere*_{4D} indicates $\mathcal{F}_{\text{VAE}}^{\text{occ}}$ in Sec. 3.1.

with $f = 3$ frames as the condition while generating long videos, we generate the 16 frames video sequence and generally use $f = 4$ frames as the condition.

7. Additional Quantitative Results

7.1. Scene Reconstruction

To validate the effectiveness of our Occupancy VAE, we conduct scene reconstruction experiments on the nuScenes validation set. These experiments evaluate both 3D and 4D scene reconstruction, providing a comprehensive analysis of the model’s capability. As shown in Tab. 5, for 3D scene reconstruction, we compare our trained 3D Occupancy VAE with SemCity, which is utilized as the occupancy tokenizer in Section 3.1. The results demonstrate that our Occupancy VAE achieves superior performance, highlighting its ability to encode and reconstruct 3D occupancy data effectively. For 4D scene reconstruction, we benchmark against OccWorld and OccSora, widely regarded as state-of-the-art methods for large-scale 4D occupancy generation. The results clearly show that our Occupancy VAE outperforms these approaches across all evaluation metrics, establishing new benchmarks for 4D scene reconstruction quality. These significant improvements are primarily attributed to the carefully designed network architecture, including the Projection Module and Expansion & Squeeze Strategy, as well as meticulously tuned experimental parameters. Such architectural innovations enable the model to capture fine-grained spatial and temporal information, ensuring accurate and efficient reconstruction of both static and dynamic elements within the scenes. This experiment underscores the robustness and effectiveness of the proposed framework in handling complex 3D and 4D scene representations.

7.2. Video Generation

To further validate the capabilities of VideoDreamer, we align our experimental settings with state-of-the-art video generation methods, ensuring a fair comparison. As presented in Tab. 6, we employ BEVFusion as the detector to quantitatively evaluate the visual fidelity of the generated videos. The results demonstrate that our method achieves superior performance, highlighting its ability to generate

Methods	FVD	mAP _(↑)	NDS _(↑)
RealData [3]	-	62.29	66.38
MagicDrive [13]	-	12.30	23.32
DriveDreamer [43]	340.8	-	-
Panacea [45]	139	11.58	22.31
Drive-WM [44]	122.7	20.66	-
<i>DrivingSphere</i> w/o <i>W</i>	121.4	17.34	26.21
<i>DrivingSphere</i>	103.4	22.71	31.79

Table 6. **Comparison of SOTA video generation methods on nuScenes validation set.** We use BEVFusion as the 3D detector. ‘w/o *W*’ indicates that the model uses no occupancy but uses the 2D sketch as the condition.

high-quality and visually coherent driving scenarios. Additionally, we conduct an ablation study by introducing the configuration ‘w/o *W*’, which uses only 2D sketches as conditions without incorporating occupancy data. This ablation effectively isolates the contribution of the 4D driving world to the video generation process. The results clearly illustrate the significant improvement in visual fidelity when occupancy data is integrated, confirming the critical role of the occupancy condition in enhancing the realism and consistency of generated video sequences. This experiment underscores the robustness of our framework in producing visually accurate driving videos and its ability to leverage multi-modal conditions effectively.

8. Additional Visualtion Results

In this section, we provide more quality visualization results and a video is also attached in the materials for better visualization of temporal results.

8.1. Scene Generation

In Fig. 7, we present a comparison between the occupancy scenes generated by our method, SemCity, and real-world data. The visual results clearly demonstrate that our method achieves significantly higher fidelity compared to SemCity, closely approximating the structural and semantic layout of real-world data. It is important to note that SemCity is an unconditional generation method, and as such, its outputs are unpaired with the real data used for comparison. In contrast, our method leverages conditions, ensuring consistency with the road structures and semantic layouts of the real data. This alignment highlights the strength of our approach in generating occupancy scenes that are not only visually realistic but also semantically coherent, demonstrating its suitability for tasks requiring precise scene understanding and reconstruction.

8.2. Video Generation

Controllable Video Generation In Fig. 11, we showcase the results of video generation spanning 40 frames. The

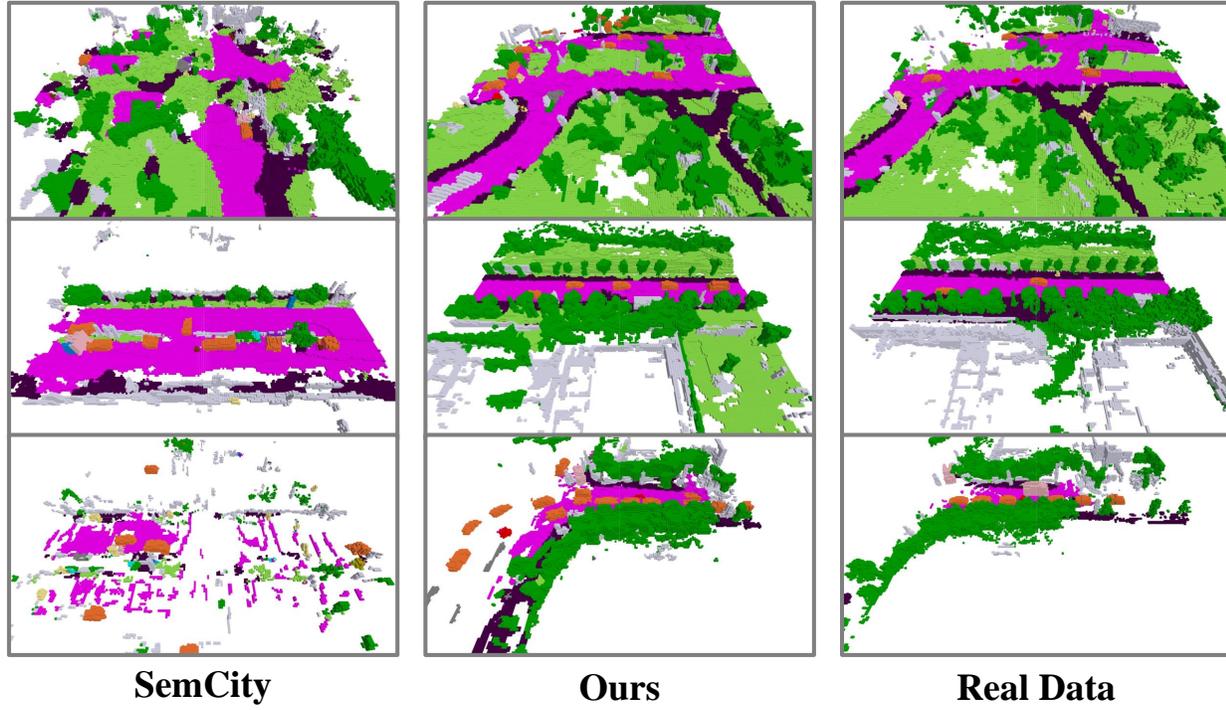


Figure 7. Comparison between Semcity [22], *DrivingSphere* and Real Data.

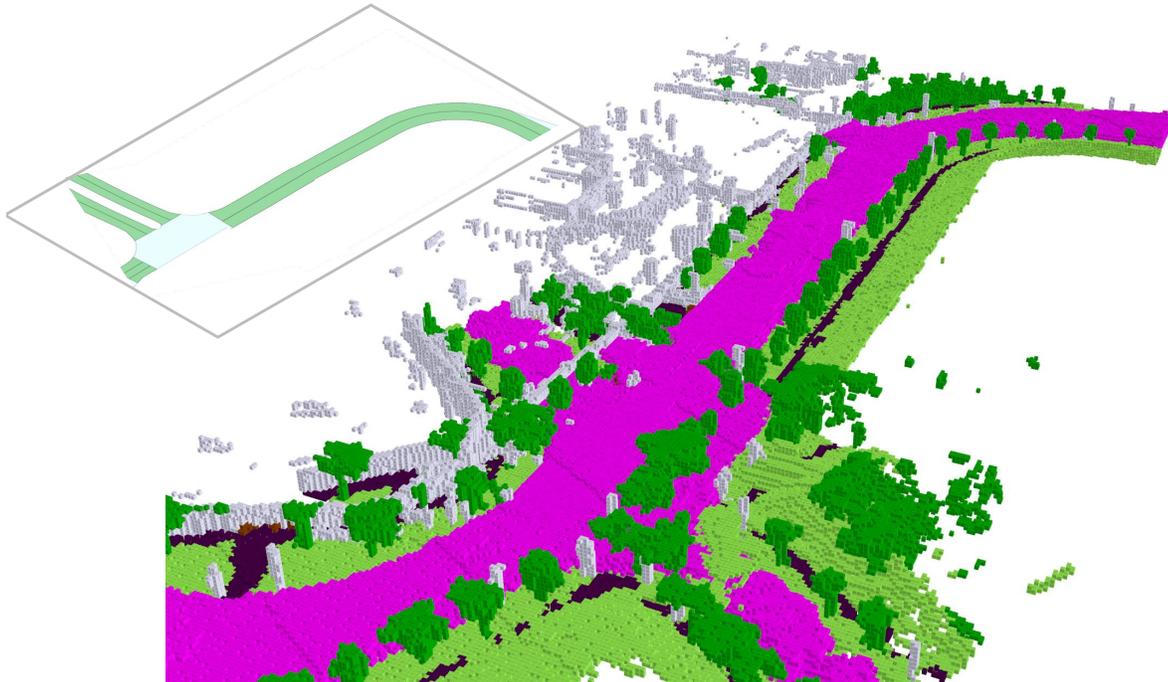


Figure 8. **Composited Driving World in a specific area.** We adopt Scene Generation and Scene Extension in Sec. 3.1 to obtain a big static background.

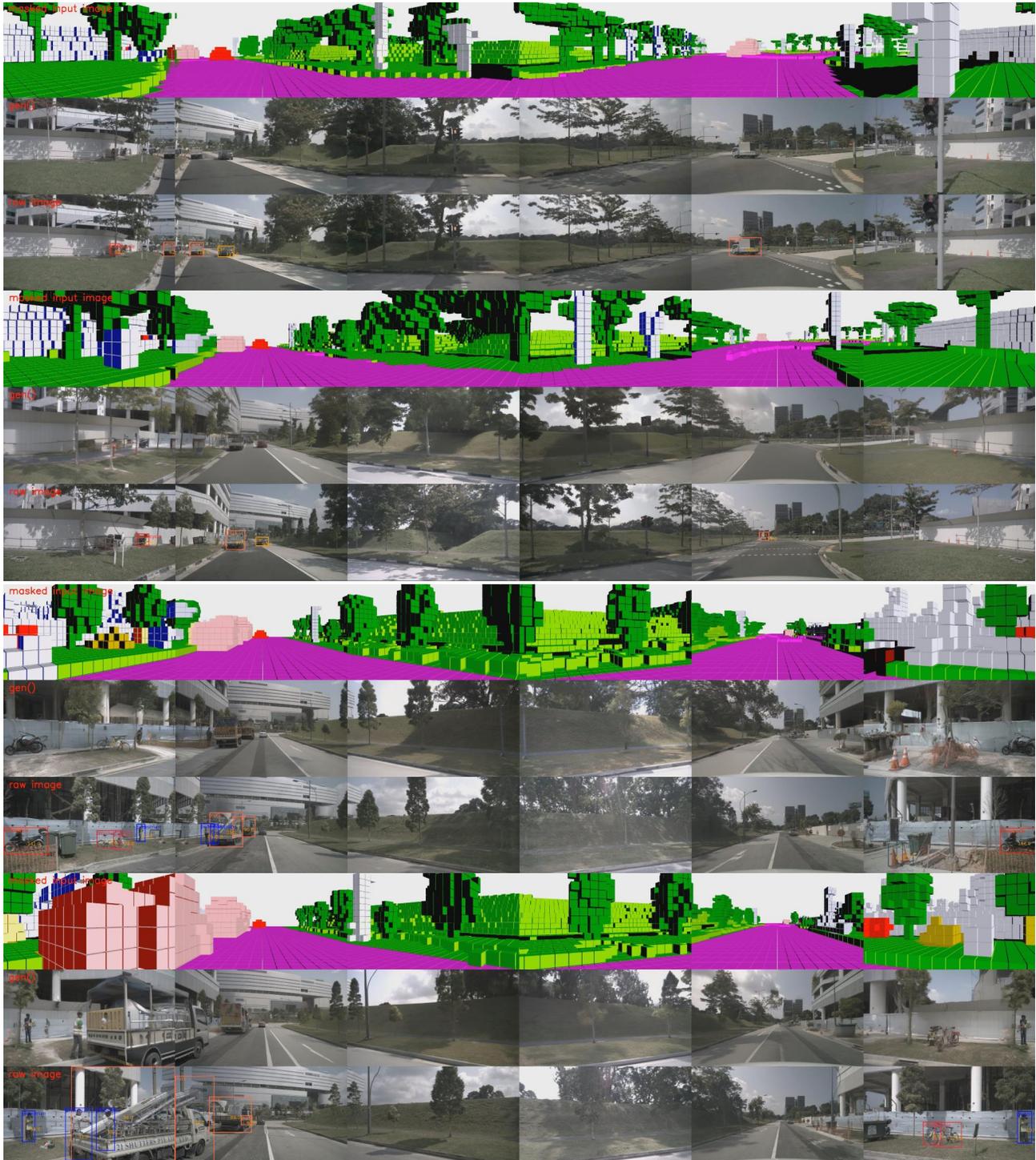


Figure 9. **Generated Video sequences in nuScene.** Top: Occupancy condition, Middle: Our generated video, Bottom: Ground truth video sequence.

987 generated sequences demonstrate the effectiveness of our
 988 method in accurately modeling not only the occlusions and
 989 depth relationships of foreground objects but also in pre-

cisely controlling the generation of non-direct traffic partic-
 ipants, such as trees, buildings, and man-made landmarks.

Our approach leverages the accurate control provided by

990
 991
 992

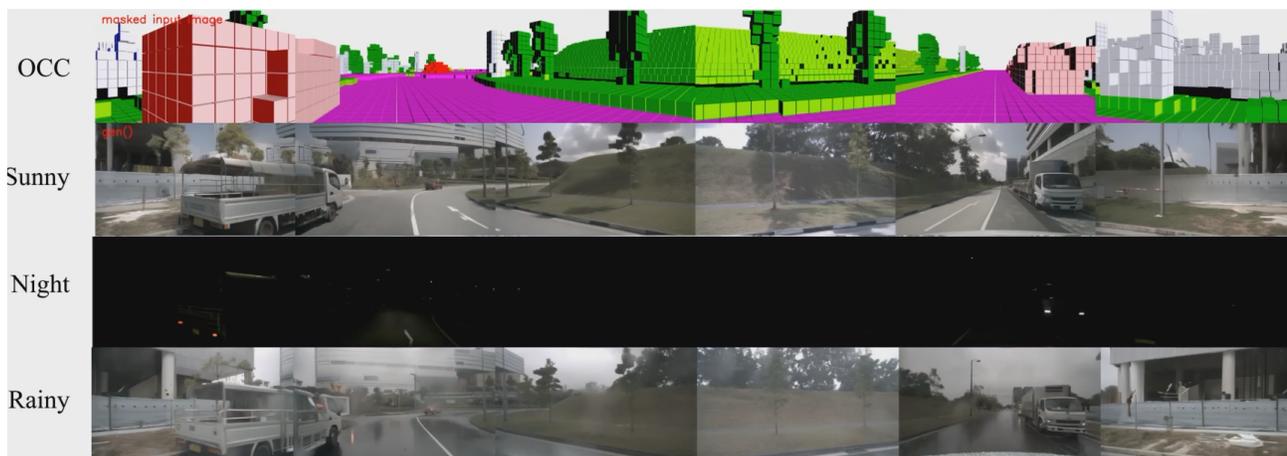


Figure 10. **Controllable Generation with scene captions.** The visual results vary with the given scene description.

993 occupancy data, enabling consistent and realistic representation of both dynamic and static elements within the scene.
994 This capability highlights the robustness of our model in
995 generating complex and coherent driving environments over
996 extended temporal horizons, making it well-suited for real-
997 world applications that require high fidelity and detailed
998 scene understanding.
999

1000 **Simulation Results**

1001 **Long-term Video Generation** We also provide a demo
1002 of ultra-long video generation on private data (refer to the
1003 attached video). The demo showcases a video generated at
1004 10 Hz over a duration of 1 minute, resulting in an impressive
1005 600 frames of continuous generation.

1006 **9. Limitations and Future Work**

1007 Optimizing the computational pipeline for generating and
1008 rendering 4D occupancy and video data will be a key fo-
1009 cus. Techniques such as model pruning, quantization, and
1010 adaptive sampling will be explored to reduce computational
1011 overhead without compromising fidelity. Additionally, en-
1012 abling real-time rendering capabilities will make the system
1013 more practical for online validation.

1014 Expanding the diversity of simulated environments is
1015 critical for robustness testing. Future work will aim to
1016 model a broader range of conditions, including extreme
1017 weather (e.g., heavy rain, snow, and fog), varying road
1018 geometries, and rare traffic scenarios. This enhancement
1019 will allow the simulation to evaluate autonomous driving
1020 systems under more comprehensive and challenging condi-
1021 tions.

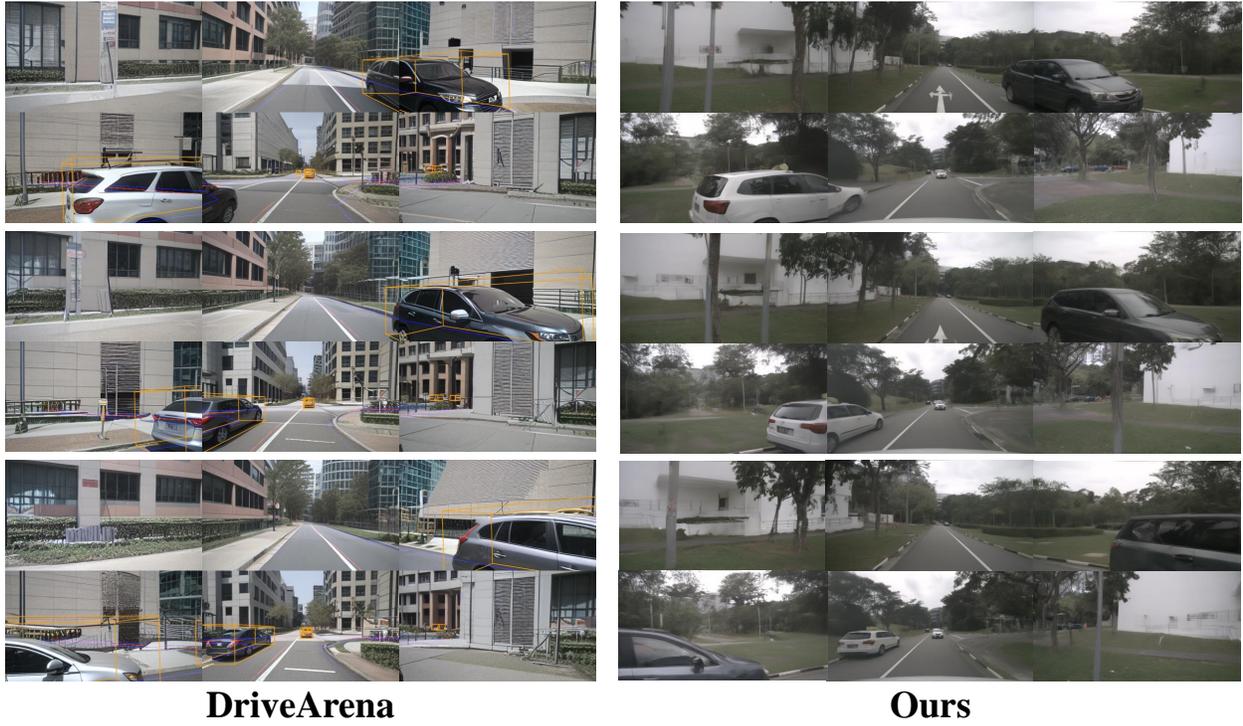


Figure 11. **Comparison with DriveArena [50].** The visual output of DriveArena and *DrivingSphere* on the same route demonstrates superior temporal and spatial consistency in generated simulations.