

GEM: Generating LiDAR World Model via Deformable Mamba

Supplementary Material

001 A. Appendix

002 In this Appendix, we provide additional information. In
003 Sec. A.1, we introduce the process of projecting LiDAR
004 point clouds onto range maps; In Sec. A.2, we provide a
005 more detailed explanation of how we utilize the Mamba al-
006 gorithm and the diffusion algorithm; In Sec. A.3, we intro-
007 duce the detailed architectures of discriminator; In Sec. A.4,
008 we provide additional visualizations of future prediction re-
009 sults generated by our GEM; In Sec. A.5, we present recon-
010 struction results from our LiDAR scene tokenizer, including
011 both 32-line and 64-line LiDAR data; In Sec. A.6, we show
012 more results on controllable generation, including counter-
013 factual reasoning and BEV-layouts controlled generation;
014 In Sec. A.7, we analyze GEM’s capability for long-term fu-
015 ture prediction and present visualizations of the results; In
016 Sec. A.8, we discuss the limitations and broader impacts of
017 GEM and provide possible solutions.

018 Additionally, we also upload videos in the supplement-
019 ary materials demonstrating the quality of reconstructed
020 data and generated sequences.

021 A.1. LiDAR Data Representation

022 LiDAR point clouds captured in outdoor environments typi-
023 cally exhibit high sparsity and irregularity. To address these
024 challenges, the range map has emerged as an effective rep-
025 resentation that successfully mitigates the issues of sparsity
026 and irregularity in LiDAR data, as demonstrated in prior
027 studies [2, 7, 8]. While efficient feature extraction for range
028 maps remains under-explored, their adoption in LiDAR
029 generation research has become mainstream [9, 10, 13, 16].

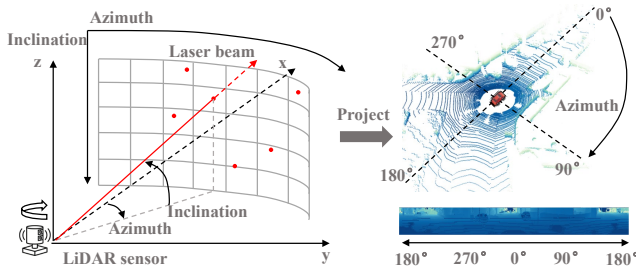


Figure 1. The conversion process from the LiDAR point cloud to the range map.

030 Compared to voxel and BEV representations that suffer
031 from information loss caused by limited resolution and
032 computational cost [15], a key advantage of the range map
033 representation is its ability to preserve nearly all original
034 information from the raw point cloud data while maintain-
035 ing low computational cost [8]. As shown in Figure 1, we

transform the coordinate space of each point, projecting it
from Cartesian coordinates $\mathbf{r} = (x, y, z) \in \mathbb{R}^3$ to spherical
coordinates $\mathbf{s} \in (\theta, \phi, d)$:

$$\begin{aligned} d &= \sqrt{x^2 + y^2 + z^2}, \\ \theta &= \arccos(z/\sqrt{x^2 + y^2 + z^2}), \\ \phi &= \text{atan2}(y, x), \end{aligned} \quad (1)$$

where θ is the inclination, ϕ is the azimuth, d is the depth.
For a 64-line LiDAR, we set the dimensions of range map
to 64×1024 , and for a 32-line LiDAR, we set the dimen-
sions of range map to 32×1024 . A row in the range map
corresponds to a beam circle obtained from a full rotation
scan by the LiDAR sensor, and a column in the range map
corresponds to the central ray acquired from the emission
of the LiDAR sensor [14].

A.2. More Detailed Algorithmic Procedures

In this section, we first elaborate on the specific implemen-
tation of Mamba in our framework, followed by a detailed
explanation of how the diffusion algorithm operates within
our GEM.

Each row of a range map corresponds to one full rota-
tion of LiDAR scanning points. Partitioning range maps
into patches lacks clear physical significance, whereas each
row naturally preserves the structural information of LiDAR
data. Therefore, unlike previous methods that apply Mamba
by scanning image patches, we perform scanning point-by-
point along the LiDAR beams, which better aligns with the
inherent characteristics of point clouds.

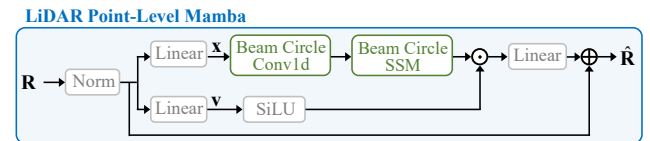


Figure 2. The architecture of the Mamba we employ. Unlike patch-level Mamba used in image processing, our Mamba performs scanning circle-by-circle directly on LiDAR points, effectively preserving the structural information of LiDAR data.

Specifically, taking the range map \mathcal{R} as input, to align
with the imaging characteristics of LiDAR data, we use
convolutional layers for preliminary feature modeling, then
flatten it pixel by pixel along rows to obtain $\mathbf{R} \in \mathbb{R}^{K \times J}$,
where K and J represent the number of points and chan-
nels. As shown in Fig. 2, within the Mamba process, the
input \mathbf{R} is first projected to embeddings \mathbf{x} and \mathbf{v} via linear

068 layers:

$$069 \quad \mathbf{x} = \text{Linear}_{\mathbf{x}}(\mathbf{R}), \mathbf{v} = \text{Linear}_{\mathbf{v}}(\mathbf{R}). \quad (2)$$

070 Then, to enable input-independent reasoning, we introduce
071 selective parameters \mathbf{B} , \mathbf{C} and a discretization parameter
072 $\Delta \in \mathbb{R}^{K \times D}$:

$$073 \quad \begin{aligned} \mathbf{B} &= \text{Linear}_{\mathbf{B}}(\mathbf{x}), \mathbf{C} = \text{Linear}_{\mathbf{C}}(\mathbf{x}) \\ \Delta &= \log(1 + \exp(\text{Linear}_{\Delta}(\mathbf{x}) + \text{bias}_{\Delta})), \end{aligned} \quad (3)$$

074 where bias_{Δ} is a learnable bias parameter of Δ . Next,
075 to apply the concept of continuous dynamic systems to dis-
076 crete inputs, we use the zero-order hold to obtain discretiza-
077 tion $\bar{\mathbf{A}}, \bar{\mathbf{B}} \in \mathbb{R}^{K \times D}$:

$$078 \quad \bar{\mathbf{A}} = \exp(\mathbf{A} \cdot \Delta), \bar{\mathbf{B}} = (\Delta \cdot \mathbf{A})^{(-1)} \cdot (\exp(\Delta \cdot \mathbf{A}) - \mathbf{I}), \quad (4)$$

079 where $\mathbf{A} \in \mathbb{R}^{K \times D}$ is the learnable parameters. Subse-
080 quently, we perform the Mamba scan for each step $k \in$
081 $[1, \dots, K]$ until $\hat{\mathbf{R}} \in \mathbb{R}^{K \times D}$ is obtained:

$$082 \quad \begin{aligned} h_0 &= 0, h_k = \bar{\mathbf{A}}_k \cdot h_{k-1} + \bar{\mathbf{B}}_k \cdot \mathbf{x}_k, \hat{\mathbf{R}}_k = \mathbf{C}_k \cdot h_k, \\ \hat{\mathbf{R}} &= \text{MLP}([\hat{\mathbf{R}}_1, \dots, \hat{\mathbf{R}}_L]^T) \cdot \text{SiLU}(\mathbf{v}) + \text{Norm}(\mathbf{R}), \end{aligned} \quad (5)$$

083 where MLP consists of linear and normalization layers. Fina-
084 lly, after several channel modification operations, we can
085 obtain the compressed latent $\mathbf{z} \in \mathbb{R}^{h \times w \times C}$, where h, w, C
086 means the number of height, width and channels.

087 Give a latent feature $\mathcal{Z} \in \mathbb{R}^{\tau \times h \times w \times C}$ with τ frames
088 of \mathbf{z} . We construct GEM following the latent diffusion
089 paradigm [11, 12], which mainly contains two parts. (1)
090 Forward process: given LiDAR series $\mathcal{Z} = [\mathcal{Z}^p, \mathcal{Z}^f]$, this
091 process transforms the future frames \mathcal{Z}^f into pure Gaus-
092 sian noise in an iterative manner. The addition of noise at
093 timestep t is defined as: $\mathcal{Z}_t^f = \sqrt{\alpha_t} \mathcal{Z}_{t-1}^f + \sqrt{1 - \alpha_t} \epsilon_t$,
094 where $\alpha_t = 1 - \beta_t$, β_t is the noise weight, and the ϵ_t is
095 the added Gaussian noise. Through iterative derivation [5],
096 we can obtain $\mathcal{Z}_t^f = \sqrt{\bar{\alpha}_t} \mathcal{Z}_0^f + \sqrt{1 - \bar{\alpha}_t} \epsilon_t$, where $\bar{\alpha}_t =$
097 $\alpha_t \cdot \alpha_{t-1} \dots \alpha_1$, and $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. This process can be fur-
098 ther expressed as: $q(\mathcal{Z}_t^f | \mathcal{Z}_0^f) = \mathcal{N}(\mathcal{Z}_t^f; \sqrt{\bar{\alpha}_t} \mathcal{Z}_0^f, (1 - \bar{\alpha}_t) \mathbf{I})$.
099 (2) Reverse process: we need to recover the data from the
100 noise distribution.

$$101 \quad q(\mathcal{Z}_{t-1}^f | \mathcal{Z}_t^f, \mathcal{Z}_0^f) = \mathcal{N}(\mathcal{Z}_{t-1}^f; \tilde{\mu}_t(\mathcal{Z}_t^f, \mathcal{Z}_0^f), \tilde{\beta}_t \mathbf{I}), \quad (6)$$

102 in which, $\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t$, $\tilde{\mu}_t(\mathcal{Z}_t^f, \mathcal{Z}_0^f) = \frac{1}{\sqrt{\alpha_t}} (\mathcal{Z}_t^f -$
103 $\frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \hat{\epsilon}_t)$, and $\hat{\epsilon}_t$ is predicted by our LiDAR world model.

104 The training objective of timestamp t is to estimate the
105 noise $\hat{\epsilon}_t$ in Eq. (6) based on $\mathcal{Z}_t = [\mathcal{Z}^p, \mathcal{Z}_t^f]$. After training,
106 given an initial noise latent $\hat{\mathcal{Z}}_T^f \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ we can progres-
107 sively recover it over T steps to produce the final prediction,
108 where T is the predefined denoising step.

A.3. Detailed Architectures of Discriminator 109

110 During the training of the LiDAR scene tokenizer, the MSE
111 loss function tends to generate blurred and averaged results,
112 as it minimizes pixel-level errors while overlooking the
113 plausibility of local structures [4]. Relying solely on MSE
114 loss makes it challenging to reconstruct high-frequency de-
115 tails in range maps, which are critical for features such as
116 object contours.

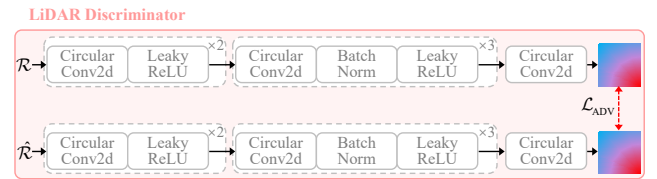


Figure 3. The architecture of the discriminator we employ. The application of a discriminator assists the LiDAR scene tokenizer in reconstructing high-frequency details, which is crucial for range maps since these details correspond to object contours, road structures, and other critical information.

117 To address this, we introduce a discriminator to im-
118 prove the model’s reconstruction accuracy. As illustrated
119 in Fig. 3, the discriminator comprises circular convolutions,
120 normalization layers, and activation functions. Since the re-
121 constructed $\hat{\mathcal{R}}$ has already undergone Mamba processing,
122 which effectively preserves LiDAR structural information
123 during feature extraction and reconstruction, the discrimi-
124 nator focuses on assessing the perceptual realism of $\hat{\mathcal{R}}$, thus
125 eliminating the need for Mamba operations that are primar-
126 ily geared toward precision reconstruction. The discrimi-
127 nator emphasizes perceptual “realism” rather than pixel-
128 wise alignment, thereby compelling the tokenizer to learn
129 the true distribution on the data manifold instead of merely
130 fitting the average behavior of the training samples.

A.4. More Visualizations of World Modeling 131

132 In this section, we present more visualizations for 3s predic-
133 tions and compare them with 4D-Occ [6]. On nuScenes [1]
134 dataset, the 3s future comprises 6 frames, while on KITTI
135 Odometry [3], it consists of 5 frames.

136 **On nuScenes.** Fig. 4 supplements the main-text visual-
137 izations by presenting the complete sequence from 0.5s to
138 3.0s at 0.5s intervals alongside the historical 6 frames ob-
139 servations, which have omitted the 0.5s, 1.5s, and 2.5s fu-
140 ture frames due to space constraints before. Additionally,
141 Figure 5 provides further visualizations on the nuScenes
142 dataset. It can be observed that our method achieves more
143 accurate predictions of objects, with clearer contours. In
144 contrast, competing methods exhibit blurred or even miss-
145 ing object boundaries. Additionally, our method provides
146 more precise predictions of LiDAR beams, while compet-
147 ing methods show significant noise. These visualizations
148 further demonstrate our superiority on nuScenes.

149 **On KITTI Odometry.** Fig. 6 supplements the main-text
150 visualizations by presenting the complete sequence from
151 0.6s to 3.0s at 0.6s intervals alongside the historical 5
152 frames observations, which have omitted the 1.2s, and 2.4s
153 future frames due to space constraints before. Addition-
154 ally, Figure 7 provides further visualizations on the KITTI
155 Odometry dataset. These visualizations demonstrate our su-
156 periority on KITTI Odometry. It is evident that our method
157 produces clearer object predictions with fewer noisy points.
158 As the distance increases, our approach maintains superior
159 LiDAR beam accuracy, whereas competing methods exhibit
160 progressively more blurring.

161 **A.5. Visualizations of Reconstruction**

162 This section demonstrates our LiDAR scene tokenizer’s
163 reconstruction capability, achieving compression ratios of
164 2.67 (32-line) and 5.32 (64-line) while maintaining high-
165 fidelity reconstruction. As shown in Figs. 8 and 9, our
166 LiDAR scene tokenizer can accurately reconstruct scene
167 structures and object details from LiDAR data, maintain-
168 ing high reconstruction accuracy even in challenging scenar-
169 ios such as rainy conditions. Specifically, in rows 1 to
170 6 of Figure 8, we demonstrate diverse road conditions, in-
171 cluding straight roads, turns, and intersections. It can be
172 seen that our method accurately reconstructs small objects
173 such as cars and pedestrians, as well as road structures. In
174 rows 7 to 10 of Fig. 8, we present the reconstruction re-
175 sults of our LiDAR scene tokenizer for rainy scenes. It can
176 be observed that, owing to our precise capture of LiDAR
177 data structural information, our method maintains high re-
178 construction accuracy even in sparser and noisier scenarios.
179 Fig. 9 demonstrates the reconstruction results of 64-line Li-
180 DAR data under various road conditions. It can be observed
181 that our method achieves exceptionally high fidelity in re-
182 constructing both object details and road structures.

183 **A.6. More Visualizations of Controlled Generation**

184 This section presents additional visualizations to demon-
185 strate our model’s capabilities in controlled generation, cov-
186 ering both counterfactual reasoning and BEV-guided gener-
187 ation. As shown in Figs. 10 and 11, we present the current
188 and past frames alongside a comparison between original
189 predictions and counterfactual reasoning. Fig. 10 illustrates
190 the transition from a left turn to emergency braking, while
191 Fig. 11 demonstrates the shift from stationary waiting to a
192 right turn. The results indicate that our method thoroughly
193 comprehends the driving world, possesses strong imaginat-
194 ive capabilities, and can generate plausible scenarios based
195 on given counterfactual actions. Fig. 12 demonstrates the
196 editing of future frames under BEV layout control, includ-
197 ing object addition, deletion, and repositioning. It can be
198 observed that GEM accurately accomplishes these editing
199 tasks.

A.7. Visualizations of Long-Sequence Generation 200

201 Generating long-term sequences is a ultimate goal for world
202 models. In this section, we present visualizations that sug-
203 gests the promising potential of GEM for long-horizon pre-
204 diction. In detail, GEM first generates a future sequence
205 based on observed past frames and subsequently uses this
206 output as input, thereby iteratively expanding the prediction
207 horizon. As shown in Fig. 13, GEM remarkably maintains
208 considerable accuracy for up to 38 seconds before a notice-
209 able performance decline at 46.8s. Notably, under the ex-
210 treme challenge of the 58 seconds horizon, the model pre-
211 serves only basic structural information. It is worth empha-
212 sizing that GEM is the only existing LiDAR world model
213 capable of supporting minute-level future observations pre-
214 diction with meaningful accuracy.

A.8. Limitations and Broader Impacts 215

216 Despite GEM achieving excellent results in LiDAR world
217 modeling, it still has some shortcomings. Firstly, controlled
218 generation including BEV-guided generation and counter-
219 factual reasoning is still not mature enough for real world
220 application. We attribute this limitation primarily to in-
221 sufficient training data, and believe that scaling the dataset
222 could potentially enhance the capacity. Secondly, despite
223 demonstrating an ability for long-term generation which is
224 not observed in previous methods, we observe a growing
225 discrepancy between predictions and ground truth as time
226 grows. We attribute this to compounding errors accumu-
227 lated during iterative generation. Addressing the impact of
228 these compounding errors on long-term performance will
229 be a key focus of our future research. Thirdly, GEM cur-
230 rently cannot predict sudden environmental changes, such
231 as abrupt rain or entering foggy areas. Collecting more
232 diverse long-sequence data and data that includes weather
233 variations would help alleviate the aforementioned issues.

234 In terms of broader impacts, GEM may potentially be
235 utilized by technical personnel in economically developed
236 regions for customized urban data simulation, indirectly re-
237 placing some labor-intensive data collection jobs. This may
238 inadvertently widen the wealth gap, disrupt regional devel-
239 opment balance, and lead to certain social inequities. As
240 a potential countermeasure, broader social impact assess-
241 ments should be conducted, and relevant authorities ought
242 to establish comprehensive technical application frame-
243 works to better harness technology for societal benefit.

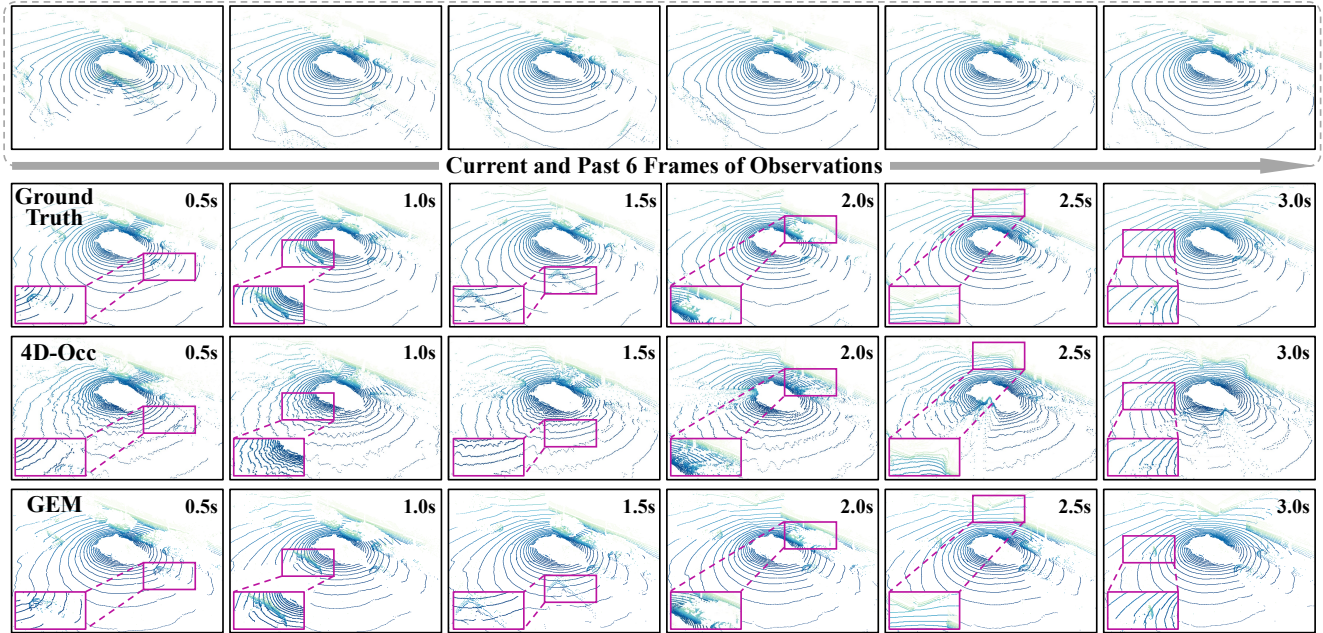


Figure 4. More visual comparisons with competitive LiDAR world models on nuScenes [1].

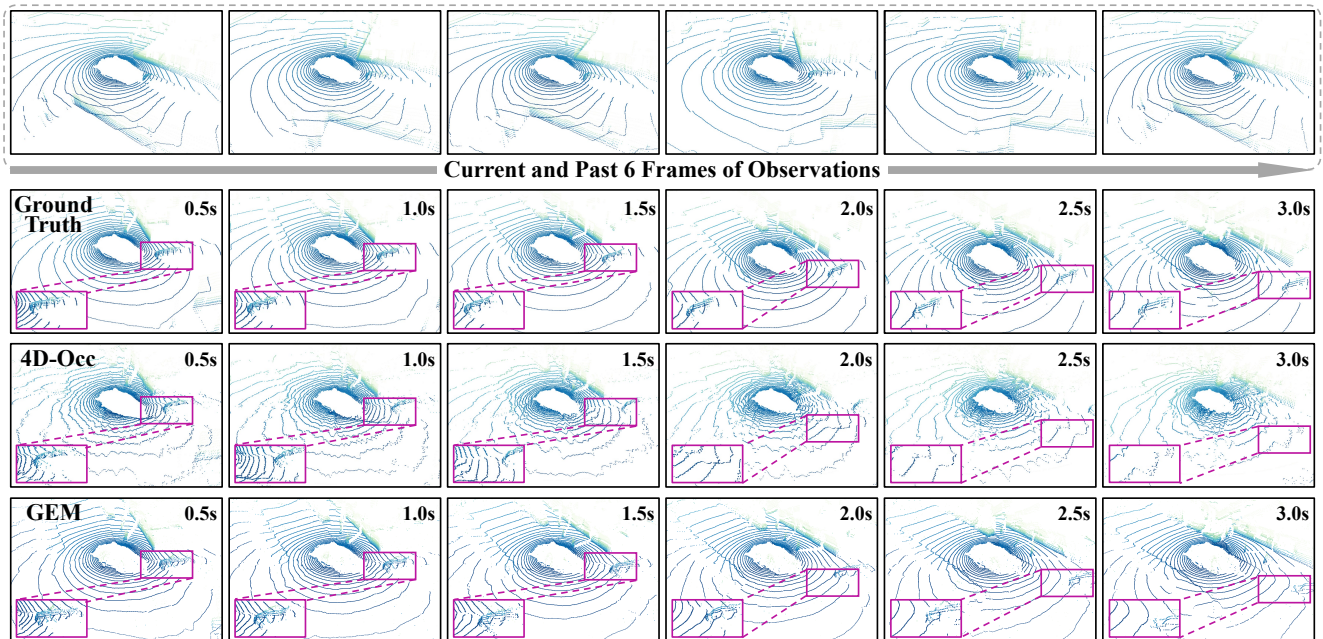


Figure 5. More visual comparisons with competitive LiDAR world models on nuScenes [1].

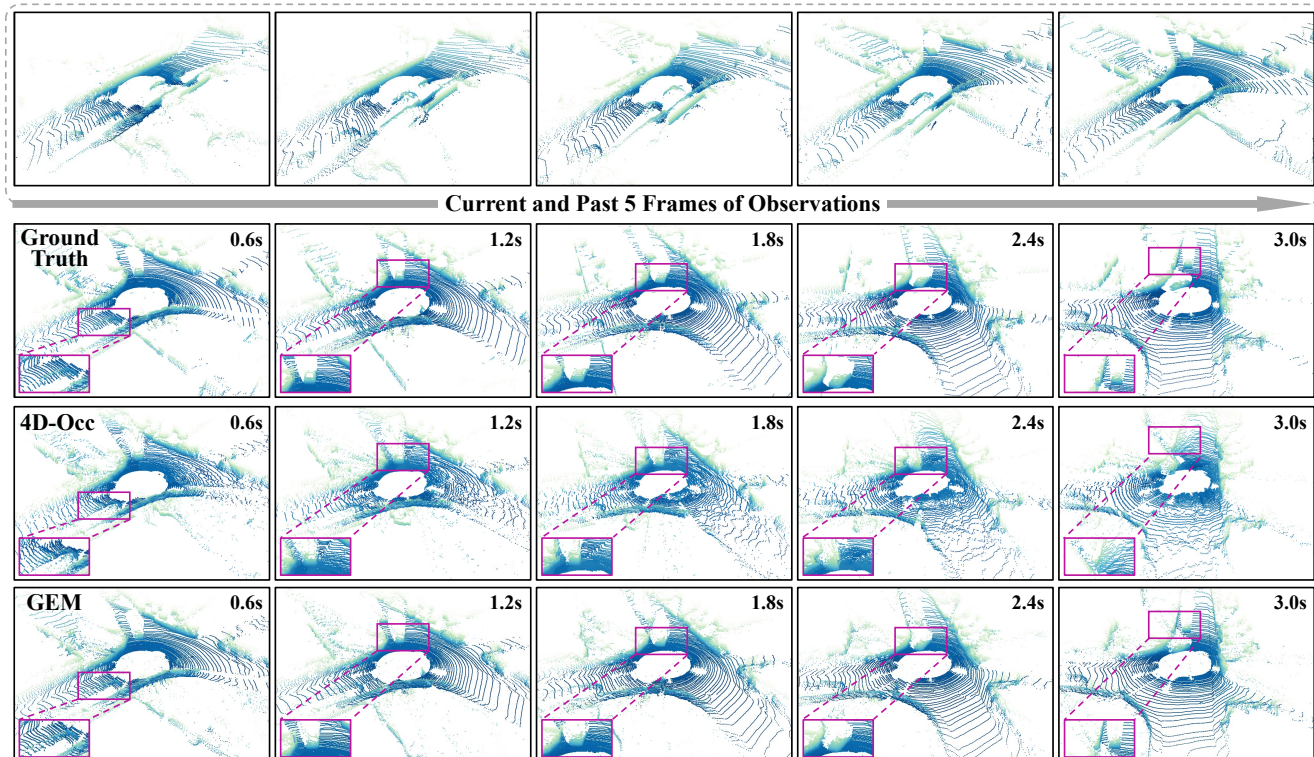


Figure 6. More visual comparisons with competitive LiDAR world models on KITTI Odometry [3].

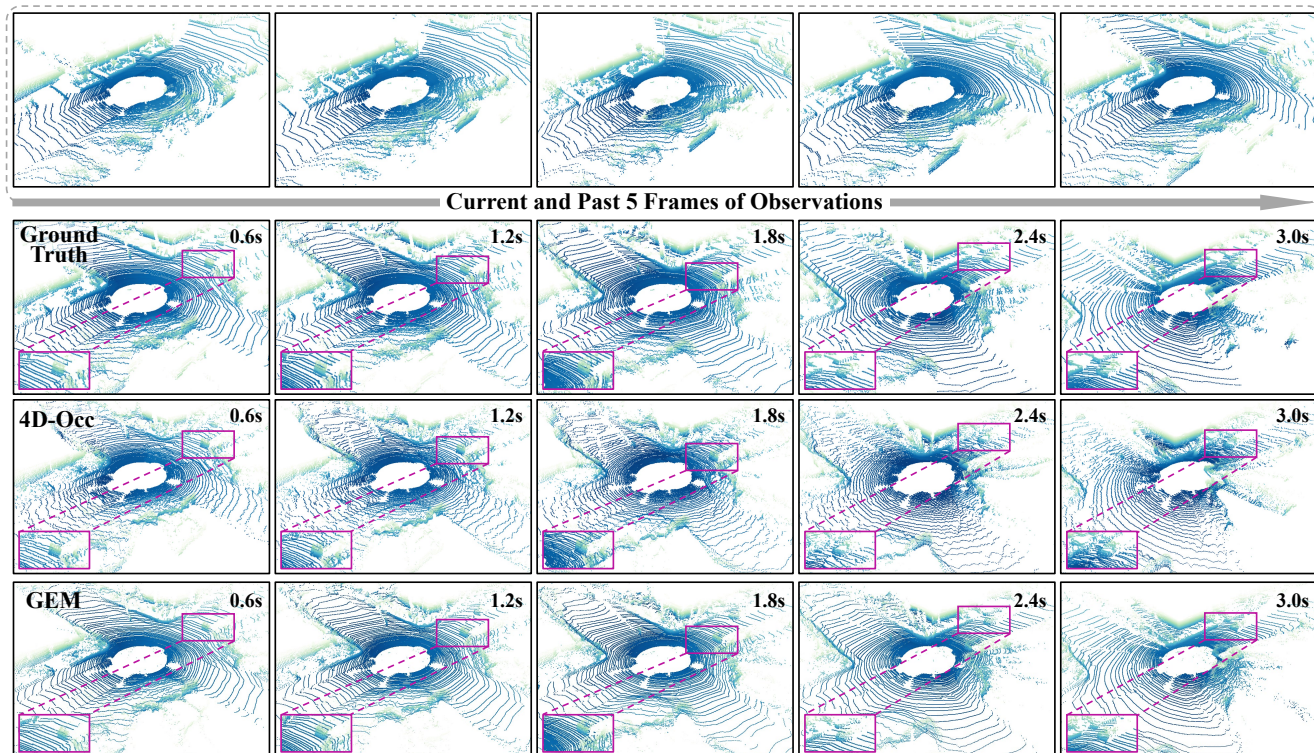


Figure 7. More visual comparisons with competitive LiDAR world models on KITTI Odometry [3].



Figure 8. Visualization of 32-line LiDAR data reconstruction by the LiDAR scene tokenizer..

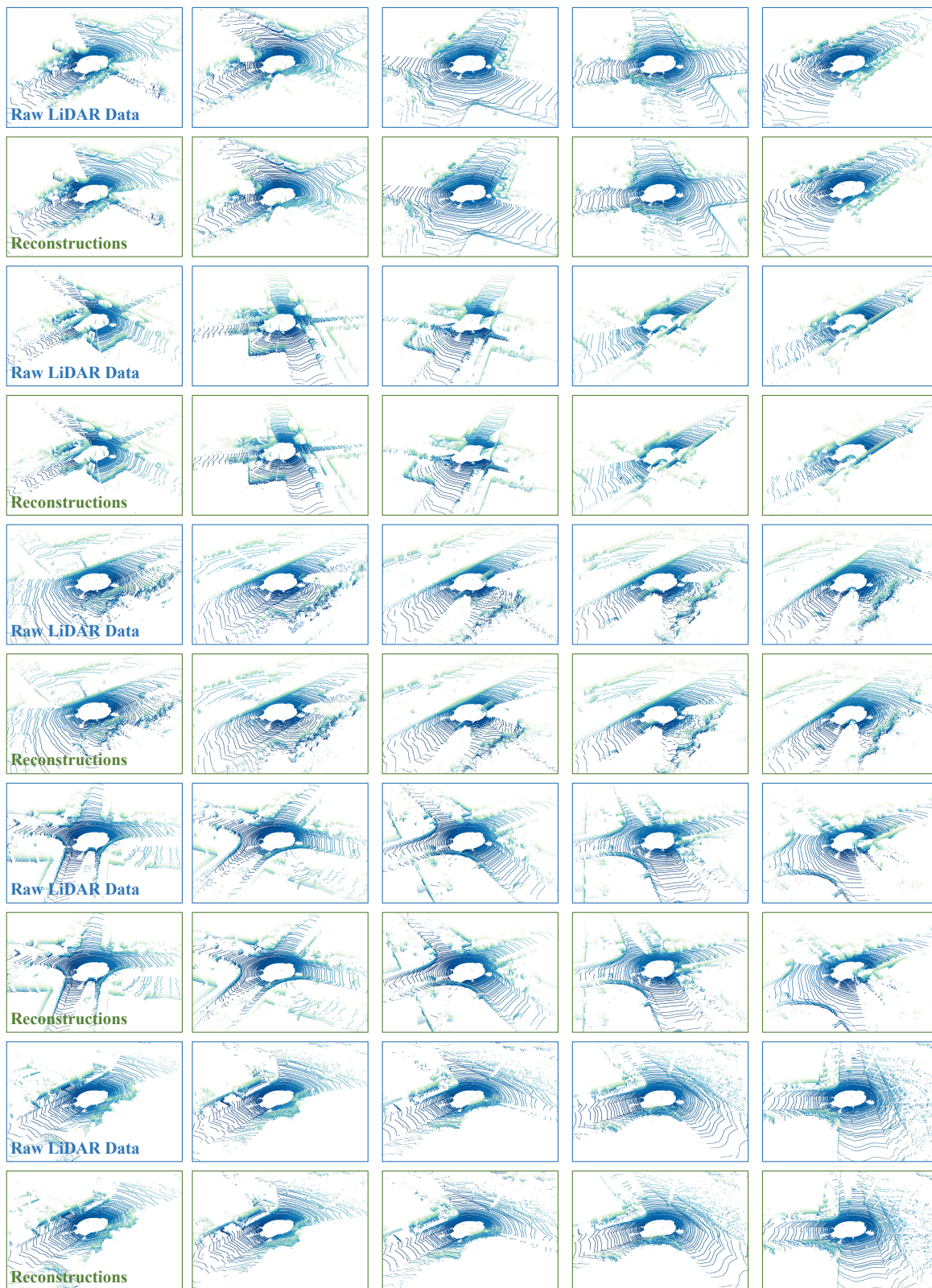


Figure 9. Visualization of 64-line LiDAR data reconstruction by the LiDAR scene tokenizer.

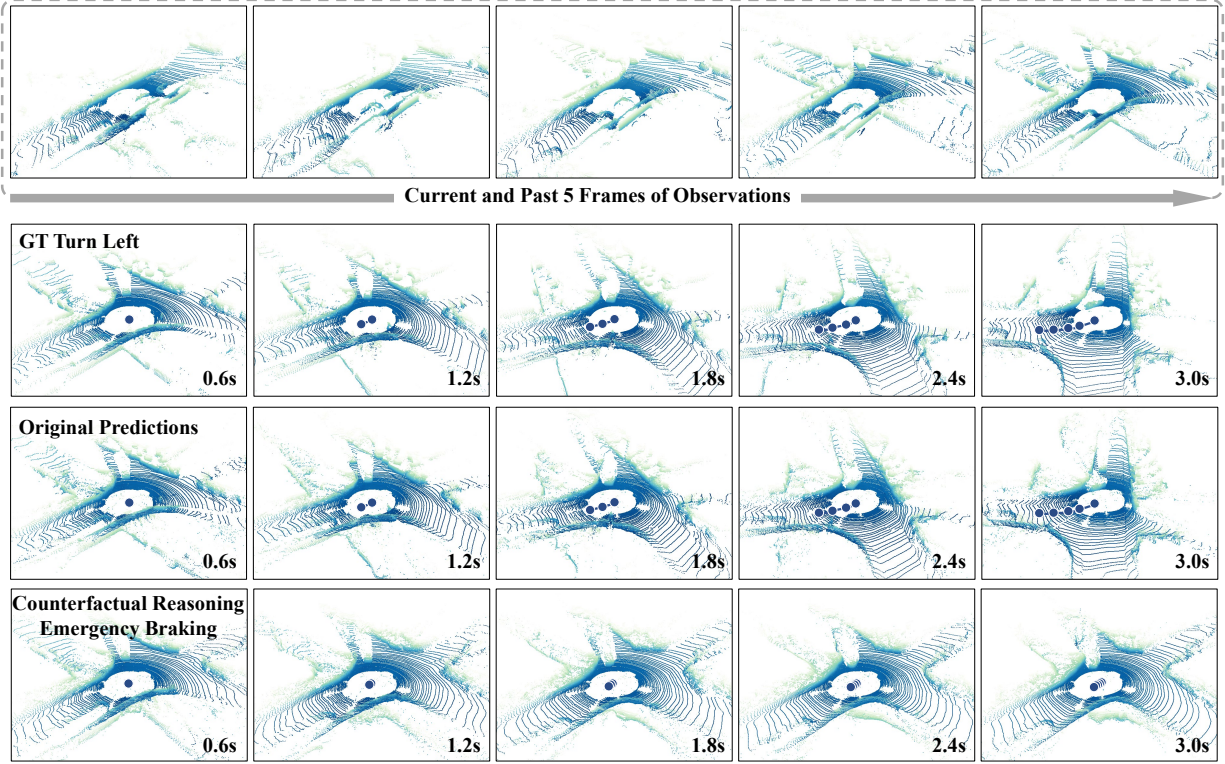


Figure 10. More visualizations of counterfactual reasoning.

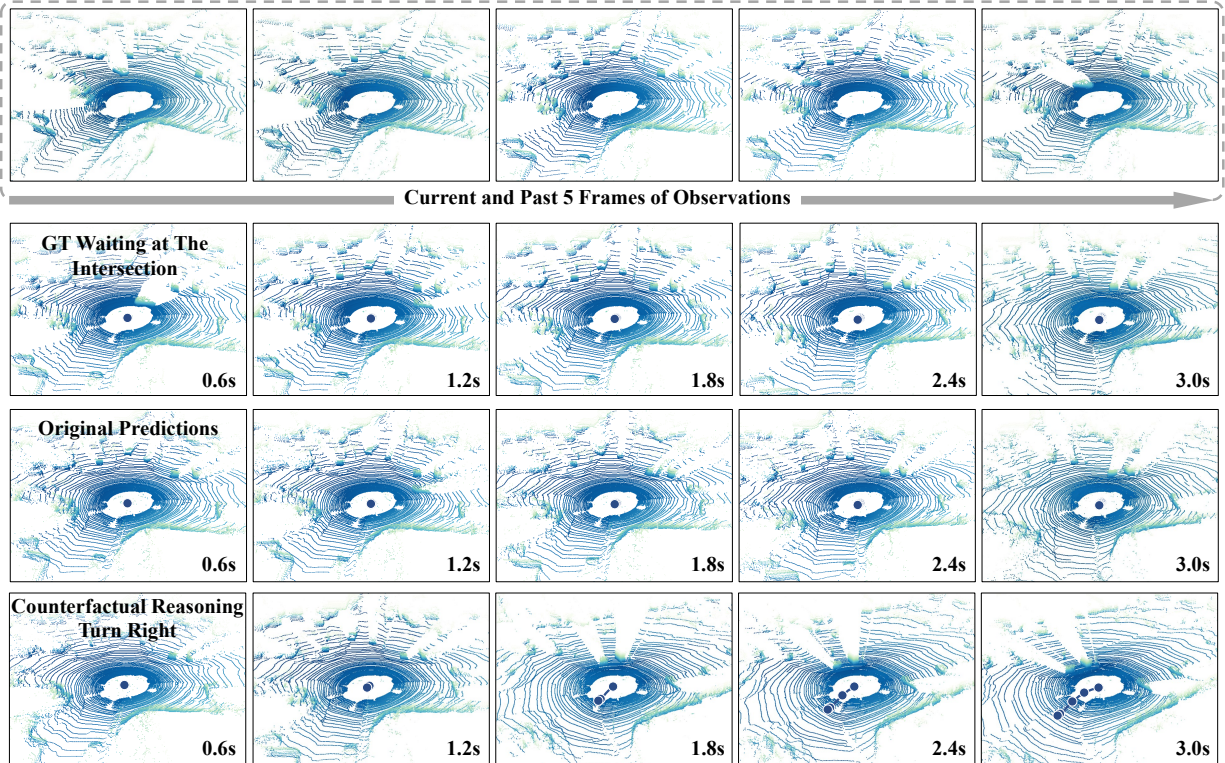


Figure 11. More visualizations of counterfactual reasoning.

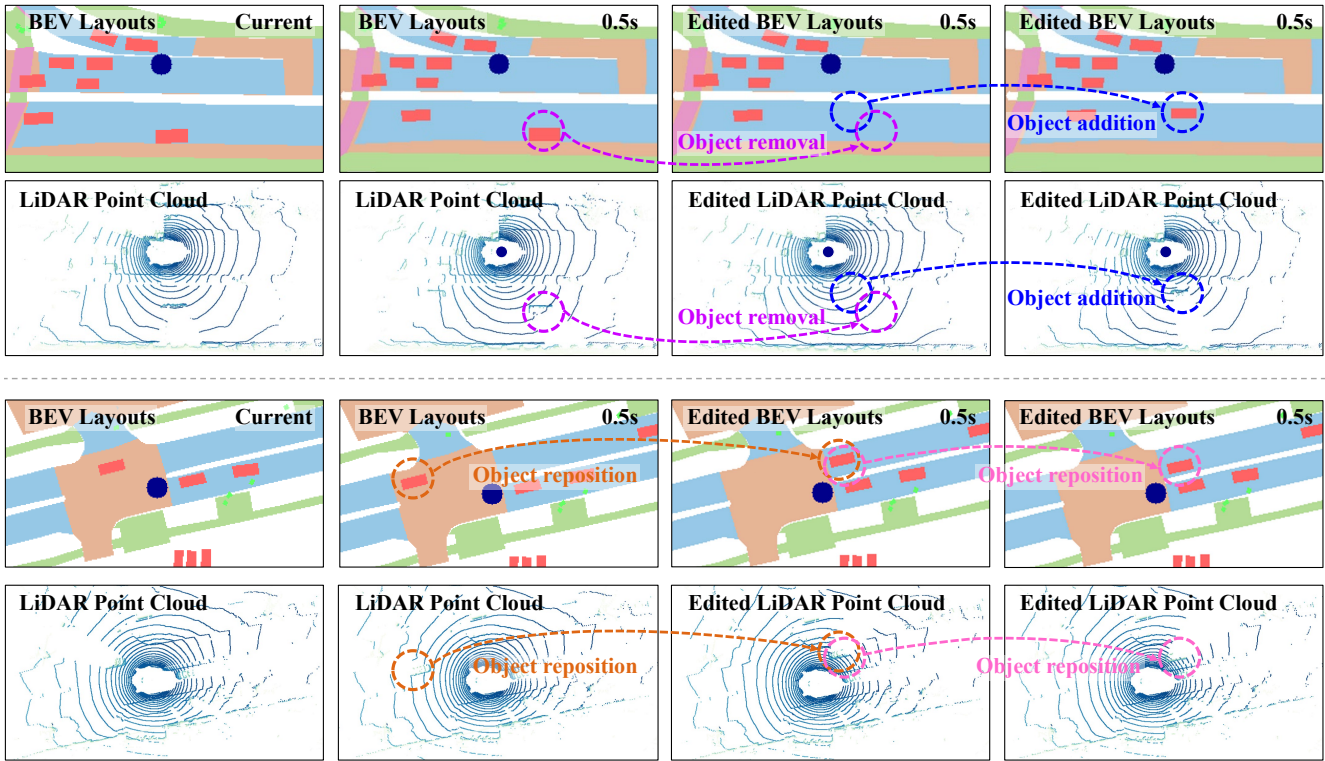


Figure 12. More visualization of bev layouts-controlled generation.

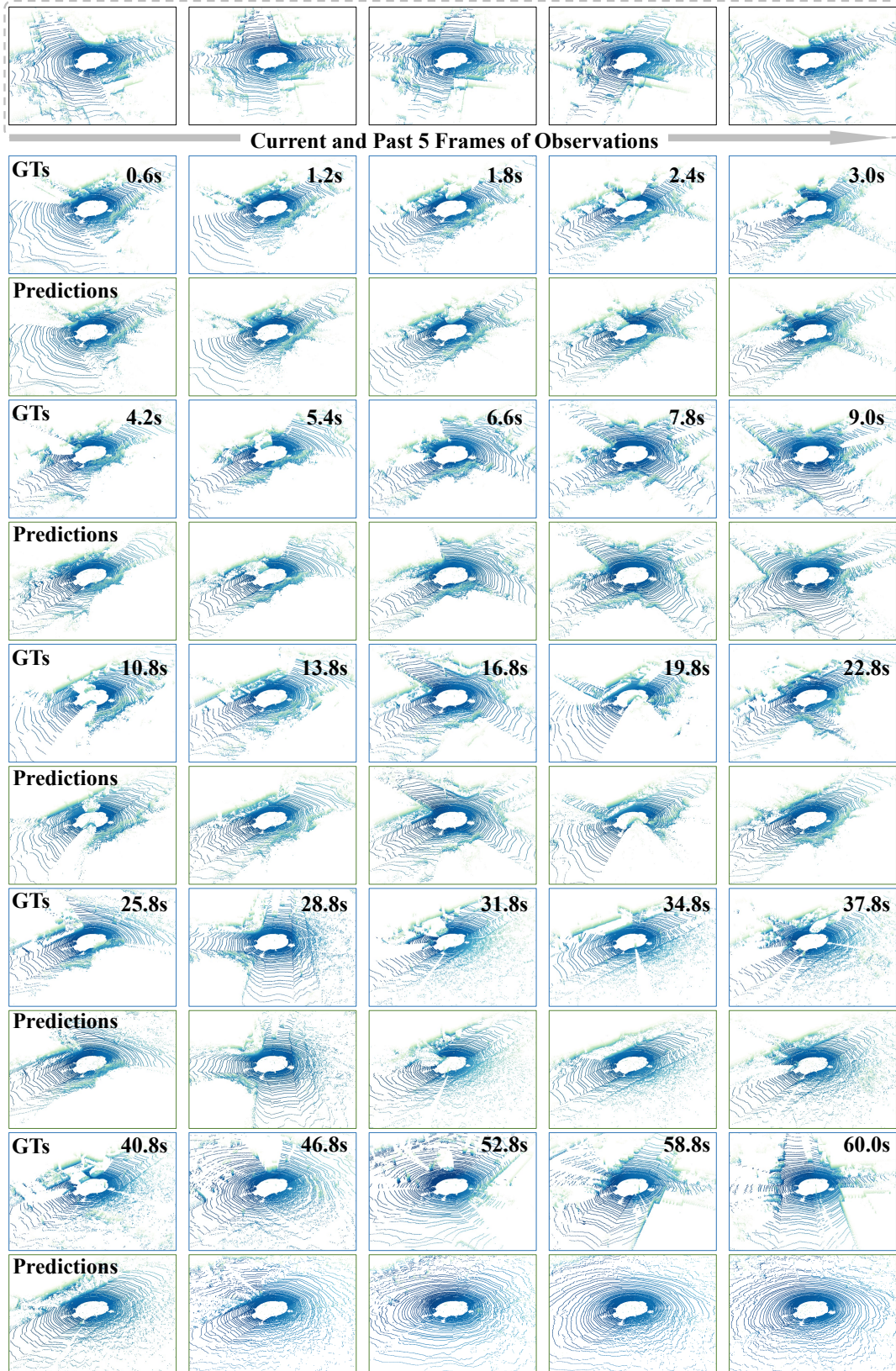


Figure 13. Visualizations of GEM's potential for minute-level future observations prediction.

244

References

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

- [1] 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*, pages 11621–11631, 2020. [2](#), [4](#)
- [2] Yuning Chai, Pei Sun, Jiquan Ngiam, Weiyue Wang, Benjamin Caine, Vijay Vasudevan, Xiao Zhang, and Dragomir Anguelov. To the point: Efficient 3d object detection in the range image with graph convolution kernels. In *CVPR*, pages 16000–16009, 2021. [1](#)
- [3] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, pages 3354–3361. IEEE, 2012. [2](#), [5](#)
- [4] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 27, 2014. [2](#)
- [5] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020. [2](#)
- [6] Tarasha Khurana, Peiyun Hu, David Held, and Deva Ramanan. Point cloud forecasting as a proxy for 4d occupancy forecasting. In *CVPR*, pages 1116–1124, 2023. [2](#)
- [7] Gregory P Meyer, Ankit Laddha, Eric Kee, Carlos Vallespi-Gonzalez, and Carl K Wellington. Lasernet: An efficient probabilistic 3d object detector for autonomous driving. In *CVPR*, pages 12677–12686, 2019. [1](#)
- [8] Andres Milioto, Ignacio Vizzo, Jens Behley, and Cyrill Stachniss. Rangenet++: Fast and accurate lidar semantic segmentation. In *IROS*, pages 4213–4220. IEEE, 2019. [1](#)
- [9] Kazuto Nakashima and Ryo Kurazume. Learning to drop points for lidar scan synthesis. In *IROS*, pages 222–229. IEEE, 2021. [1](#)
- [10] Kazuto Nakashima and Ryo Kurazume. Lidar data synthesis with denoising diffusion probabilistic models. *arXiv preprint arXiv:2309.09256*, 2023. [1](#)
- [11] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, pages 4195–4205, 2023. [2](#)
- [12] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. [2](#)
- [13] Yang Wu, Kaihua Zhang, Jianjun Qian, Jin Xie, and Jian Yang. Text2lidar: Text-guided lidar point cloud generation via equirectangular transformer. *arXiv preprint arXiv:2407.19628*, 2024. [1](#)
- [14] Yang Wu, Yun Zhu, Kaihua Zhang, Jianjun Qian, Jin Xie, and Jian Yang. Weathergen: A unified diverse weather generator for lidar point clouds via spider mamba diffusion. In *CVPR*, pages 17019–17028, 2025. [1](#)
- [15] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *CVPR*, pages 4490–4499, 2018. [1](#)
- [16] Vlas Zyrianov, Xiyue Zhu, and Shenlong Wang. Learning to generate realistic lidar point clouds. In *ECCV*, pages 17–35. Springer, 2022. [1](#)