

# Li-AutoFlow: Autoregressive Flow Matching for Continuous AV Scene Prediction

Anonymous CVPR submission

Paper ID \*\*\*\*

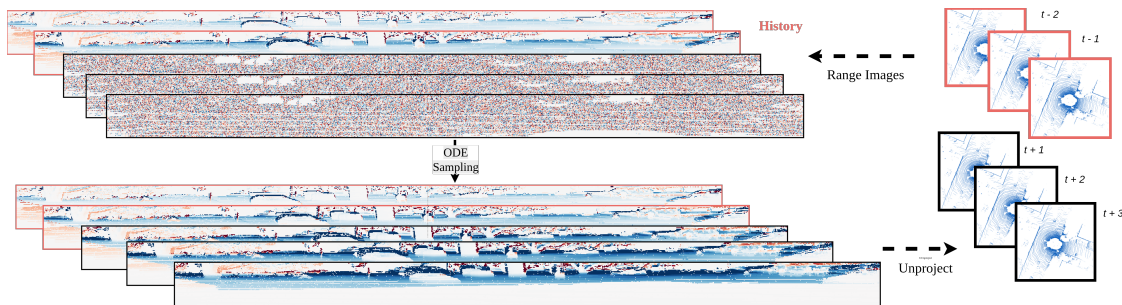


Figure 1. **Li-AutoFlow**: An autoregressive flow based method that predicts future LiDAR scans by conditioning on previous frames, utilizing an intermediate range image representation to ensure efficient generation.

## Abstract

001 Predicting the evolution of complex 3D scenes is fundamental  
 002 to safe autonomous navigation. While LiDAR provides  
 003 the necessary geometric precision, existing forecasting meth-  
 004 ods often struggle with the inherent multimodality of dynamic  
 005 environments. Deterministic models typically collapse to the  
 006 conditional mean, resulting in "ghosting" artifacts and non-  
 007 physical spatial interpolations that fail to represent discrete  
 008 future contingencies.

009 In this paper, we introduce *Li-AutoFlow*, a generative  
 010 framework that unifies Autoregressive Sequence Modeling  
 011 with Flow Matching to produce high-fidelity, multimodal  
 012 LiDAR forecasts. By operating directly on the continuous  
 013 2D range image manifold, our approach bypasses the limita-  
 014 tions of discrete tokenization and lossy codebook bottlenecks  
 015 found in prior generative works. Our formulation leverages  
 016 straight-line probability flow ODEs to transport a noise prior  
 017 into physically plausible future scenes, ensuring temporal  
 018 causality through an autoregressive conditioning scheme.

019 Crucially, our tokenizer-free architecture maintains a  
 020 fully differentiable pipeline from the flow-matching objec-  
 021 tive to the final reconstructed 3D point cloud. This enables  
 022 direct optimization against 3D geometric metrics, such as  
 023 Chamfer distance, ensuring centimeter-level precision. Ex-  
 024 perimental results on KITTI demonstrate that *Li-AutoFlow*  
 025 significantly outperforms state-of-the-art deterministic and

*discrete-diffusion baselines, providing the coherent, disjoint*  
*scene samples necessary for robust motion planning in un-*  
*certain environments.*

## 1. Introduction

Autonomous driving demands the ability to anticipate how  
 the surrounding 3D scene will evolve over the coming sec-  
 onds. Pedestrians change direction, vehicles accelerate or  
 yield, and the observed scene geometry shifts continuously  
 as both the ego vehicle and surrounding agents move. Ac-  
 curate prediction of this evolution is critical: downstream  
 tasks such as motion planning, collision avoidance, and sce-  
 nario simulation all require reliable forecasts of future scene  
 geometry [7, 20].

LiDAR provides a geometrically precise representa-  
 tion of the 3D environment, encoding surface depths with  
 centimeter-level accuracy. Forecasting future LiDAR obser-  
 vations from past sweeps is therefore a key capability for  
 autonomous systems. However, directly forecasting dense,  
 unordered 3D point clouds is computationally prohibitive.  
 To resolve this we formulate the problem in the native 2D  
 range image projection [48]. By mapping sparse 3D coordi-  
 nates into a dense, spherical grid, the range image natively  
 encodes the physical constraints of the sensor while translat-  
 ing the problem into a structured manifold amenable to deep  
 generative architectures. Operating in this space, LiDAR

051	forecasting presents two strict requirements. First, forecasts	104
052	must maintain <i>geometric precision</i> , ensuring that predicted	105
053	points lie on plausible physical surfaces with metric accuracy.	106
054	Second, the model must capture <i>distributional uncertainty</i> ,	107
055	since future scenes are inherently multimodal due to the	108
056	unpredictable behavior of dynamic agents.	109
057	While probabilistic forecasting is standard for object-	110
058	level trajectories, the high dimensionality of dense LiDAR	111
059	forecasting has largely restricted prior work [28, 29] to deter-	112
060	ministic architectures optimized via regression losses. Con-	113
061	sequently, these model regress towards conditional mean of	114
062	the future distribution [36]. In dense 3D space, this expected-	115
063	value optimization results in spatial interpolation between	116
064	divergent futures, creating geometric artifacts: such as non-	117
065	physical smeared surfaces and interpolated ghost objects.	118
066	Modern AV planners require coherent, disjoint scene sam-	119
067	ples to evaluate discrete contingency plans (e.g., a car turning	120
068	left versus going straight). A single, mean-collapsed predic-	121
069	tion fails to provide the discrete hypotheses necessary for	122
070	safe, multimodal contingency planning.	123
071	Generative models address this limitation by representing	124
072	the full distribution of possible futures rather than collaps-	125
073	ing to a single prediction [52]. We build on Flow Match-	126
074	ing [2, 24, 26], a scalable generative framework that learns	127
075	a continuous vector field transporting a simple noise prior	128
076	to the target data distribution. This yields simulation-free	129
077	training and efficient inference via straight-line probability	130
078	flow ODEs, and has demonstrated strong results in image	131
079	synthesis [15], video generation [5, 44], and robot policy	132
080	learning [3].	133
081	Forecasting, however, imposes an additional constraint:	134
082	temporal causality. Each predicted frame must condition	135
083	on the observed history $x_{<t}$ , with predictions fed back au-	136
084	to-regressively into subsequent generations. A single-shot	137
085	conditional model cannot capture this sequential dependence.	138
086	Neither treatment captures the natural asymmetry of a fore-	139
087	casting problem, where past observations are known and	140
088	clean while future frames must be generated stochastically.	141
089	Recent work has demonstrated that autoregressive gener-	142
090	ation can be unified with diffusion-based priors to achieve	143
091	precisely this form of causally consistent sequence model-	144
092	ing [10, 38]. We bring this formulation to LiDAR forecasting	145
093	for the first time: autoregressive conditioning enforces tem-	146
094	poral consistency across the predicted sequence, ensuring	147
095	that each frame is generated in the context of all preceding	148
096	ones, while flow matching preserves the continuous structure	149
097	of depth uncertainty at each step. The varying noise levels	150
098	inherent in this autoregressive diffusion formulation confer	151
099	an additional benefit: the model is trained to generate from	152
100	observations at multiple fidelity levels, and consequently	153
101	develops robustness to the measurement noise present in real	154
102	LiDAR sensors.	155
103	The prior generative approach, Copilot4D [52], dis-	
	cretizes the scene through a VQ-VAE [43] and applies dis-	
	crete diffusion over token indices. This design introduces a	
	compounding three-fold limitation: a lossy codebook bot-	
	tleneck that discards fine structural detail, a combinatorial	
	complexity that ignores the smooth distributional structure of	
	LiDAR data, and a structural decoupling between generation	
	and 3D geometry. We move beyond this discrete bottleneck	
	by operating directly on full-resolution range images without	
	intermediate compression or discretization. Drawing on re-	
	cent evidence that tokenizer-free generation is both feasible	
	and effective [22], our approach preserves the continuous	
	geometric relationships, sensor geometry, surface continuity,	
	and occlusion patterns inherent to the range image manifold.	
	Crucially, by maintaining a continuous pipeline from the	
	flow-matching ODE to the final reconstructed point cloud,	
	our framework is fully differentiable end-to-end. This al-	
	lows us to fine-tune the generative model directly against	
	3D Chamfer distance, a level of geometric supervision that	
	requires non-trivial approximations (e.g., straight-through	
	estimators [43]) in discrete-token approaches, whereas our	
	continuous formulation permits direct end-to-end optimiza-	
	tion.	
	We present <b>Li-AutoFlow</b> , the first continuous, tokenizer-	
	free, autoregressive flow matching model for LiDAR scene	
	forecasting, operating directly in full-resolution range-image	
	space. Our contributions are as follows:	
	• Autoregressive Flow Matching: We introduce an autore-	
	gressive flow matching scheme for multi-step prediction	
	that maintains temporal coherence and expressiveness.	
	• Tokenizer-free Forecasting: We introduce the first	
	tokenizer-free flow matching framework for LiDAR fore-	
	casting, generating full-resolution range images that pre-	
	serve metric depth and sensor geometry without lossy	
	compression.	
	• End-to-End Geometric Supervision: We enable 3D Cham-	
	fer fine-tuning through a fully differentiable reconstruction	
	path, a capability unique to our continuous-space genera-	
	tion.	
	<b>2. Related Work</b>	
	<b>LiDAR point cloud forecasting.</b> Early approaches op-	
	erated directly on unordered point sets using recurrent ar-	
	chitectures [16], or exploited scene flow as a proxy for fu-	
	ture state [27, 46]. Projecting each scan into a range im-	
	age [11, 29, 30] reduces the problem of a structured 2D se-	
	quence prediction task, enabling convolutional and attention-	
	-based backbones at a fraction of the memory cost of volu-	
	-metric representations [21, 53]. Subsequent work explored	
	transformer-based temporal encoders [12, 28, 31], motion-	
	-augmented spatio-temporal convolutions [13], stochastic	
	multi-modal prediction [47], state-space sequence mod-	
	els [40], and visual forecasting as a pre-training objective	
	for camera-based driving [50]. All of these methods are	

156 deterministic regressors, or at best model only a restricted  
157 parametric distribution; they collapse the future to a point  
158 estimate and cannot represent the full distribution over plau-  
159 sible scenes. Copilot4D [52] introduced discrete diffusion  
160 over VQ-VAE tokens as a generative alternative, but the  
161 VQ-VAE bottleneck discards fine metric structure before  
162 any generation occurs, degrading the geometric fidelity of  
163 predicted range images. We avoid this lossy compression  
164 and perform generative modeling directly in the pixel space  
165 of range-difference images, preserving full metric resolution  
166 throughout.

167 **Diffusion and flow matching models.** Denoising diffu-  
168 sion probabilistic models [19] and score-based models [39]  
169 learn to synthesize data by reversing a Gaussian corruption  
170 process, enabling high-fidelity generation through iterative  
171 denoising. Operating directly in pixel space is computa-  
172 tionally expensive, so Latent Diffusion Models (LDM) [37]  
173 compress images with a VAE and perform all denoising in  
174 the resulting low-dimensional latent space; VQVAE-based  
175 tokenizers [35] push this further by quantizing the latent to  
176 a discrete codebook, enabling discrete diffusion [52] and  
177 autoregressive priors over compact token sequences. This  
178 compression pipeline dominates the field [4, 15] as the re-  
179 construction loss of the VAE is small and it enables more  
180 efficient learnig. Flow Matching [1, 24, 25] reformulates  
181 generative modeling as regression onto an optimal transport  
182 problem. Compared to score-based diffusion, this elimi-  
183 nates the need for iterative Langevin corrections and reduces  
184 sampling to a single ODE integration, while providing more  
185 stable training targets and straighter trajectories. Flow match-  
186 ing has since become the backbone of several state-of-the-art  
187 generative model: for images in SD3 [14] and FLUX [15],  
188 and for video in Movie Gen [33], CogVideoX [51], and  
189 Wan [45].

190 However, these methods work on a pretrained latent space.  
191 Just image Transformers (JiT) [22] demonstrated that flow  
192 matching on raw pixels is viable at full resolution by replac-  
193 ing noise prediction with direct clean-image ( $x$ -)prediction.  
194 Under the manifold assumption, predicting the clean data  
195 point rather than the noised quantity allows a plain ViT  
196 operating on large pixel patches to achieve competitive gen-  
197 eration quality without a pretrained tokeniser.

198 **Autoregressive generative models.** Large language mod-  
199 els (LLMs) [6, 42] model joint distributions over discrete  
200 sequences by factorizing them into products of next-token  
201 conditionals, enabling scalable training and open-ended gen-  
202 eration. This paradigm was extended to visual data by rep-  
203 resenting images as sequences of discrete VQ-VAE tokens  
204 and training autoregressive transformers [34, 35] or paral-  
205 lel masked predictors [8]. Visual AutoRegressive Model-  
206 ing (VAR) [41] further reformulated image generation as

coarse-to-fine next-scale prediction over multi-resolution to-  
ken maps, achieving diffusion-level image quality with sub-  
stantially faster inference. These approaches rely on a tok-  
enizer that maps continuous signals into a discrete codebook,  
introducing an inherently lossy compression step. Masked  
Autoregressive (MAR) models [23] mitigate this limitation  
by replacing discrete token prediction with a per-token diffu-  
sion objective, enabling autoregressive generation directly  
in the continuous latent space of a VAE and improving re-  
construction fidelity.

**Autoregressive diffusion models.** Integrating diffusion-  
based generation with autoregressive temporal structure has  
attracted growing interest. Diffusion Forcing [10] is the en-  
abling insight: by assigning an independent, per-token noise  
level at training time, the model learns to condition on any  
mixture of clean context and partially-noised future tokens.  
This subsumes both teacher-forced sequence training and  
standard unconditional diffusion as special cases, and sup-  
ports flexible inference strategies such as sliding-window  
generation and variable-horizon rollout. The Diffusion For-  
cing Transformer (DFoT) [9] instantiated this framework for  
high-resolution video generation with a transformer back-  
bone. Related work on video prediction has also explored  
causal diffusion rollouts [18] and masked autoregressive  
generation with per-token diffusion heads [23], but these  
methods operate on natural-image pixels or VAE latents  
where metric precision is not a constraint. Our work is the  
first to apply autoregressive flow matching to LiDAR range  
image sequences, where the metric semantics of pixel values  
impose qualitatively different fidelity requirements on the  
generative transport process.

## 3. Background

### 3.1. Flow Matching

We formulate our generative model from the perspective  
of flow matching [2, 24, 26]. Consider a data distribution  
 $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$  and a noise distribution  $\epsilon \sim p_{\text{noise}}(\epsilon)$ , where  
 $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . During training, a noisy sample  $\mathbf{z}_t$  is defined  
as an interpolation between data and noise:

$$\mathbf{z}_t = t \mathbf{x} + (1 - t) \epsilon, \quad t \in [0, 1], \quad (1)$$

using a linear schedule [2, 24, 26], so that  $\mathbf{z}_t \sim p_{\text{data}}$  when  
 $t = 1$  and  $\mathbf{z}_t \sim p_{\text{noise}}$  when  $t = 0$ . We sample  $t$  from a  
logit-normal distribution [15],  $\text{logit}(t) \sim \mathcal{N}(\mu, \sigma^2)$ , which  
concentrates training on intermediate noise levels where the  
velocity field is hardest to learn.

The flow velocity  $\mathbf{v}$  is defined as the time-derivative of  
 $\mathbf{z}_t$ :

$$\mathbf{v} = \frac{d\mathbf{z}_t}{dt} = \mathbf{x} - \epsilon. \quad (2)$$

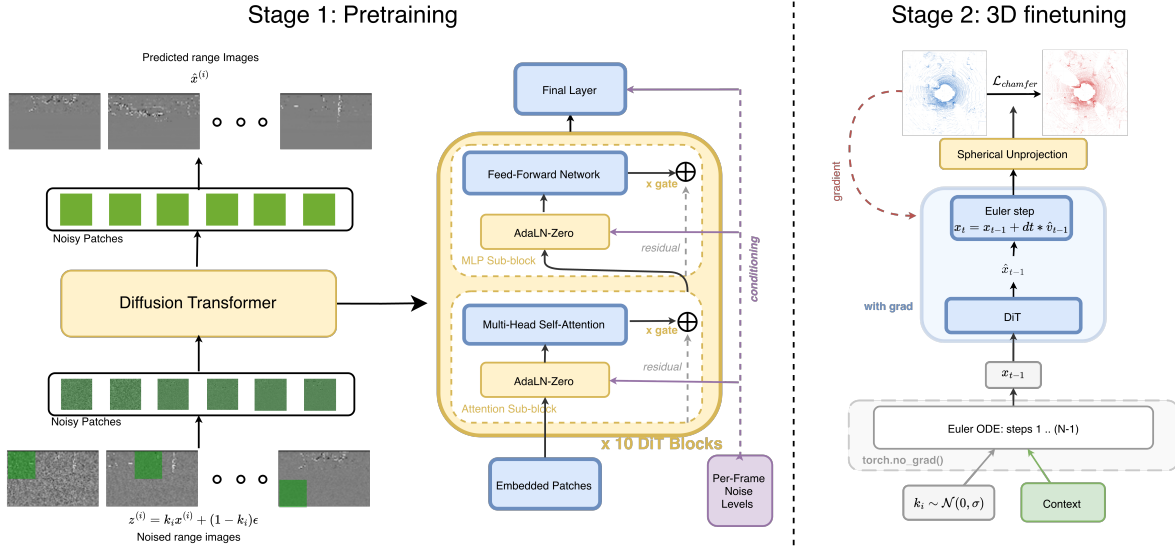


Figure 2. **Overview of the proposed two-stage pipeline for LiDAR point cloud forecasting.** Our architecture (center) utilizes a Diffusion Transformer (DiT) backbone with 10 layers, conditioned on noise levels via AdaLN-Zero. **Stage 1 (Pre-training):** The model learns to denoise range image samples across varying noise schedules (Eq. 8). **Stage 2 (Fine-tuning):** The pipeline incorporates a differentiable Euler sampling transition (right) supervised by a 3D Chamfer loss (Eq. 10), facilitating end-to-end optimization for geometrically consistent 3D reconstruction.

254 A network  $\mathbf{v}_\theta$  is trained to predict this velocity by minimising  
255 ing:

$$256 \quad \mathcal{L}_{\text{FM}} = \mathbb{E}_{t, \mathbf{x}, \epsilon} \left[ \left\| \mathbf{v}_\theta(\mathbf{z}_t, t) - (\mathbf{x} - \epsilon) \right\|_2^2 \right]. \quad (3)$$

257 Sampling is performed by solving the ODE:

$$258 \quad \frac{d\mathbf{z}_t}{dt} = \mathbf{v}_\theta(\mathbf{z}_t, t), \quad (4)$$

259 starting from  $\mathbf{z}_0 \sim p_{\text{noise}}$  and integrating to  $t = 1$

### 260 3.2. Autoregressive Diffusion

261 Standard diffusion models apply a uniform noise level across  
262 all tokens in a sequence. Diffusion Forcing [10] relaxes this  
263 by assigning *independent* noise levels  $k_{1:T} \in [0, 1]^T$  to  
264 individual tokens. Each token is noised independently via  
265 Eq. (1):

$$266 \quad \mathbf{z}_t^{(i)} = k_i \mathbf{x}^{(i)} + (1 - k_i) \epsilon^{(i)}, \quad \epsilon^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (5)$$

267 and a single network is trained to denoise any combination  
268 of noise levels simultaneously:

$$269 \quad \mathcal{L}_{\text{DF}} = \mathbb{E}_{k_{1:T}, \mathbf{x}_{1:T}, \epsilon_{1:T}} \left[ \frac{1}{T} \sum_{i=1}^T \left\| \mathbf{v}_\theta(\mathbf{z}_{1:T}, k_{1:T})^{(i)} - (\mathbf{x}^{(i)} - \epsilon^{(i)}) \right\|_2^2 \right]. \quad (6)$$

270 This enables flexible conditioning at inference: any subset  
271 of tokens can be pinned at  $k_i = 1$  (fully clean) while others  
272 are initialized at  $k_i = 0$  and refined via Eq. (4), with

no architectural difference between training and inference. 273  
The Diffusion Forcing Transformer (DFoT) [9] scales this 274  
to high-resolution sequences via a transformer backbone, 275  
injecting per-token noise levels through AdaLN-Zero modulation 276  
[32]. 277

## 278 4. Method

279 We present **Li-AutoFlow** (Fig 2, an autoregressive generative 280  
model for LiDAR scene forecasting that operates directly on 281  
full-resolution range difference images. Given a reference 282  
range image  $\mathbf{r}_0$  and  $T_c$  past range differences  $\Delta \mathbf{r}_{1:T_c}$ , the 283  
model predicts  $K$  future differences  $\Delta \mathbf{r}_{T_c+1:T_c+K}$  using a 284  
flow-matching generative process. Future range images are 285  
then recovered through a differentiable reconstruction step.

286 Our model is implemented using a Diffusion Forcing 287  
Transformer (DFoT) that jointly models spatial geometry 288  
within each frame and temporal dynamics across frames. 289  
Unlike prior approaches based on latent diffusion [52], we 290  
operate directly in the native sensor resolution, preserving 291  
the one-to-one correspondence between pixels and LiDAR 292  
measurements. Training proceeds in two stages: we first 293  
pretrain the model on clean next-frame predictions using the 294  
flow-matching objective, then fine-tune on a 3D-reprojection 295  
loss to improve geometric fidelity in the output point cloud.

### 296 4.1. Preprocessing

297 Each LiDAR sweep  $\mathcal{P}_k \subset \mathbb{R}^3$  is projected onto a spherical 298  
range image  $\mathbf{r}_k \in \mathbb{R}^{H \times W}$ , where pixel  $(i, j)$  stores the

299 Euclidean range  $r_{ij} = \|\mathbf{p}_{ij}\|_2$ . Pixels with no LiDAR  
300 return are marked invalid and assigned  $r_{ij} = 0$ . To align  
301 frames, we warp  $\mathbf{r}_k$  into the coordinate frame of the ref-  
302 erence image  $\mathbf{r}_0$  using the relative pose  $\mathbf{T}_{0\leftarrow k}$ , yielding  
303  $\mathbf{r}_k^{(0)} = \text{Warp}(\mathbf{r}_k, \mathbf{T}_{0\leftarrow k})$ . Rather than modeling raw range  
304 images directly, we operate on the *range difference*

$$305 \quad x_t = \Delta \mathbf{r}_k = \mathbf{r}_k^{(0)} - \mathbf{r}_0, \quad k = 1, \dots, N, \quad (7)$$

306 which removes static background and concentrates signal  
307 on dynamic objects. The resulting representation is sparse  
308 and near-zero for the majority of pixels, substantially reduc-  
309 ing the complexity of the generative manifold compared to  
310 modeling raw range images directly.

## 311 4.2. Full-Resolution Autoregressive Flow Matching

312 Each pixel of a range image corresponds to a physically  
313 grounded depth measurement: its value determines the 3D  
314 location of a LiDAR return through the spherical projection.  
315 Spatial compression disrupts this one-to-one correspondence.  
316 When convolutional VAEs downsample the range image,  
317 neighbouring laser beams are merged into shared latent cells,  
318 erasing fine angular structures such as sharp object bound-  
319 aries, thin objects, and precise range discontinuities. To  
320 preserve this geometric fidelity, we perform generative mod-  
321 elling directly in the full-resolution range difference space  
322  $\mathbf{x}_k \in \mathbb{R}^{H \times W}$ .

323 We learn this distribution using flow matching. For a  
324 ground-truth range difference frame  $\mathbf{x}_k$ , we construct noisy  
325 interpolants

$$326 \quad \mathbf{z}_t = t\mathbf{x}_k + (1-t)\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

327 and train a denoising network  $D_\theta$  to predict the clean  
328 signal  $\hat{\mathbf{x}}^{(i)}$ . This  $x$ -prediction parameterization provides  
329 a well-conditioned regression target across the noise spec-  
330 trum, avoiding the large target variance encountered by noise-  
331 or velocity-prediction objectives in high-dimensional pixel  
332 spaces [22]. The loss however, is evaluated in velocity space,  
333 which proves to be a more favorable learnable objective [22].  
334 Training then minimizes the velocity-space regression error

$$335 \quad \mathcal{L} = \mathbb{E} \left[ \frac{1}{K} \sum_{i=T_c+1}^{T_c+K} \left\| \hat{\mathbf{v}}^{(i)} - \mathbf{v}^{(i)} \right\|_2^2 \right], \quad (8)$$

336 where the loss is applied only to future frames  $i > T_c$ .

337 Operating at full resolution remains tractable due to the  
338 properties of the range difference representation. Subtract-  
339 ing a reference frame suppresses the static background, pro-  
340 ducing sparse signals that concentrate on dynamic scene  
341 elements. The resulting generative manifold is therefore  
342 substantially simpler than that of raw range images, enabling  
343 stable training without a latent compression bottleneck.

**Pacification.** To process frames with a transformer back-  
bone, each  $\mathbf{x}_k \in \mathbb{R}^{H \times W}$  is partitioned into non-overlapping  
 $p \times p$  patches and linearly projected into  $C$ -dimensional  
tokens, yielding  $\frac{H}{p} \times \frac{W}{p}$  tokens per frame. This patch projec-  
tion reorganizes the computation for transformer processing  
but does not alter the underlying signal resolution, and thus  
preserves the full geometric fidelity of the range image.

**Backbone.** We instantiate the denoiser  $D_\theta$  using a Diffu-  
sion Forcing Transformer (DFoT) [9]. The model processes  
frame-token sequences through alternating spatial and tem-  
poral attention layers. Spatial attention operates within each  
frame, capturing local geometric structure and the relation-  
ships between neighbouring laser beams. Temporal attention  
propagates information across frames, enabling the model  
to track the motion of surfaces over time. Causality is en-  
forced by masking future frames during temporal attention  
so that each generated frame conditions only on its preced-  
ing history. This factorized design scales more efficiently  
with sequence length than full spatiotemporal attention while  
explicitly separating spatial coherence from temporal dynam-  
ics.

**Per-frame noise conditioning.** Diffusion Forcing assigns  
independent noise levels to different frames within the same  
forward pass. Context frames are fixed at  $k_i = 1$  (clean),  
while future frames are sampled with  $k_i \in [0, 1)$  from the  
training noise schedule. Each frame’s noise level is encoded  
using a sinusoidal embedding and injected into every atten-  
tion block via AdaLN-Zero modulation [32]. This condition-  
ing provides the network explicit awareness of each frame’s  
noise state, allowing it to exploit clean context frames while  
denoising noisy future frames during autoregressive genera-  
tion.

## 510 4.3. 3D Geometric Fine-Tuning

511 The flow matching stage optimises a purely 2D surrogate  
512 objective: the network receives no direct signal about the  
513 3D geometry its predictions imply. Since  $\hat{\mathbf{x}}_k$  lives in full-  
514 resolution range image space, differentiable reconstruction  
515 to  $\hat{\mathcal{P}}_k$  is exact and requires no decoder—we exploit this  
516 to introduce a fine-tuning stage that closes the supervision  
517 gap by minimising Chamfer distance between predicted and  
518 ground-truth point clouds.

519 Given  $\hat{\mathbf{x}}_k$ , we recover  $\hat{\mathcal{P}}_k$  by applying Eq. (7) to obtain  
520  $\hat{\mathbf{r}}_k^{(0)}$ , then unprojecting and applying  $\hat{\mathbf{T}}_{k\leftarrow 0}$ . Now, backprop-  
521 agating through all  $S$  ODE steps is memory-prohibitive. We  
522 instead run  $S - 1$  steps under `torch.no_grad()`, detach  
523 the penultimate state  $\bar{\mathbf{z}} = \text{sg}(\mathbf{z}_{t_{S-1}})$ , and perform a single  
524 gradient-enabled final step:

$$525 \quad \hat{\mathbf{x}}_k = \bar{\mathbf{z}} + \Delta t \cdot \mathbf{v}_\theta(\bar{\mathbf{z}}, t_{S-1}, \mathbf{c}), \quad (9)$$

so that gradients from  $\mathcal{L}_{3D}$  flow into  $\mathbf{v}_\theta$  through exactly one denoising step, keeping memory overhead constant in  $S$ .

The fine-tuning objective combines the Chamfer distance with an invalid-pixel regularizer:

$$\mathcal{L}_{3D} = \text{CD}(\hat{\mathcal{P}}_k, \mathcal{P}_k) + \lambda_{\text{inv}} \mathcal{L}_{\text{inv}}, \quad (10)$$

where the Chamfer distance is:

$$\begin{aligned} \text{CD}(\hat{\mathcal{P}}_k, \mathcal{P}_k) &= \frac{1}{|\hat{\mathcal{P}}_k|} \sum_{\hat{\mathbf{p}} \in \hat{\mathcal{P}}_k} \min_{\mathbf{p} \in \mathcal{P}_k} \|\hat{\mathbf{p}} - \mathbf{p}\|_2 \\ &+ \frac{1}{|\mathcal{P}_k|} \sum_{\mathbf{p} \in \mathcal{P}_k} \min_{\hat{\mathbf{p}} \in \hat{\mathcal{P}}_k} \|\mathbf{p} - \hat{\mathbf{p}}\|_2, \end{aligned} \quad (11)$$

and  $\mathcal{L}_{\text{inv}}$  penalises non-zero predictions at pixels where the ground-truth has no valid return:

$$\mathcal{L}_{\text{inv}} = \frac{1}{|\mathcal{I}|} \sum_{(i,j) \in \mathcal{I}} H_\delta(\hat{r}_{ij}), \quad \mathcal{I} = \{(i,j) : r_{ij} = 0\}, \quad (12)$$

where  $H_\delta$  is the Huber loss,

$$H_\delta(x) = \begin{cases} \frac{1}{2}x^2, & |x| \leq \delta, \\ \delta(|x| - \frac{1}{2}\delta), & |x| > \delta, \end{cases} \quad (13)$$

suppressing hallucinated returns in occluded and empty regions that Chamfer distance alone does not penalise. Unlike a squared penalty, the linear tail of  $H_\delta$  prevents the loss from being dominated by the few pixels with large spurious predictions, while still providing a smooth  $L_2$  signal near zero.

#### 4.4. Implementation Details

**Architecture.** The DFoT backbone uses hidden dimension 384, patch size  $8 \times 8$ , depth 8, and 6 attention heads with MLP ratio 4. Temporal position embeddings are learned 1D embeddings along the sequence axis.

**Datasets and splits.** We evaluate on KITTI Odometry [17], using sequences 00–05 for training, 06–07 for validation, and 08–10 for testing, and on nuScenes [7]. Range images are at resolution  $64 \times 2048$ . We use  $T_c = 2$  context frames and  $K = 3$  future frames.

**Training.** The flow model is optimised with AdamW ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , weight decay 0) at learning rate  $10^{-3}$  with gradient norm clipping at 1.0. The flow model is trained for 300k steps with batch size 8 on L40s GPUs. The timestep sampling follows a logit norm based noise schedule [22],  $k_i \sim \mathcal{N}(\mu_p, \sigma_p^2)$ , which biases training toward intermediate noise levels. 3D fine-tuning runs for up to 10k steps at learning rate  $10^{-4}$  with  $S = 5$  ODE steps with 5 steps of ODE iterations and  $\lambda_{\text{inv}} = 0.1$ .

**Autoregressive inference.** During inference, the model generates  $K$  future frames jointly within a single flow trajectory. The context frames are fixed at  $k_i = 1$  and future frames are initialised from standard Gaussian noise at  $k_i = 0$ . The probability flow ODE of Eq. (4) is then integrated from  $t = 0$  to  $t = 1$  using an ODE solver [num of samples = 5], producing the predicted range differences  $\hat{\mathbf{x}}_{T_c+1:T_c+K}$ .

## 5. Experiments

### 5.1. Evaluation

**Dataset.** We evaluate on the KITTI Odometry benchmark [17], a standard testbed for LiDAR-based autonomous driving research. The dataset provides 22 sequences of outdoor driving data captured with a Velodyne HDL-64E at 10 Hz. We follow the standard split used in prior forecasting work [28, 29, 31]: sequences 00–07 for training and 08–10 for testing.

**Metrics.** We evaluate performance using two complementary metrics over the  $K = 3$  predicted future steps. *Chamfer Distance (CD)* measures the mean bidirectional nearest-neighbour distance between predicted and ground-truth 3D point clouds (Eq. 11). We report both per-step values and the mean across the prediction horizon; lower values indicate better geometric accuracy. *Range Image Loss (LI)* computes the mean absolute error between predicted and ground-truth range images, Lower values indicate improved reconstruction fidelity.

### 5.2. Baselines

We compare against three range-image-based point cloud forecasting methods that span the major architectural families. **Mersch et al. [29]** project each LiDAR scan to a range image and stack consecutive frames into a 3D tensor, processing it with an encoder-decoder of 3D spatio-temporal convolutions in a fully self-supervised manner, establishing the foundational range-image baseline.

**PCPNet [28]** replaces convolutions with a Transformer that operates on range image tokens compressed along both spatial dimensions, and augments training with a semantic consistency loss derived from SemanticKITTI labels to improve downstream utility.

**ATPPNet [31]** processes range image sequences through Conv-LSTM blocks gated by dual channel-wise and spatial attention, complemented by a parallel 3D-CNN branch for global spatio-temporal context, yielding the strongest recurrent baseline.

**CoPilot4D [52]** is a generative world model that first tokenizes point clouds into discrete BEV tokens via a VQ-VAE, then forecasts future tokens using a masked discrete diffusion Transformer, enabling scalable unsupervised training

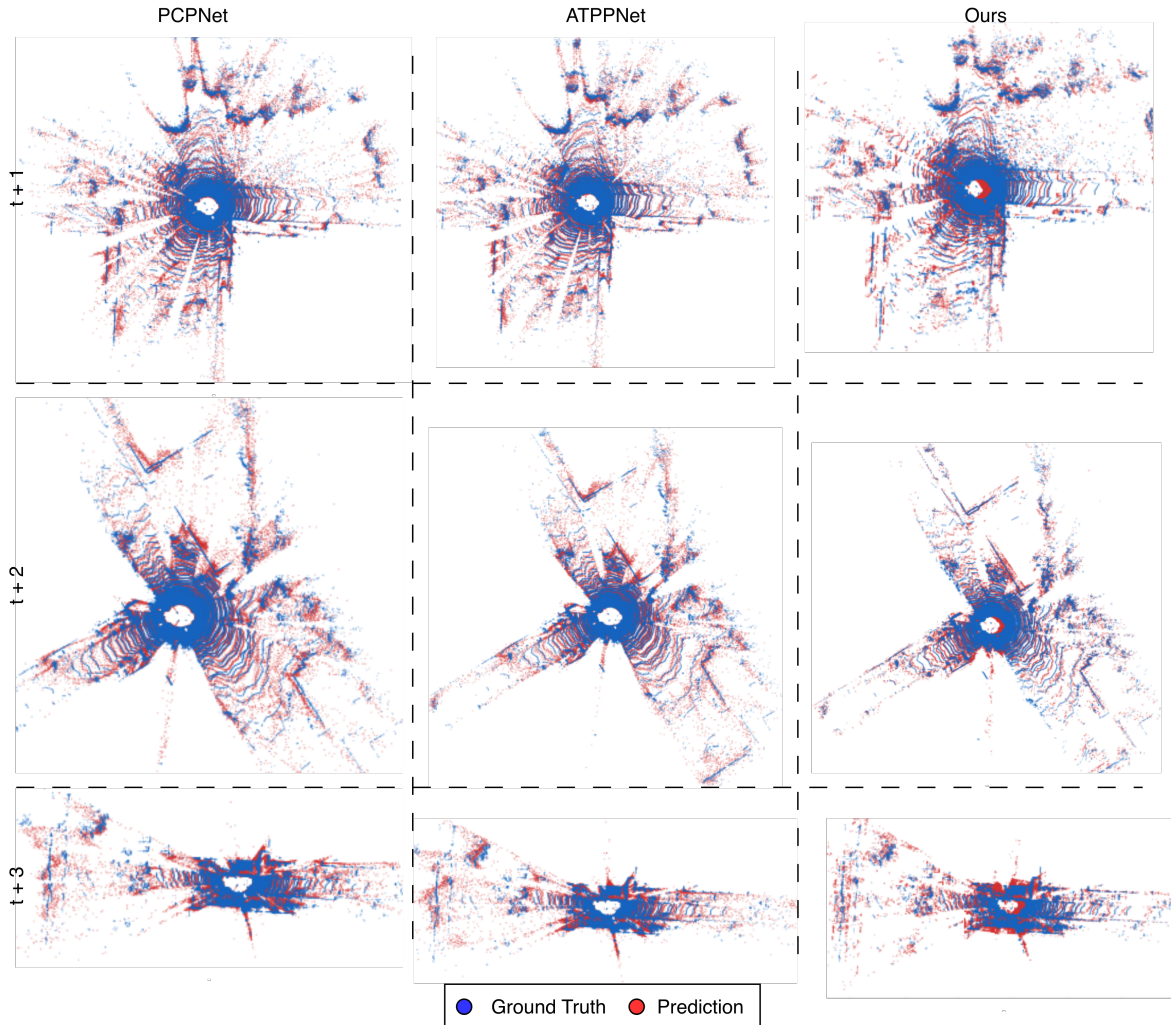


Figure 3. Qualitative comparison of predicted (red) versus ground truth (blue) point clouds at horizons  $t+1$ ,  $t+2$ , and  $t+3$  (rows) for PCPNet [28], ATPNet [31], and our method (columns, left to right).

478 across multiple datasets. \*As the code is not publicly avail- 493  
 479 able, we report the mean Chamfer distance for the 1-second 494  
 480 prediction horizon directly from their paper. 495

### 481 5.3. Quantitative Results 496

482 Table 1 compares our method with prior point cloud fore- 497  
 483 casting approaches on KITTI Odometry. Our approach con- 498  
 484 sistentlly outperforms all baselines across all future hori- 499  
 485 zons. At  $t+1$ , our model achieves a CD of 0.087, improv- 500  
 486 ing over PCPNet (0.280) and ATPNet (0.221) by **68.9%** 501  
 487 and **60.6%**, respectively. At  $t+2$ , the gap further widens, 502  
 488 where our model obtains 0.117, corresponding to **65.6%** and 503  
 489 **57.3%** improvements over PCPNet (0.340) and ATPNet 504  
 490 (0.274). At  $t+3$ , our method achieves 0.155, reducing the 505  
 491 error by **62.4%** and **54.9%** compared to PCPNet (0.412) 506  
 492 and ATPNet (0.344). Against CoPilot4D, a strong gen- 507

erative baseline based on discrete diffusion, our method 493  
 achieves a mean CD of 0.120 compared to their reported 494  
 mean of 0.180, a **33.3%** improvement, despite CoPilot4D 495  
 operating on tokenized BEV representations trained across 496  
 multiple large-scale datasets. Overall, our approach substan- 497  
 tially outperforms all prior methods while maintaining stable 498  
 performance as the prediction horizon increases. 499

### 500 5.4. Qualitative Analysis 501

Figure 3 shows predicted range images and their correspond- 501  
 ing 3D back-projections for representative test sequences. 502  
 Li-AutoFlow produces sharper object boundaries and more 503  
 coherent predictions of dynamic objects (e.g., vehicles and 504  
 cyclists) compared to baseline methods, which tend to blur 505  
 high-motion regions due to the averaging effects of determin- 506  
 istic regression. Our method preserves fine local geometric 507

Table 1. Comparison of point cloud forecasting performance on KITTI Odometry (sequences 08–10). † denotes our method

Method	Chamfer Distance (m) ↓			
	$t+1$	$t+2$	$t+3$	Mean
PCPNet [28]	0.280	0.340	0.412	0.344
ATPPNet [31]	0.221	0.274	0.344	0.280
Copilot4D* [52]	-	-	-	0.180
Li-AutoFlow (ours)†	<b>0.087</b>	<b>0.117</b>	<b>0.155</b>	<b>0.120</b>
w/o 3D fine-tuning	0.177	0.234	0.271	0.234

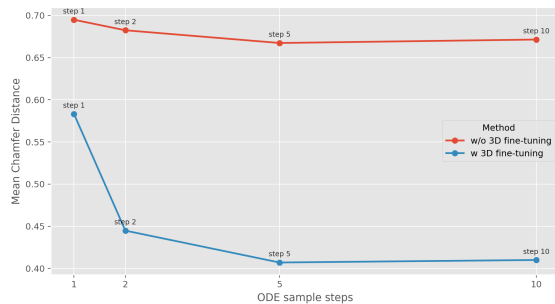


Figure 4. Figure show the Chamfer distance with varying number of ODE sampling steps. Results are computed on a single test sequence.

508 changes and demonstrates improved understanding of scene  
 509 dynamics, resulting in more accurate motion structures in  
 510 the predicted frames. Furthermore, the stochastic nature of  
 511 the flow model enables the generation of geometrically plau-  
 512 sible future states, particularly around occluded or uncertain  
 513 regions.

## 514 5.5. Ablation Studies

515 **Effect of chamfer distance fine-tuning** Table 1 shows  
 516 fine-tuning with a 3D Chamfer distance objective yields sub-  
 517 stantial improvements across all prediction horizons, reduc-  
 518 ing mean CD from 0.227 to 0.120, a relative improvement  
 519 of over 47%. Qualitatively, the fine-tuned model produces  
 520 significantly fewer outlier points, with predictions adher-  
 521 ing much more tightly to planar surfaces such as roads and  
 522 building facades. Without fine-tuning, stray points are scat-  
 523 tered around object boundaries and flat regions, whereas  
 524 fine-tuning encourages the predicted geometry to collapse  
 525 onto the correct surfaces. The gains are consistent across  
 526 time steps, with CD improving from 0.177 to 0.087 at  $t+1$   
 527 and from 0.271 to 0.155 at  $t+3$ , suggesting that 3D supervi-  
 528 sion corrects systematic geometric errors rather than simply  
 529 reducing noise at a single horizon.

530 **Number of ODE sampling steps.** In Figure 4 we evaluate  
 531 generation quality as a function of the number of ODE inte-

gration steps (1, 2, 5, 10) with and without 3D fine-tuning, 532  
 reporting Mean Chamfer Distance in Figure 4. Without 3D 533  
 fine-tuning, performance improves only marginally as the 534  
 number of steps increases, plateauing beyond step 5 with 535  
 consistently high CD throughout. With 3D fine-tuning, qual- 536  
 ity improves sharply between steps 1 and 5, after which gains 537  
 become negligible. This suggests that 5 ODE steps strike a 538  
 favorable balance between inference speed and prediction 539  
 accuracy, and that 3D fine-tuning is the dominant factor in 540  
 performance improvement rather than increasing the number 541  
 of sampling steps. 542

## 543 6. Conclusion

In this work, we introduced Li-AutoFlow, a LiDAR scene 544  
 forecasting framework designed to bypass the quantization 545  
 artifacts and information loss typical of discrete latent-space 546  
 architectures. By formulating the prediction task as a flow- 547  
 matching process directly on the full-resolution range-image 548  
 manifold, we maintain the inherent sensor topology and 549  
 metric depth necessary for high-fidelity 3D reconstruction. 550  
 Our experiments demonstrate that a tokenizer-free, Diffu- 551  
 sion Forcing Transformer can effectively capture complex 552  
 scene dynamics, while the fully differentiable nature of our 553  
 pipeline allows for direct optimization via a 3D Chamfer 554  
 distance (CD) objective. This dense geometric supervision 555  
 results in a 47% improvement in CD over baseline models 556  
 that lack 3D-aware fine-tuning, highlighting the importance 557  
 of alignment between the generative space and the physical 558  
 3D world. 559

**Limitations and future work.** The current system relies 560  
 on ground-truth odometry during training (KITTI poses) and 561  
 a pose estimator at test time, introducing an external depen- 562  
 dency that should ideally be removed by jointly learning 563  
 ego-motion estimation directly from raw scans. Future work 564  
 will transition to a fully end-to-end architecture via joint 565  
 odometry learning directly from raw scans, while adopting 566  
 consistency distillation to collapse the sampling process into 567  
 1–2 steps without sacrificing structural fidelity. Finally, ex- 568  
 tending evaluation to nuScenes [7] and Waymo [49] would 569  
 strengthen generalization claims beyond the KITTI odome- 570  
 try benchmark. We believe continuous-space autoregressive 571  
 generation on raw sensor data represents a promising direc- 572  
 tion for scalable world modeling in autonomous systems. 573  
 The absence of a tokenization bottleneck makes this class 574  
 of architectures well-suited for large-scale pre-training on 575  
 unlabeled LiDAR corpora, with subsequent 3D fine-tuning 576  
 for safety-critical multi-step forecasting. 577

578

**References**

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

- [1] Michael Samuel Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *International Conference on Learning Representations (ICLR)*, 2023. 3
- [2] Michael S. Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *International Conference on Learning Representations*, 2023. 2, 3
- [3] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al.  $\pi_0$ : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024. 2
- [4] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 3
- [5] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023. 2
- [6] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 2020. 3
- [7] Holger Caesar et al. nuScenes: A multimodal dataset for autonomous driving. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 1, 6, 8
- [8] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T. Freeman. MaskGIT: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [9] Boyuan Chen et al. Diffusion forcing transformer for long-horizon video generation. *arXiv preprint*, 2024. 3, 4, 5
- [10] Boyuan Chen et al. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *arXiv preprint arXiv:2407.01392*, 2024. 2, 3, 4
- [11] Tiago Cortinhal, George Tzelepis, and Eren Erdal Aksoy. SalsaNext: Fast, uncertainty-aware semantic segmentation of LiDAR point clouds. In *Int. Symp. Visual Computing*, 2020. 2
- [12] Soham Dasgupta, Kshitij Aphale, Kaustab Pal, and Avinash Sharma. Accurate and real-time lidar point cloud forecasting for autonomous driving. New York, NY, USA, 2025. Association for Computing Machinery. 2
- [13] Jue Dong et al. Self-supervised point cloud prediction for autonomous driving. *IEEE Transactions on Intelligent Transportation Systems (T-ITS)*, 25(2), 2024. 2
- [14] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 3
- [15] Patrick Esser, Sumith Kulal, Andreas Blattmann, Robin Rombach, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *International Conference on Machine Learning (ICML)*, 2024. 2, 3
- [16] Hehe Fan and Yi Yang. PointRNN: Point recurrent neural network for moving point cloud processing. In *arXiv preprint arXiv:1910.08287*, 2019. 2
- [17] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2012. 6
- [18] William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weillbach, and Frank Wood. Flexible diffusion modeling of long videos. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 3
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Adv. Neural Inform. Process. Syst.*, 2020. 3
- [20] Tarasha Khurana, Peiyun Hu, David Held, and Deva Ramanan. Point cloud forecasting as a proxy for 4d occupancy forecasting. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1
- [21] Alex H Lang, Sourabh Vora, Holger Caesar, Lubomir Lublin, Oscar Spong, Susien Peters, and Juan Nuñez-Iglesias. PointPillars: Fast encoders for object detection from point clouds. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 2
- [22] Tianhong Li and Kaiming He. Back to basics: Let denoising generative models denoise. *arXiv preprint arXiv:2511.13720*, 2025. 2, 3, 5, 6
- [23] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation without vector quantization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 3
- [24] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. In *Int. Conf. Learn. Represent.*, 2023. 2, 3
- [25] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *International Conference on Learning Representations (ICLR)*, 2023. 3
- [26] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *International Conference on Learning Representations*, 2023. 2, 3
- [27] Fan Lu, Guang Chen, Zhijun Li, Lijun Zhang, Yinlong Liu, Sanqing Qu, and Alois Knoll. MoNet: Motion-based point cloud prediction network. *IEEE Trans. Intell. Transp. Syst.*, 23(8):13794–13804, 2021. 2
- [28] Zhen Luo, Junyi Ma, Zijie Zhou, and Guangming Xiong. PCPNet: An efficient and semantic-enhanced transformer network for point cloud prediction. *IEEE Robot. Autom. Lett.*, 8(7):4267–4274, 2023. 2, 6, 7, 8
- [29] Benedikt Mersch, Xieyuanli Chen, Jens Behley, and Cyrill Stachniss. Self-supervised point cloud prediction using 3d spatio-temporal convolutional networks. In *Conf. on Robot Learning*, 2022. 2, 6
- [30] Andres Milioto, Ignacio Vizzo, Jens Behley, and Cyrill Stachniss. RangeNet++: Fast and accurate LiDAR semantic seg-

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

- 691 mentation. In *IEEE/RSJ Int. Conf. Intell. Robot. Syst.*, 2019. 2
- 692
- 693 [31] Kaustab Pal, Aditya Sharma, Avinash Sharma, and K Mad- 748  
694 hava Krishna. ATPNet: Attention based temporal point 749  
695 cloud prediction network. *arXiv preprint arXiv:2401.17399*, 750  
696 2024. 2, 6, 7, 8 751
- 697 [32] William Peebles and Saining Xie. Scalable diffusion mod- 752  
698 els with transformers. In *Proceedings of the IEEE/CVF In- 753  
699 ternational Conference on Computer Vision (ICCV)*, pages 754  
700 4195–4205, 2023. 4, 5 755
- 701 [33] Adam Polyak, Amit Zohar, Andrew Brown, Andros Swami- 756  
702 nathan, Oran Metzger, et al. Movie gen: A cast of media 757  
703 foundation models. *arXiv preprint arXiv:2410.13720*, 2024. 758  
704 3 759
- 705 [34] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, 760  
706 Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 761  
707 Zero-shot text-to-image generation. In *International Confer- 762  
708 ence on Machine Learning (ICML)*, 2021. 3 763
- 709 [35] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generat- 764  
710 ing diverse high-fidelity images with VQ-VAE-2. In *Advances 765  
711 in Neural Information Processing Systems (NeurIPS)*, 2019. 766  
712 3 767
- 713 [36] Nicholas Rhinehart, Kris M. Kitani, and Paul Vernaza. R2p2: 768  
714 A Reparameterized pushforward policy for diverse, precise 769  
715 generative path forecasting. In *European Conference on Com- 770  
716 puter Vision*, 2018. 2 771
- 717 [37] Robin Rombach, Andreas Blattmann, Dominik Lorenz, 772  
718 Patrick Esser, and Björn Ömmer. High-resolution image syn- 773  
719 thesis with latent diffusion models. In *IEEE Conf. Comput. 774  
720 Vis. Pattern Recog.*, 2022. 3 775
- 721 [38] Kiwhan Song, Boyuan Chen, Max Simchowitz, Yilun Du, 776  
722 Russ Tedrake, and Vincent Sitzmann. History-guided video 777  
723 diffusion, 2025. 2 778
- 724 [39] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Ab- 779  
725 hishek Kumar, Stefano Ermon, and Ben Poole. Score-based 780  
726 generative modeling through stochastic differential equations. 781  
727 In *Int. Conf. Learn. Represent.*, 2021. 3 782
- 728 [40] Yue Song et al. PCPMamba: Point cloud prediction with 783  
729 mamba-based state space models. *arXiv preprint*, 2025. 2 784
- 730 [41] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Li- 785  
731 wei Wang. Visual autoregressive modeling: Scalable image 786  
732 generation via next-scale prediction. In *Advances in Neural 787  
733 Information Processing Systems (NeurIPS)*, 2024. Best Paper 788  
734 Award. 3 789
- 735 [42] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Mar- 790  
736 tinet, Marie-Anne Lachaux, Timothée Lacroix, et al. LLaMA: 791  
737 Open and efficient foundation language models. *arXiv 792  
738 preprint arXiv:2302.13971*, 2023. 3 793
- 739 [43] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete 794  
740 representation learning. *Advances in neural information pro- 795  
741 cessing systems*, 30, 2017. 2 796
- 742 [44] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, 797  
743 Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao 798  
744 Yang, et al. Wan: Open and advanced large-scale video 799  
745 generative models. *arXiv preprint arXiv:2503.20314*, 2025. 2 800
- 746 [45] Wan Team. Wan: Open and advanced large-scale video gen- 801  
747 erative models. *arXiv preprint arXiv:2503.01502*, 2025. 3 802
- [46] Xinshuo Weng, Jianren Wang, Sergey Levine, Kris Kitani, 748  
and Nicholas Rhinehart. Inverting the pose forecasting 749  
pipeline with SPF2: Sequential pointcloud forecasting for 750  
sequential pose forecasting. In *Conf. on Robot Learning*, 751  
2021. 2 752
- [47] Xinshuo Weng, Junyu Nan, Kuan-Hui Lee, Rowan McAllister, 753  
Adrien Gaidon, Nicholas Rhinehart, and Kris M Kitani. 754  
S2Net: Stochastic sequential pointcloud forecasting. In *Eur. 755  
Conf. Comput. Vis.*, 2022. 2 756
- [48] Bichen Wu, Alvin Wan, Xiangyu Yue, and Kurt Keutzer. 757  
Squeezeseg: Convolutional neural nets with recurrent crf 758  
for real-time road-object segmentation from 3d lidar point 759  
cloud. In *2018 IEEE International Conference on Robotics 760  
and Automation (ICRA)*, pages 1887–1893, 2018. 1 761
- [49] Runsheng Xu, Hubert Lin, Wonseok Jeon, Hao Feng, Yu- 762  
liang Zou, Liting Sun, John Gorman, Ekaterina Tolstaya, 763  
Sarah Tang, Brandyn White, et al. Wod-e2e: Waymo open 764  
dataset for end-to-end driving in challenging long-tail scenar- 765  
ios. *arXiv preprint arXiv:2510.26125*, 2025. 8 766
- [50] Zetong Yang, Li Chen, Yanan Sun, and Hongyang Li. Visual 767  
point cloud forecasting enables scalable autonomous driving. 768  
In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2024. 2 769
- [51] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu 770  
Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaoh- 771  
an Zhang, Guanyu Feng, et al. CogVideoX: Text-to-video 772  
diffusion models with an expert transformer. *arXiv preprint 773  
arXiv:2408.06072*, 2024. 3 774
- [52] Anthony Zhang et al. Copilot4D: Learning unsupervised 775  
world models for autonomous driving via discrete diffusion, 776  
2023. 2, 3, 4, 6, 8 777
- [53] Yin Zhou and Oncel Tuzel. VoxelNet: End-to-end learning 778  
for point cloud based 3d object detection. In *IEEE Conf. 779  
Comput. Vis. Pattern Recog.*, 2018. 2 780