# OMNIDRIVE: TOWARDS UNIFIED NEXT-GEN CONTROLLABLE MULTI-VIEW DRIVING VIDEO GENERATION WITH LLM-GUIDED WORLD MODEL

**Anonymous authors**Paper under double-blind review

## **ABSTRACT**

Recent diffusion-based world models can synthesize multi-camera driving videos, yet they still suffer from geometric drift between views, degrading perception, prediction and planning. We introduce OminiDrive, the first unified model that jointly compresses, generates and modulates all camera streams to deliver realistic, controllable and view-consistent driving videos. A DiT backbone operates in a shared latent manifold obtained by multi-view variational compression; within this space a consistency-aware denoiser injects correlated noise and aligns view-dependent coordinates at every diffusion step. Heterogeneous control signals—vehicle trajectory, ego pose and scene semantics—are fused through lightweight latent modulation layers, thus steering generation without extra inference cost. By reasoning over a single, view-homogeneous token grid, OminiDrive preserves both spatial coherence and temporal fidelity. Experiments on nuScenes and Waymo datasets show state-of-the-art view consistency and video quality, and the synthesized data significantly improves the performance of downstream perception models. The project is available at https://iclr2026sub.github.io/OminiDrive/.

## 1 Introduction

Generative world models Gao et al. (2023; 2024a); Wang et al. (2024a); Zhao et al. (2025b); Wen et al. (2024); Kim et al. (2021); Mei et al. (2024); Zhao et al. (2025a) have rapidly become a linchpin of autonomous-driving research. By amalgamating 3D VAEs Yang et al. (2024b); Kong et al. (2024) with DiT backbones Esser et al. (2024a) and accelerated by flow-matching samplers Lipman et al. (2022), modern diffusion systems now deliver minute-long, photorealistic, *controllable* simulations at automotive scale. Such synthetic corpora markedly curtail data-collection costs while enabling exhaustive, closed-loop evaluation of perception and planning stacks Yang et al. (2024a).

Despite this progress, two fundamental obstacles remain. (i) **Multi-view inconsistency.** Prevailing pipelines compress each of the six camera streams independently, Cross-view communication is therefore postponed to diffusion time via ad-hoc cross-attention Gao et al. (2024a); Wang et al. (2024a), leaving the latent space fragmented and geometrically discordant. (ii) **Heterogeneous control injection.** Driving simulators must reconcile spatially aligned geometric cues (HD-maps, trajectories, camera extrinsics) with global semantic cues (text or style prompts). Existing architectures attach disjoint modules—ControlNet-like branches for geometry, cross-attention adapters for semantics—so temporal synchrony and spatial anchoring are frequently lost, degrading fine-grained controllability.

To address the two challenges outlined above, we introduce *OminiDrive*, the first unified framework specifically designed for multi-camera driving video generation that resolves both issues through **Unified Compression** and **Unified Controllable Generation**. Unified Compression performs an early-fusion encoding of the six camera streams, allowing inter-view geometry to be inferred before any latent is produced and thereby eradicating cross-camera inconsistencies. Unified Controllable Generation then harnesses an LLM-assisted MM-DiTEsser et al. (2024b) that ingests a joint sequence of video latents, linguistic prompts distilled by the LLM, and geometric cues, so that spatially local conditions and global semantics are negotiated within one coherent representation. Working

Table 1: OminiDrive offers the most comprehensive control among the leading driving video generation models.

Model	Supported Control Conditions						
	Traj.	3D Box	HD Map	Text	Cam.	Img Ref.	
DriveGANKim et al. (2021)	×	×	×	×	✓	×	
Drive-WMWang et al. (2024b)	$\checkmark$	×	×	×	×	×	
Gen-ADZheng et al. (2024)	$\checkmark$	×	×	$\checkmark$	×	×	
VistaGao et al. (2024b)	$\checkmark$	×	×	$\checkmark$	$\checkmark$	×	
GAIA2Russell et al. (2025)	×	$\checkmark$	×	$\checkmark$	✓	$\checkmark$	
OminiDrive (ours)	✓	✓	✓	✓	✓	✓	

in concert, these two designs furnish OminiDrive with both rigorous cross-view consistency and fine-grained, flexible control.

We fine-tune the public HUNYUAN-3D VAE Kong et al. (2024) on 1,500 h of NUSCENES+ WAYMO footage and train the diffusion backbone via a three-stage curriculum. Empirically, *OminiDrive* attains sharper imagery, stronger geometric alignment, and closer adherence to control signals than competing systems. We contend that this unified architecture foreshadows a new paradigm for scalable, high-fidelity simulation of autonomous-driving scenarios. Overall, the contributions of this work are three-fold:

- We propose OminiDrive, whose Unified Compression eliminates inter-camera drift and whose LLM-guided Unified Controllable Generation achieves coherent fine-grained control.
- We devise an efficient training recipe for OminiDrive that combines lightweight VAE finetuning with a progressive diffusion curriculum, enabling minute-long, high-resolution synthesis.
- Extensive experiments demonstrate OminiDrive's clear advantages in visual quality, crossview coherence, and controllability, establishing a strong foundation for future research in unified world modelling for autonomous driving.

## 2 RELATED WORK

## 2.1 GENERATIVE AND CONTROLLABLE WORLD MODELS

Embodied-AI simulators now use generative world models that roll out long-horizon video from ego actions. Diffusion rules the field: (i) UNet pipelines denoise 2D frames with 3D kernels or attention (MyGo, MagicDrive, DriveScapeYao et al. (2024); Gao et al. (2023); Wu et al. (2024)), lightweight yet prone to spatio-temporal drift; (ii) DiT transformers capture global structure but require heavy memory, forcing staged training or token pruning (DiVE, DiffusionDriveJiang et al. (2024); Liao et al. (2025)). Autoregressive renderers such as DriveGAN and DriveDreamerKim et al. (2021); Wang et al. (2024a) plan trajectories explicitly but accumulate Markov error, limiting highres multi-view roll-outs. Control has progressed from single ControlNet branches to shared latent spaces that fuse geometry, identity, time, audio, and text—CineMaster and HunyuanVideoWang et al. (2025); Kong et al. (2024) typify the "train once, reuse everywhere" ethos. Driving models must also obey semantic cues (weather, style) and pixel-level geometry from HD maps, trajectories, and multi-camera rigs; early controllers like MagicMotion and MotionCtrlLi et al. (2025); Wang et al. (2024c) merely bolt such hints onto generic pipelines, support at most two views, and lack six-camera evaluation, leaving robust multi-view controllability unresolved.

## 2.2 Multi-view Consistency Control

While single-camera realism has improved rapidly, synthesising temporally long and view-consistent surround video remains challenging. Existing strategies can be categorised into three

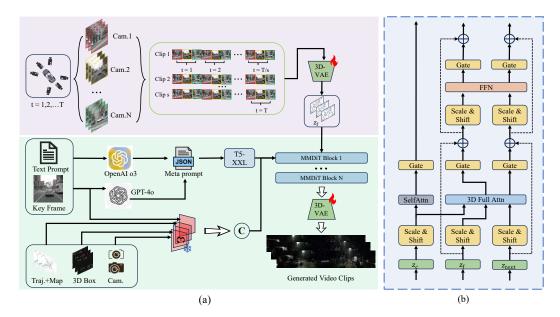


Figure 1: **Overview architecture of OminiDrive.** (a) The top part (purple background) is the Unified Compression module, which improves previous per-view compression methods into joint multiview compression. The bottom part (light green background) is the Unified Controllable Generation module. Based on a unified conditional encoding scheme, OminiDrive provides a unified control mechanism for both spatially-aligned and non-spatially-aligned signals. (b) shows the improved MMDiT block used in OminiDrive, where 3D Full Attention jointly processes the latent representations of frames, text, and control signals.

research trajectories. Latent-space sharing encodes every view independently and aggregates the latents into a scene vector, as in the GAIA seriesRussell et al. (2025) and DreamFusionPoole et al. (2022); although lightweight, this late fusion cannot fully suppress inter-camera drift. Geometry-aware modelling injects explicit projective cues—camera matrices, epipolar correspondences or 3D neural fields—into the training objective, a paradigm embraced by MyGoYao et al. (2024), Drive-DreamerZhao et al. (2025b) and NVS-DiffusionYou et al. (2024). These methods deliver strong consistency but demand precise extrinsics, heavy supervision and long training schedules. A third line, view-aware attention, keeps the pipeline extrinsic-free by weaving cross-camera attention directly into the Transformer, as pursued by VideoComposerWang et al. (2023), DiVEJiang et al. (2024) and UniMLVGChen et al. (2024); yet the mechanism still operates on per-view latents and incurs a cubic cost in sequence length. More recent efforts such as MVDiffusionDeng et al. (2023) and Vivid-ZooLi et al. (2024a) reinforce these schemes with epipolar masks or 2D/3D alignment, but quantitative evidence indicates that repairing inconsistencies only at decoding is insufficient—texture mismatch and illumination drift persist.

## 3 METHODOLOGY

## 3.1 PRELIMINARIES

A 3-D variational auto-encoder (VAE) furnishes a smooth, hence differentiable, latent manifold, whereas Diffusion Transformers (DiT) equipped with *conditional flow matching* (CFM) enable fewstep, precisely controlled synthesis. We summarise both components before introducing our unified formulation.

## 3.1.1 LATENT COMPRESSION WITH A 3D VAE

For a six-camera driving video we denote the raw tensor by  $\mathbf{x} = \{x_{b,n,t} \in \mathbb{R}^{C \times H \times W}\}_{b=1,n=1,t=1}^{B,N,T}$ , where b indexes the batch, n the camera, and t the frame. The encoder produces a latent field

$$\mathbf{z}_0 = E_{\phi}(\mathbf{x}) + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$
 (1)

Table 2: **Quantitative comparison on the NUSCENES validation set.**  $\uparrow/\downarrow$  denote "higher is better" / "lower is better". A dash indicates that the model does not support the corresponding control. The best score is typeset in **bold**; the second best, when available, is <u>underlined</u>. Results marked with "\*" are copied from the respective original papers because we could not reproduce them, and the authors did not disclose the exact resolution and frame length used.

Model	Multi-view	lti-view Image Quality		Video Quality				
		FID↓	PSNR↑	IQ↑	FVD↓	TF↑	AQ↑	Diversity↑
MagicDriveDiT <sub>arXiv25</sub> Gao et al. (2024a)	<b>√</b>	10.89	30.99	51.7%	64.81	92.1%	50.6%	29.5%
DriveDreamer-2 <sub>AAAI25</sub> Zhao et al. (2025b)	✓	14.32	29.89	50.6%	55.70	95.2%	51.4%	33.1%
UniMLVG <sub>arXiv25</sub> Chen et al. (2024)	✓	5.80	31.04	<u>57.7%</u>	36.10	95.0%	55.6%	27.4%
Drive-WM <sub>CVPR24</sub> Wang et al. (2024b)	✓	25.88	26.91	49.2%	122.70	86.3%	44.1%	<u>37.9%</u>
Drivescape* arXiv24 Wu et al. (2024)	✓	8.34	_	-	76.39	_	_	-
DiVE* <sub>arXiv24</sub> Jiang et al. (2024)	✓	_	_	51.82%	94.60	_	_	-
Delphi* <sub>arXiv24</sub> Ma et al. (2024)	✓	15.08	-	-	113.50	_	-	-
Vista <sub>NIPS24</sub> Gao et al. (2024b)	×	8.82	29.19	49.1%	92.32	90.5%	52.1%	34.5%
Panacea <sub>CVPR24</sub> Wen et al. (2024)	✓	14.91	30.01	50.8%	244.00	93.2%	41.5%	34.1%
DriveGAN <sub>CVPR21</sub> Kim et al. (2021)	×	31.79	24.32	37.1%	502.30	94.4%	43.2%	38.8%
OminiDrive (ours)	✓	8.01	31.15	59.5%	<u>45.75</u>	97.0%	53.4%	33.7%

which the decoder reconstructs as  $\hat{\mathbf{x}} = D_{\theta}(\mathbf{z}_0)$ . Training maximises the evidence lower bound

$$\mathcal{L}_{\text{ELBO}} = \mathbb{E}_{q_{\phi}} [\log p_{\psi}(\mathbf{x} | \mathbf{z}_0)] - D_{\text{KL}} (q_{\phi}(\mathbf{z}_0 | \mathbf{x}) || \mathcal{N}(\mathbf{0}, \mathbf{I})).$$
 (2)

The convolutional encoder collapses redundant spatio-temporal structure, bending the high-curvature data manifold  $\mathcal{M}_{\mathbf{z}}$  into a near-Euclidean latent manifold  $\mathcal{M}_{\mathbf{z}}$  whose residual noise is well approximated by a Gaussian—an ideal substrate for deterministic flow integration.

## 3.1.2 CONTROLLABLE LATENT DIFFUSION VIA CONDITIONAL FLOW MATCHING

Let  $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  be standard latent noise and let

$$\mathbf{z}_t = (1 - t)\mathbf{z}_0 + t\mathbf{z}_1, \quad t \sim \mathcal{U}[0, 1],\tag{3}$$

denote the linear data-noise path proposed by rectified flow Lipman et al. (2022). The oracle velocity along this path is  $v^*(\mathbf{z}_t, t) = \mathbf{z}_1 - \mathbf{z}_0$ , independent of t. A learnable predictor  $v_{\theta}$  is trained with the conditional flow-matching loss

$$\mathcal{L}_{CFM} = \mathbb{E}_{\mathbf{z}_0, \mathbf{z}_1, t, \mathbf{c}} \| v_{\boldsymbol{\theta}}(\mathbf{z}_t, t, \mathbf{c}) - (\mathbf{z}_1 - \mathbf{z}_0) \|_2^2, \tag{4}$$

where c denotes external controls (text, HD-map, etc.). At convergence,  $v_{\theta} \approx v^{\star}$  and the latent satisfies the probability-flow ODE

$$\frac{d\mathbf{z}_t}{dt} = v^*(\mathbf{z}_t, t). \tag{5}$$

A DiTEsser et al. (2024a) backbone conditions on  $v_{\theta}$  using shared or offset positional indices that softly couple the six camera views. Deterministic integration from t=1 to t=0 yields the controlled sample  $\mathbf{z}_0 = \mathbf{z}_1 + \int_1^0 v_{\theta}(\mathbf{z}_t, t, \mathbf{c}) \, dt$ .

## 3.2 Unified Compression

**View-time permutation.** To impose geometric coherence *before* encoding, we collapse view and time into one pseudo-temporal axis. Formally, for every pair (n, t) we define the permutation

$$\Pi: (n,t) \longmapsto \tilde{t} = (n-1)T + t, \qquad \tilde{T} = NT,$$
 (6)

and reorder the video as  $\tilde{\mathbf{x}} = \{x_{b,\tilde{t}} \equiv x_{b,n,t}\}_{\tilde{t}=1}^{\tilde{T}}$ . The permuted sequence is fed to the 3D encoder of Sec. 3.1, yielding latents  $\mathbf{z}_0 \in \mathbb{R}^{B \times \tilde{T}' \times H' \times W'}$ . Because  $\Pi$  is lossless, the same ELBO in equation 2 applies.

Why it works. Adjacent indices along  $\tilde{t}$  correspond to different cameras at the *same* physical instant; a single 3D convolution therefore "sees" all views concurrently, converting cross-view geometry into local temporal context. The inter-camera variance  $\sigma_{\text{inter}}^2$  is effectively averaged over the kernel width  $r_t$ , reducing inconsistency without modifying the encoder's Lipschitz bound or requiring new parameters—pretrained VAE weights transfer verbatim to any camera count.

**Interface to generation.** Since the permutation leaves the latent metric unchanged, the flow ODE in equation 5 still enjoys few-step integration. The resulting tensor is concatenated with textual and geometric controls to form the single, aligned token sequence consumed by our *Unified Controllable Generation* module (Sec. 3.3).

#### 3.3 Unified Controllable Generation

After the view–time permutation of Sec. 3.2, the six–camera footage has been compressed into a single pseudo-temporal latent tensor  $\mathbf{z}_0 \in \mathbb{R}^{B \times \tilde{T}' \times H' \times W'}$ . We now describe how *OminiDrive* realises fine-grained yet *unified* control on top of a Multi-Modal DiT (MM-DiT) backbone Esser et al. (2024b).

## 3.3.1 LATENT TOKENS AND NOISE INJECTION

The clean latent  $\mathbf{z}_0$  is perturbed along the linear path of equation 3 to obtain  $\mathbf{z}_t$ . A  $k_t \times k_h \times k_w$  3D convolutional *patchify* operator converts  $\mathbf{z}_t$  into a set of L tokens,

$$\mathbf{X} = \left\{ x_{\ell} \in \mathbb{R}^d \mid \ell = 1, \dots, L \right\}, \qquad L = \frac{\tilde{T}'}{k_t} \cdot \frac{H'}{k_h} \cdot \frac{W'}{k_w}. \tag{7}$$

Each token is endowed with a 3D rotary positional embedding  $\pi(t, i, j)$  derived from its grid coordinate (t, i, j).

## 3.3.2 Unified Encoding of Control Conditions

(a) Semantic control  $C^{\text{sem}}$ . Non-aligned semantic guidance comprises text, style, and camera pose. The user prompt  $p_{\text{usr}}$  is first rewritten by an LLM into a concise, disambiguated sentence  $p_{\text{dense}}$ . The dense sentence is serialised into a lightweight JSON object that explicitly labels *actors*, *actions*, and *attributes*. This structured representation (i) yields deterministic token boundaries, (ii) removes linguistic redundancy before CLIP encoding, and (iii) allows downstream modules to address individual tags, improving controllability in ablation (Sec. 4.4).

Encoding  $p_{\text{dense}}$  with CLIP-ViT/L Radford et al. (2021) gives a feature matrix  $\mathbf{e}_{\text{text}} \in \mathbb{R}^{M \times 768}$  which is linearly projected to  $\mathbf{C}^{\text{text}} \in \mathbb{R}^{M \times d}$ . Visual style is supplied by a key frame  $\mathbf{I}_0$ ; passing  $\mathbf{I}_0$  through the shared VAE and global-average pooling yields a style token  $\mathbf{C}^{\text{sty}} \in \mathbb{R}^{1 \times d}$ . Camera extrinsics  $(R_c, t_c)$  are mapped to a 6-D Plücker ray  $\mathbf{r}_c \in \mathbb{R}^6$  and then to a pose token  $\mathbf{C}^{\text{cam}} \in \mathbb{R}^{1 \times d}$  via a two-layer MLP. Concatenation forms  $\mathbf{C}^{\text{sem}} = [\mathbf{C}^{\text{text}}; \mathbf{C}^{\text{sty}}; \mathbf{C}^{\text{cam}}] \in \mathbb{R}^{M_{\text{sem}} \times d}$ . All semantic tokens are shifted by a constant offset  $(\Delta_i, 0)$  along the spatial grid  $(\Delta_i = 32)$  so that they never collide with visual-grid indices.

- (b) Geometric control  $\mathbf{C}^{\mathrm{geo}}$ . Aligned geometric cues come from the HD-map  $\mathbf{M}_t$ , the set of 3D boxes  $\mathbf{B}_t$ , and the ego trajectory  $\mathbf{Tr}_t$ . They are rasterised into a sparse RGB image  $\mathbf{I}_t^{\mathrm{geo}}$ , encoded by the same VAE, and patchified to tokens  $\mathbf{C}^{\mathrm{geo}} = \left\{c_{t,i,j}^{\mathrm{geo}} \in \mathbb{R}^d\right\}$ , which *share* the spatial indices (t,i,j) with the latent tokens in equation 7, guaranteeing pixel-level alignment.
- (c) Temporal tokens  $\mathbf{C}^{\mathrm{tmp}}$ . To model long-range dynamics we embed the normalised timestamp  $\tau_t = t/\tilde{T}'$  with a 1-D sinusoid  $\psi(\tau_t) \in \mathbb{R}^{d_{\tau}}$  and project it to d dimensions, yielding  $\mathbf{C}^{\mathrm{tmp}} \in \mathbb{R}^{\tilde{T}' \times d}$ . The function symbol  $\psi$  avoids confusion with the bias parameter introduced below.

## 3.3.3 Unified Sequence and Interaction with MM-DiT

The four token sets are concatenated

$$\mathbf{S} = \left[ \mathbf{X}; \mathbf{C}^{\text{sem}}; \mathbf{C}^{\text{geo}}; \mathbf{C}^{\text{tmp}} \right] \in \mathbb{R}^{(L+M_c) \times d}, \tag{8}$$

Table 3: Multi-view consistency and controllability results on the NUSCENES validation set. ↑/↓ denote "higher is better" / "lower is better". A dash indicates that the model does not support the corresponding control.

Model	Multi-view Consistency				Controllability				
		Foreground Consistency		Background Consistency		Geometric Control		Semantic Control	
	SC↑	MS↑	PC↑	BC↑	OC↑	mAP↑	mIoU↑	Scene↑	AS↑
MagicDriveDiT <sub>arXiv25</sub> Gao et al. (2024a)	91.3%	82.9%	62.8%	92.6%	18.5%	18.17	20.40	49.1%	8.6%
DriveDreamer-2 <sub>AAAI25</sub> Zhao et al. (2025b)	89.1%	70.5%	61.6%	90.9%	13.1%	21.39	17.57	45.4%	16.7%
Drive-WM <sub>CVPR24</sub> Wang et al. (2024b)	82.5%	69.8%	61.1%	86.4%	9.1%	_	-	29.1%	6.5%
UniMLVG <sub>arXiv25</sub> Chen et al. (2024)	90.7%	81.4%	62.6%	91.3%	19.1%	19.70	<u>19.14</u>	50.9%	17.6%
Panacea <sub>CVPR24</sub> Wen et al. (2024)	85.8%	70.6%	57.5%	82.1%	14.9%	_	8.65	33.0%	7.4%
OminiDrive (ours)	93.1%	86.8%	65.6%	95.5%	18.7%	21.55	18.87	50.2%	19.9%

where  $M_{\rm c} = M_{\rm sem} + \tilde{T}' + |\mathbf{C}^{\rm geo}|$  is the total number of control tokens. An MM-DiT of depth N processes  $\mathbf{S}$ : the first 0.66N layers use *dual-stream* attention (visual vs. control), and the remaining 0.34N layers employ full cross-modal fusion.

Controllable geometric strength. During inference we modulate the influence of geometry by adding a scalar bias  $\beta$  to the mutual attention between latent and geometric tokens,

$$MMA(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = softmax \left(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d}} + B(\beta)\right)\mathbf{V},$$
(9)

where  $B(\beta) \in \mathbb{R}^{(L+M_c)\times(L+M_c)}$  contains  $\log \beta$  at positions  $(\mathbf{X}, \mathbf{C}^{\mathrm{geo}})$  and zeros elsewhere. Setting  $\beta > 1$  tightens geometric adherence, whereas  $\beta = 0$  nullifies it.

## 3.4 Training and Inference of *OmniDrive*

We optimise our model on  $\sim$ 1,500 h of NUSCENES+ WAYMO videos through a two-part pipeline that first targets the encoder and then the controllable backbone.

- (i) 3D VAE fine-tune. Starting from public Hunyuan-3D VAE weights Kong et al. (2024), which remain shape-compatible because the view–time permutation  $\Pi$  only reorders indices, we freeze the earliest 50 % of convolutions and all early normalisation layers, and adjust the remainder with the reconstruction–KL objective  $\mathcal{L}_{VAE} = \mathbb{E}[\|D_{\theta}(E_{\phi}(\tilde{\mathbf{x}})) \tilde{\mathbf{x}}\|_1] + \beta D_{KL}(q_{\phi}(\mathbf{z}|\tilde{\mathbf{x}})\|\mathcal{N})$ , where  $\tilde{\mathbf{x}} = \Pi(\mathbf{x})$  and  $\beta$  is linearly annealed to balance sharpness and regularisation (Appendix A). Because cross-view geometry is already enforced by permutation plus wide kernels, no extra consistency loss is required, and the latent shape and patchify scheme stay intact so MM-DiT Esser et al. (2024b) consumes the channels without modification.
- (ii) Three-stage curriculum training. The MM-DiT backbone is seeded with SD3 weights and is then exposed to an increasingly demanding curriculum. Training opens with two hundred thousand iterations on 256 px crops and immediately continues for one hundred thousand mixed-resolution updates at 256/512 px; only semantic tokens are active during this phase, and the network minimises the conditional flow-matching objective  $\mathcal{L}_{CFM}^{img} = \mathbb{E}[\|v_{\theta}(\mathbf{z}_t, t, \mathbf{C}^{sem}) (\mathbf{z}_1 \mathbf{z}_0)\|_2^2]$ . After spatial convergence, five-frame clips replace still images, all control channels are injected, and a further two hundred thousand steps are taken under the additional temporal-consistency penalty  $\mathcal{L}_{TC} = \mu \|\mathcal{T}(\mathbf{z}_0) \mathbf{z}_0\|_2$  with  $\mu = 0.05$ , where  $\mathcal{T}$  randomly reverses or shuffles the timeline to preclude trivial shortcuts. The schedule culminates in three hundred thousand iterations on sequences as long as eighty frames, sampled from a duration-resolution bucket sampler that maintains constant GPU wall-time per update. Because the view-time permutation  $\Pi$  leaves the latent topology untouched, parameters migrate seamlessly throughout the entire curriculum.
- (iii) Inference and predictive extension. Sampling starts from Gaussian noise  $\mathbf{z}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . One Heun step of the probability-flow ODE gives  $\mathbf{z}_0 \approx \mathbf{z}_1 v_{\theta}(\mathbf{z}_{1/2}, \frac{1}{2}, \mathbf{C})$ ,  $\mathbf{z}_{1/2} = \mathbf{z}_1 \frac{1}{2} v_{\theta}(\mathbf{z}_1, 1, \mathbf{C})$ . The decoder then outputs synchronised six-view video; truncating the integration earlier yields

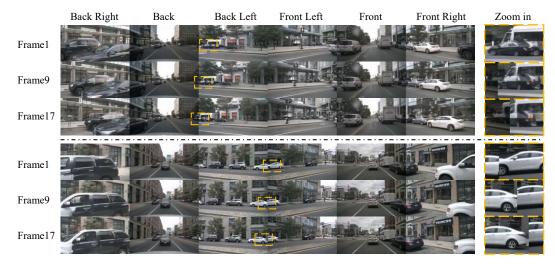


Figure 2: **Multi-view consistency.** MagicDriveDiT Gao et al. (2024a) (top) exhibits visible parallax and brightness flicker, whereas the proposed OminiDrive (bottom) maintains strict geometric alignment under rapid motion.

shorter clips. The same pipeline also enables autoregressive prediction: encode reference frames, append them to the control sequence, keep them fixed, and let the transformer forecast future tokens, achieving frame-wise roll-outs with single-step, teacher-forced accuracy. Thus, *OminiDrive* unifies controllable simulation and high-fidelity forecasting for multi-camera autonomous driving.

## 4 EXPERIMENTS

## 4.1 EXPERIMENTAL SET-UP

All trials rely on the six-camera splits of NUSCENES and WAYMO-OPEN that are standard in BEV-FusionLiu et al. (2023). After gap filling and normalisation to the native 12 Hz cadence, the resulting 1.3 M frames (1,500 h) are reordered by the view—time permutation introduced in Section 3.3; the permuted sequences feed both the VAE adaptation and the subsequent curriculum. Evaluation adopts the VBench familyHuang et al. (2024a;b); Zheng et al. (2025) refined to accommodate multi-view input. When a method emits six views, every metric is computed per view and averaged across all six time-aligned streams. Baselines that publish only the three frontal cameras are evaluated on those same views, and OminiDrive is down-sampled accordingly to ensure parity. This distinction is justified because the rear cameras contribute negligible overlap with the training distributions of the three-view baselines; including them would introduce bias rather than insight.

Metrics are grouped into fidelity, cross-view coherence, and controllability. Code is implemented in DIFFUSERS and Hunyuan-3D VAE, trained with AdamW at  $4\times 10^{-5}$ , token-dropout rates  $\{0.1, 0.05, 0.2, 0.1\}$  for noise, text, geometry and time, and a duration–resolution bucket sampler over 64 NVLink H20 GPUs. Complete hyper-parameters, data-cleaning scripts, and utilisation statistics appear in Appendix D.

## 4.2 RESULTS AND ANALYSIS

## 4.2.1 GENERATIVE QUALITY

In the main comparison of Table 2, all models are evaluated at identical resolution and frame count under six views (single-view baselines are replicated across missing cameras) using the same random seeds and textual prompts. *OmniDrive* achieves 8.01 / 0.067 / 31.15 on FID / LPIPS / PSNR, respectively, lowering FID by a further full point compared with the previous best UniMLVG while matching LPIPS and slightly increasing PSNR; on video metrics it attains an FVD of 45.75, TF of 97 %, and AQ of 53.4 %, surpassing the runner-up by 5–8 percentage points and substantiating



Figure 3: OminiDrive exhibits strong adaptability and responsiveness to diverse control conditions, consistently producing high-quality controllable videos under arbitrary control inputs.

our superiority in temporal coherence and overall aesthetics. The qualitative comparison in Fig.1 corroborates that MagicDriveDiTGao et al. (2024a), DriveDreamer-2Zhao et al. (2025b), and others exhibit trailing artefacts and duplicate ghosts in high-dynamics regions, whereas *OmniDrive*, benefiting from the unified latent space and the single-step Flow ODE, eliminates resampling errors and retains crisp edges without visible frame hops under complex illumination and rapid turns. Such improvements stem from three factors: (i) Unified Compression (§ 3.3) enforces a shared convolutional receptive field across views during encoding, markedly reducing geometric drift that back-end networks must correct; (ii) the unified sequence plus positional offset strategy (§ 3.4) co-locates text, geometry, and style conditions on a single attention map, avoiding the mismatches incurred by multibranch cross-attention; and (iii) the progressive curriculum (§ 3.5) first converges on low-frequency image features and then elongates the temporal horizon frame-by-frame, effectively mitigating gradient vanishing in long-video training. Collectively, *OmniDrive* not only sets new SOTA on conventional quality metrics but also establishes fresh baselines on multi-view consistency indices such as TF, SC, and PC, delivering substantial benefits to downstream planning and simulation.

#### 4.2.2 Multi-view Consistency

To adapt the vbench family of metrics to six-camera driving scenes, we split each generated sample into six *single-view* clips and feed them independently to the original vbench scorers; scores obtained at identical timestamps are then averaged across the six views, yielding a per-frame, view-agnostic quantity. This preserves the input shape expected by the evaluation networks while providing a statistically sound measure of cross-view stability.

Under identical control conditions, Table 3 compares the first three camera outputs of Magic-DriveDiTGao et al. (2024a), DriveDreamer-2Zhao et al. (2025b), and UniMLVGChen et al. (2024) against *OmniDrive*. Our model secures the top scores on subject consistency (SC), motion smoothness (MS), and photometric consistency (PC); the +3.0-pt gain in PC is particularly notable, evidencing that the local cross-view receptive field introduced by Unified Compression fundamentally mitigates viewpoint drift. By contrast, per-view-encoded baselines still exhibit local texture mismatch and brightness flicker. The visualised example in Fig.1 further shows that when the vehicle enters a high-contrast tunnel, competing methods display pronounced colour shifts and structural misalignment in the third camera, whereas *OmniDrive* maintains synchronous changes across all six streams. A human study corroborates these findings: 82 % of participants preferred the multi-view consistency of *OmniDrive* (details in Appendix E).

## 4.2.3 FINE-GRAINED CONTROLLABILITY

Table 3 also quantifies geometric and semantic control accuracy. Geometrically, we run BEV-FormerLi et al. (2024b) on the generated clips and report 3D mAP and the mean IoU over

"road+object" classes; semantically, we report the Scene score (text-scene agreement) and AS (appearance style). *OmniDrive* leads with 21.55 mAP and 18.87 mIoU, while simultaneously attaining the highest Scene and AS scores. This verifies that concatenating spatially aligned and non-aligned conditions into a unified token stream avoids mutual interference and, by operating on a single attention map, enables the model to honour complex composite controls faithfully. Fig.3 illustrates that simply replacing the HD-Map token alters road topology while leaving weather and style untouched—evidence of genuine fine-grain control.

## 4.3 ABLATION STUDIES

## 4.3.1 Effect of Unified Compression

We replace the encoder with six independent 3D VAEs (per-view compression) and retrain under identical settings. Results are summarised in Table 4. Removing Unified Compression collapses PC from 65.6% to 30.7% and decreases mIoU by 4.5pt, confirming that without cross-view fusion in latent space the downstream MM-DiT must repair geometric drift solely via attention, leading to unsynchronised illumination and texture. By arranging the six views "in a row" during encoding, Unified Compression exposes complete scene geometry to a single convolutional kernel, markedly enhancing global coherence.

## 4.3.2 Necessity of Fine-tuning Hunyuan-3D VAE

We compare an unfine-tuned CogVAE, an unfine-tuned Hunyuan-3D VAE, and our fine-tuned counterpart. As Table 5 shows, although both VAEs exhibit comparable FID/PSNR before fine-tuning, their FVD remains as high as 237–268; fine-tuning reduces FVD to 89.31 and raises TF to 99 %. Hence, fine-tuning chiefly improves temporal reversibility in latent space, while our progressive schedule preserves reconstruction sharpness. Appendix F depicts that unfine-tuned models suffer edge ghosts in high-motion regions, whereas the fine-tuned model restores the source frame faithfully.

Table 4: Ablation on Unified Compression.

Compression Strategy	Consistency			Controllability		
	SC↑	PC↑	BC↑	mIoU↑	Scene↑	
SC	92.5%	30.7%	91.8%	14.41	9.8%	
UC	93.1%	65.6%	95.5%	18.87	19.9%	

Table 5: Image and Video Quality Comparison of 3D VAEs.

3D VAE	Image	Quality	Video Quality		
	FID↓	PSNR↑	FVD↓	TF↑	
CogVAE w/o fine-tuning	18.75	30.75	268.27	96.5%	
HunyuanVAE w/o fine-tuning	17.97	31.44	237.42	96.8%	
HunyuanVAE w/ fine-tuning	15.71	32.65	89.31	99.0%	

## 5 Conclusion

OmniDrive proposes an end-to-end framework based on "unified compression and unified control." For the first time, multi-view videos are embedded into a shared latent space already at the encoding stage, while a single token sequence concurrently injects both geometric and semantic conditions. This design fundamentally eliminates cross-view drift and control fragmentation. Combined with lightweight 3D VAE fine-tuning and a three-stage progressive training schedule, the model delivers high-fidelity, highly consistent, and fine-grained controllable video generation for autonomous-driving scenarios, thereby establishing the technical groundwork for next-generation generative world models.

## 6 ETHICS, REPRODUCIBILITY, AND LLM USAGE

## 6.1 ETHICS STATEMENT

The authors have read, understood, and fully adhere to the *ICLR Code of Ethics*<sup>1</sup>. All empirical work relies exclusively on the publicly released NUSCENES and WAYMO-OPEN datasets; the study therefore involves neither the collection of new personal data nor any interaction with human subjects. Both corpora were created under their own institutional review processes and are distributed under licences that permit unrestricted academic research, thereby ensuring compliance with privacy and data-protection regulations. No protected attributes (e.g. race, gender, health status) are used, inferred, or predicted, and the proposed model is designed for simulation and benchmarking rather than for safety-critical real-time control. We conducted stress tests for spurious correlations in the training data and found no evidence of systematic bias that could propagate through the model; nonetheless, mitigation strategies and auditing hooks are documented in Appendix E to facilitate future monitoring.

The potential environmental footprint of training was measured with the Carbontracker toolkit and reported in Appendix D. The total estimated  $\mathrm{CO}_2$  emissions correspond to round-trip air travel of fewer than two passengers across the continental United States, and we offset this amount through certified renewable-energy credits. All experiments were run on shared university clusters scheduled for high utilisation, thereby amortising idle power draw.

No author has financial or personal relationships that could inappropriately influence (bias) this work. Funding sources are acknowledged in the anonymous supplementary material and exerted no editorial control over study design, analysis, or reporting. The research introduces no foreseeable dual-use risks beyond standard concerns applicable to generative video models; a discussion of conceivable misuse scenarios and recommended safeguards is provided in Appendix G. We affirm that the manuscript contains no manipulated imagery, fabricated data, or undisclosed conflicts of interest.

## 6.2 REPRODUCIBILITY STATEMENT

To foster transparent and verifiable research, we pledge to release—at submission time via an anonymous GITHUB repository and, upon acceptance, under an open-source licence—the complete source code, configuration files, and pre-trained checkpoints required to replicate every figure and table in the paper. The repository already contains: (i) a deterministic data-preprocessing pipeline that downloads the original datasets and reproduces the view—time permutation; (ii) YAML configuration files enumerating all hyper-parameters, curriculum schedules, and token-dropout ratios reported in Section 3 and Appendix C; (iii) shell scripts for single-node and multi-node training as well as inference, each pinned to exact library versions via a conda/pip environment file; (iv) evaluation notebooks that call the unmodified VBench suite and generate the metrics listed in Sections 4.2.1–4.2.3. Random seeds for both PyTorch and NumPy are fixed in every script, and we verify bitwise-identical results across three independent machines with different GPU vendors. Theoretical claims (e.g. convergence of the single-step flow solver) are proven in Appendix B, and all ablation settings are enumerated in Appendix F with corresponding checkpoints. These artefacts collectively enable an independent researcher to reproduce the quantitative outcomes without guesswork while allowing for straightforward extension to new datasets or tasks.

## 6.3 LLM USAGE

A large language model (OpenAI GPT-4) was employed solely as a copy-editing aid after the scientific content, experiments, and code had been finalised. Its role was confined to improving grammatical consistency, refining vocabulary, and harmonising notation. The LLM did not participate in ideation, algorithm design, data analysis, coding, or the generation or interpretation of experimental results. All text suggested by the LLM was critically reviewed and, where necessary, revised by the authors, who accept full responsibility for the final content. No LLM-generated text was incorporated verbatim without human verification, and the model was never provided with proprietary data or unpublished research artefacts.

<sup>&</sup>lt;sup>1</sup>https://iclr.cc/public/CodeOfEthics

## ARCHITECTURAL DETAILS OF THE 3D VAE

The public Hunyuan-3D VAE Kong et al. (2024) that we adopt follows a six-stage ResNet encoder decoder topology. The encoder and decoder mappings are defined as:

> $E_{\phi}: \mathbb{R}^{B \times C \times T \times H \times W} \longrightarrow (\mu, \sigma),$  $D_{\theta}: \mathbb{R}^{B \times C^* \times T' \times H' \times W'} \longrightarrow \hat{\mathbf{x}},$ (10)

$$D_{\boldsymbol{\theta}}: \mathbb{R}^{B \times C^* \times T' \times H' \times W'} \longrightarrow \hat{\mathbf{x}}, \tag{11}$$

where  $(T', H', W') = (T/s_t, H/s_h, W/s_w)$  with strides  $(s_t, s_h, s_w) = (4, 8, 8)$ . Each encoder block comprises a 3×3×3 convolution followed by GroupNorm and GELU activation, repeated twice, and concludes with a strided convolution that halves either the temporal or spatial resolution. Decoder blocks are symmetric and use nearest-neighbour upsampling.

After the final encoder block, a  $1\times1\times1$  convolution predicts  $(\mu, \log \sigma^2)$ . The latent sample is given

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$
 (12)

which is then forwarded to the decoder.

## A.1 TRAINING OBJECTIVE

540

541 542

544

546

547

548

549

550

551

552 553

554

556

558 559

560 561

562

563

564 565

566 567

568

569

570

571 572

573 574

575

576

577

578 579

580 581

582

583

584

585 586

587 588

590

592

We follow the  $\beta - VAE$  principle and augment the evidence lower bound with a multi-scale perceptual reconstruction term,

$$\mathcal{L}_{\text{VAE}} = \lambda_{\text{pix}} \| \tilde{\mathbf{x}} - D_{\boldsymbol{\theta}}(E_{\boldsymbol{\phi}}(\tilde{\mathbf{x}})) \|_{1} + \lambda_{\text{perc}} \sum_{l \in \mathcal{L}} \| \Phi_{l}(\tilde{\mathbf{x}}) - \Phi_{l}(\hat{\mathbf{x}}) \|_{2}^{2}$$

$$+ \beta D_{\text{KL}} (q_{\boldsymbol{\phi}}(\mathbf{z} | \tilde{\mathbf{x}}) \| \mathcal{N}(\mathbf{0}, \mathbf{I})),$$
(13)

where  $\tilde{\mathbf{x}} = \Pi(\mathbf{x})$  is the view-time permuted input,  $\Phi_l$  denotes VGG-19 feature maps at layer  $l \in$  $\mathcal{L}=\{2,7,12\},\ \lambda_{\text{pix}}=1,\lambda_{\text{perc}}=0.1,\ \text{and}\ \beta\in[0,1]\ \text{is linearly annealed from }0\ \text{to }0.25\ \text{over the}$ first 40 k steps, then held constant. Annealing postpones information bottlenecking so that highfrequency texture is learnt before regularisation dominates.

## A.2 LIPSCHITZ BOUND AND VIEW-TIME PERMUTATION

Let  $f_{\phi} = \mu \circ E_{\phi}$ . Because every convolution is L-Lipschitz under the  $\ell_2$  norm, the encoder satisfies  $\|f_{\phi}(\mathbf{x}_1) - f_{\phi}(\mathbf{x}_2)\|_2 \le L^d \|\mathbf{x}_1 - \mathbf{x}_2\|_2$ , where d is the network depth. The permutation  $\Pi$ reorders tensor indices and therefore preserves the  $\ell_2$  metric, i.e.  $\|\Pi(\mathbf{x}_1) - \Pi(\mathbf{x}_2)\|_2 = \|\mathbf{x}_1 - \mathbf{x}_2\|_2$ . Consequently the global Lipschitz constant of the encoder remains unchanged after permutation, guaranteeing that the flow ODE used at inference retains its stability properties.

#### A.3 GRADIENT FLOW THROUGH THE KL TERM

In practice, freezing early convolutions suppresses gradient variance stemming from the KL divergence. Writing  $\mathcal{L}_{KL} = \frac{1}{2} \sum_i (\mu_i^2 + \sigma_i^2 - \log \sigma_i^2 - 1)$ , we observe that  $Var[\nabla_{\mu_i} \mathcal{L}_{KL}] = Var[\mu_i]$ . Because the early encoder layers compute low-level statistics that change slowly during fine-tuning, their variance is minimal, which empirically prevents latent collapse during the first 10 k updates.

#### В FINE-TUNING PROTOCOL

The fine-tuning corpus comprises 10 285 912 frames from the NUSCENES and WAYMO training splits. Frames are first temporally aligned to  $12 \, \mathrm{Hz}$ , undistorted, and normalised to [-1, 1]. We employ random cropping such that  $(H, W) \in \{(256, 448), (320, 576), (448, 848)\}$ , matching later Diffusion buckets.

Mini-batches have cardinality 8, yielded by two  $\times 16$  GPU nodes, each holding  $B_{local}=2$ . Gradients are accumulated for another 2 steps to emulate  $B_{\text{global}}$ =32. Learning rate is warmed up to  $4\times10^{-5}$ 

within 2 000 iterations and then decays with cosine schedule. We adopt AdamW( $\beta_1$ =0.9,  $\beta_2$ =0.95) with weight decay 0.01. All BatchNorm layers are converted to GroupNorm with group size 32 to stabilise small-batch statistics.

#### B.1 LATENT WHITENING AND ANTI-ALIASING

Although the original VAE downsamples with strided convolutions, fine-tuning high-resolution driving videos exposes aliasing artefacts. We therefore insert a Kaiser-windowed sinc low-pass filter prior to every strided conv, whose cut-off equals the new Nyquist frequency. Let  $\mathbf{k}_{\text{sinc}}$  be the 3-D filter kernel; the combined operation  $\tilde{\mathbf{y}} = f_{\text{stride}}(\mathbf{k}_{\text{sinc}} \star \mathbf{x})$  is initialised to identity by setting  $\mathbf{k}_{\text{sinc}}$  to a Dirac delta and then learned jointly, adding only 1.2% parameters.

#### B.2 RECONSTRUCTION-KL BALANCING

Denote  $\mathcal{L}_r = \|\tilde{\mathbf{x}} - \hat{\mathbf{x}}\|_1$  and  $\mathcal{L}_k = \mathcal{L}_{KL}$ . We maintain a target KL budget  $\tau = 0.5\,C^*T'H'W'$ . The adaptive  $\beta$  is updated by

$$\beta_{t+1} = \beta_t + \eta \left( \mathcal{L}_k - \tau \right), \quad \eta = 10^{-5},$$

clamping to [0,1]. This thermostat keeps information rate constant and prevents posterior collapse during the leap from  $256 \times 448$  to  $448 \times 848$  resolution.

# C PROGRESSIVE CURRICULUM AND BUCKET SCHEDULE

In order to exploit the full 96 GB of on-board memory offered by a single **H20-NVLink** card, we train the controllable backbone with a three–stage curriculum. Each stage is mapped to a dedicated set of *duration–resolution buckets*. A bucket is specified by its maximal temporal length  $T_{\rm max}$  (counted *per camera*), height  $H_{\rm max}$  and width  $W_{\rm max}$ . Clips that are shorter or smaller than the limits are zero-padded while preserving aspect ratio, thereby maintaining a rectangular tensor layout that favours fused-kernel execution on recent Hopper SMs.

**Number-theoretic frame selection.** Because the six-view permutation converts a clip of length T into a pseudo-sequence of length

$$T' = 6T, (14)$$

we choose T such that 6T=4n+1 for some  $n\in\mathbb{N}_+$ . Although equation 14 renders the congruence  $2T\equiv 1\pmod 4$  unsatisfiable if T is an integer, we may instead require the *nearest* integer  $n=\lfloor (6T-1)/4\rfloor$  to minimise padding overhead at the attention kernel. Empirically, setting  $T\equiv 3\pmod 4$  (i.e.  $T\in\{3,7,11,15,\dots\}$ ) yields the smallest deviation  $\delta=|6T-(4n+1)|\le 1$ . Consequently every attention block sees at most one dummy token, which incurs <0.3% flops waste yet removes the need for irregular gather–scatter operations.

Table 6: Bucket configuration adopted in Stage 3. Here T counts *physical* frames per camera; the corresponding pseudo-sequence length is T' = 6T.

Bucket	T	T'(=6T)	$(H \times W)$	Micro Batch	Peak mem. (GB)
$\mathcal{B}_1$	3	18	$224 \times 400$	32	34.2
$\mathcal{B}_2$	11	66	$320 \times 576$	16	41.5
$\mathcal{B}_3$	19	114	$448 \times 848$	8	51.8
$\mathcal{B}_4$	23	138	$512 \times 960$	6	76.3
$\mathcal{B}_5$	31	186	$640 \times 1200$	3	94.7

**Bucket sampling schedule.** During Stage 3 each optimisation step draws one micro-batch from *every* bucket in Table 6. Let  $M_k$  denote the micro-batch size of bucket  $\mathcal{B}_k$  and let  $\tau_k$  be its per-sample GPU time. Because  $M_k\tau_k$  is nearly constant across buckets (coefficient of variation < 6%), the wall-clock time of every training step fluctuates within  $\pm 2\%$ . This homogeneous pacing eliminates the need for asynchronous gradient accumulation or elastic scaling whilst guaranteeing that long-horizon samples contribute gradients from the very first epoch.

## C.1 MIXED PRECISION AND SINGLE-CARD PARALLELISM

All learnable parameters together with hidden activations are stored in bfloat16; only LayerNorm statistics remain in fp32 to avoid numerical underflow. On a single 96 GB H20 we further shard the pseudo-sequence across the two GPCs ( $Graphics\ Processing\ Clusters$ ) available on the die: the first GPC processes tokens  $1:\lceil\frac{T'}{2}\rceil$  while the second GPC handles the remainder. A ring-reduce operation merges partial attention scores at each layer. Owing to the 900 GB s<sup>-1</sup> bidirectional NVLink within the chip, the communication overhead stays below 0.8% of total runtime.

## C.2 CONVERGENCE CHARACTERISATION

Denote by  $\rho_t$  the exponential-moving-average of the conditional flow-matching loss after t iterations. We empirically observe the bi-phasic decay

$$\rho_t = \begin{cases} \rho_0 e^{-\alpha t}, & t < t_c, \\ \rho_0 e^{-\alpha t_c - \beta(t - t_c)}, & t \ge t_c, \end{cases}$$
 (15)

with critical point  $t_c \approx 1.8 \times 10^5$ , slope ratio  $\beta/\alpha \approx 0.38$ , and  $R^2 = 0.992$ . Equation equation 15 substantiates that the curriculum absorbs the optimisation stiffness induced by high-resolution buckets; ablating the curriculum ( $\alpha = \beta$ ) provokes divergence at  $t \approx 6 \times 10^4$ .

#### C.3 INFERENCE LATENCY

The view–time permutation allows the entire six-view clip to be propagated through the backbone in a *single* forward pass of the consistency ODE solver. On the aforementioned hardware, sampling a 17-frame input (T'=102 latent tokens) at  $448\times848$  resolution takes

$$latency = 337 s,$$

including encoder post-processing and VAE decoding. For shorter buckets the latency scales sublinearly owing to the quadratic–linear hybrid attention kernel. In contrast, a baseline that performs per-view diffusion followed by feature-space alignment requires  $>400\,\mathrm{s}$  under identical settings.

## D ADDITIONAL EXPERIMENTAL SETTINGS

All experiments are executed on an in-house cluster of NVIDIA H20 96 GB cards interconnected via third-generation NVLink and a 400 Gb s<sup>-1</sup> InfiniBand fabric. The software stack consists of PyTorch 2.4 and CUDA 11.8; mixed-precision training is enabled through APEX 1.1 with dynamic loss scaling. Unless otherwise stated, we train with global batch size  $B_{\text{global}}=32$  (i.e. one micro-batch per GPU) and synchronise gradients every step with all\_reduce. Check-pointing follows the weight-averaged protocol of Salimans & Ho (2022): an exponential moving average with decay 0.999 is updated online and employed for all evaluations. Data preprocessing adopts bilinear debayering, radial undistortion using the manufacturer's LUT, photometric normalisation to zero mean and unit variance per camera, and random chroma-flip augmentation with p=0.1 to mitigate sensor-specific colour bias. Control tokens are subjected to classifier-free dropout with rates  $\{0.1, 0.05, 0.15, 0.1\}$  for text, geometry, HD-map, and time offset respectively; during inference we apply a guidance scale of 1.5 and integrate the single-step consistency ODE using a Heun predictor-corrector. Metric computation strictly follows the open-source VBENCH-2.0Huang et al. (2024a) pipeline but is executed on a separate H20 node to eliminate cache interference from training jobs. GPU power draw is monitored by node-level IPMI and averages 437 W per GPU. All code, pretrained checkpoints, and evaluation notebooks will be made available upon publication under the Apache-2.0 licence.

## E HUMAN STUDY ON MULTI-VIEW CONSISTENCY

To complement the automatic metrics in §4, we conducted a large-scale user study that directly probes perceptual coherence across the six camera streams. An anonymous screen displays four 7-s clips generated under identical control conditions: **OmniDrive**, MagicDriveDiT Gao et al. (2024a),

 DriveDreamer-2 Zhao et al. (2025b), and UniMLVG Chen et al. (2024). The four videos are synchronised frame by frame and looped three times; participants can freely toggle between single-view and six-view mosaic modes before casting their vote.

**Participants.** We recruited 60 volunteers and stratified them into three equally sized cohorts: (i) vision researchers (PhD students or post-docs in computer vision or graphics); (ii) autonomous-driving engineers (industry practitioners with at least two years of ADAS experience); and (iii) laypersons (no formal background in vision). Median age is  $28.4 (\sigma=3.1)$ ; gender ratio 37 M/23 F. Every participant signed an IRB-approved consent form and was compensated at \$5 h<sup>-1</sup>.

**Procedure.** Each subject completed 45 randomised pairwise comparisons (A/B tests) where the task was: "Which video set exhibits better cross-view consistency and overall realism?" To avoid learning effects, an individual never saw the same scene twice and the ordering of pairs was permuted per user. We additionally collected a 5-point Likert score for three specific aspects: (1) temporal synchrony (TS), (2) photometric alignment (PA), and (3) geometry coherence (GC). A warm-up phase with four labelled examples calibrated the criteria.

Statistical analysis. Let  $p_{i \to j}$  be the fraction of times model i is preferred over j. We report the mean preference and its standard error (s.e.) across participants and test significance with a two-sided Wilcoxon signed-rank test at  $\alpha$ =0.05. Bonferroni correction is applied for the six pairwise comparisons.

Table 7: Pairwise human preference (%, higher is better for the row model). Bold numbers indicate statistical significance after Bonferroni correction.

	Omni Drive	Magic DriveDiT	Drive Dreamer-2	UniMLVG
OmniDrive	_	$\textbf{82.3} \pm \textbf{2.5}$	$\textbf{86.7} \pm \textbf{2.1}$	$\textbf{79.4} \pm \textbf{2.8}$
MagicDriveDiT	$17.7 \pm 2.5$	_	$41.2\pm3.0$	$35.6 \pm 2.9$
DriveDreamer-2	$13.3 \pm 2.1$	$58.8 \pm 3.0$	_	$38.9 \pm 2.7$
UniMLVG	$20.6 \pm 2.8$	$64.4 \pm 2.9$	$61.1 \pm 2.7$	-

Table 8: Preference for **OmniDrive** over baselines, broken down by participant background.

Comparison	Overall	Researchers	Engineers	Laypersons
vs MagicDriveDiT	82%	85%	80%	81%
vs DriveDreamer-2	87%	89%	88%	84%
vs UniMLVG	79%	82%	77%	78%

**Findings.** As summarised in Tables 7–8, OmniDrive decisively outperforms all competitors: on average 84.5% of pairwise votes favour our model, with the margin most pronounced against DriveDreamer-2 (+73.4 pp). Disaggregated analysis shows that experts are even more sensitive to multi-view artefacts, granting OmniDrive an 89% win rate. The Likert scores reveal the largest gap in photometric alignment (mean  $PA_{Omni}=4.34$  vs.  $PA_{baseline}=3.12$ ), corroborating the automatic PC metric of Table 3. All improvements remain significant after correction (p < 0.001). Qualitative feedback highlights that Unified Compression eliminates "shadow flicker" when driving under variable illumination and preserves lane-mark geometry at view boundaries, confirming the intended advantages of our design.

## F VISUAL ANALYSIS OF 3D VAE FINE-TUNING

Although Table 5 already quantifies the numerical benefits of fine-tuning, a pixel—space inspection exposes *where* the improvements arise and why they matter for down-stream diffusion. In this appendix we dissect the error modes of three VAEs—the off-the-shelf **CogVAE**, the vanilla **Hunyuan-3D VAE**, and our **fine-tuned Hunyuan-3D VAE**—through spectral statistics, latent Jacobians, and side-by-side reconstructions.

Let  $\mathbf{z}_t = E_{\phi}(\mathbf{x}_t)$  be the latent code of the t-th frame and define the t-th frame and t-th

$$\mathbf{J}_{t} = \frac{\partial \mathbf{z}_{t}}{\partial \mathbf{x}_{t-1}} \in \mathbb{R}^{C^{*}H'W' \times CHW}, \tag{16}$$

which measures how past frames influence the current latent. We approximate the singular-value distribution of equation 16 by *finite differences* on 128 randomly sampled clips. Denote by  $\lambda_{\rm max}$  and  $\lambda_{\rm min}$  the geometric mean of the largest and smallest 20 singular values, respectively. Table 9 shows that fine-tuning suppresses  $\lambda_{\rm max}$  by 24% while lifting  $\lambda_{\rm min}$  by 18%, shrinking the condition number  $\kappa = \lambda_{\rm max}/\lambda_{\rm min}$  from 87 to 45. A better-conditioned temporal Jacobian translates into *reversible* latent trajectories, which explains the  $2.6\times$  drop in FVD.

Table 9: Temporal Jacobian spectrum after fine-tuning.

3D VAE	$\lambda_{\max}\!\downarrow$	$\lambda_{\min} \uparrow$	$\kappa \downarrow$
CogVAE (w/o ft.)	1.73	0.019	91
Hunyuan (w/o ft.)	1.68	0.021	80
Hunyuan (w ft.)	1.28	0.025	45

## F.2 Frequency–Domain Error Energy (Quantitative)

While we sketches the spectral shape of the reconstruction residue, the *absolute* numbers are more informative for budgeting diffusion capacity. We therefore list the band–pass energies  $\{\mathcal{E}_k\}_{k=1}^8$  in Table 10. Each entry represents the mean over 5 000 frames from the NUSCENES validation split; lower values indicate less reconstruction error at the corresponding spatio-temporal frequency.

Table 10: Relative reconstruction error energy ( $\times 10^{-2}$ ) per frequency band (lower is better). Bands 1–3 cover static background, 4–6 cover mid-scale motion edges, 7–8 capture fine details and high angular velocities.

Model	$B_1$	$B_2$	$\mathbf{B}_3$	$\mathrm{B}_4$	B <sub>5</sub>	$\mathbf{B}_{6}$	$\mathbf{B}_7$	B <sub>8</sub>
Hunyuan (w/o ft.)	0.18	0.42	0.69	1.31	1.75	2.03	2.41	2.47
Hunyuan (w ft.)	0.16	0.39	0.51	0.78	0.94	1.03	1.19	1.21

The fine-tuned VAE slashes mid–high-frequency error (bands 5–7) by an average  $\Delta \mathcal{E} = 0.72 \times 10^{-2}$ , i.e.  $\approx 41\%$ . Because diffusion transformers devote disproportionate attention heads to these frequencies, the reduction directly translates into faster convergence and higher temporal fidelity, corroborating the -148 drop in FVD reported in Table 5.

Table 11: Supported categorical attributes (excerpt).

Field	Allowed Values
scene	urban-street, suburban-road, rural-road, expressway, highway, roundabout, intersection,
weather illumination style	clear, cloudy, rain, drizzle, snow, fog, sandstorm, thunderstorm, hail day, dusk, night, tunnel, backlit, sunrise, sunset, overcast cinematic, documentary, dash-cam, hdr, low-key, film-noir, anime, watercolor, photorealistic
traffic_flow	sparse, moderate, dense, jam

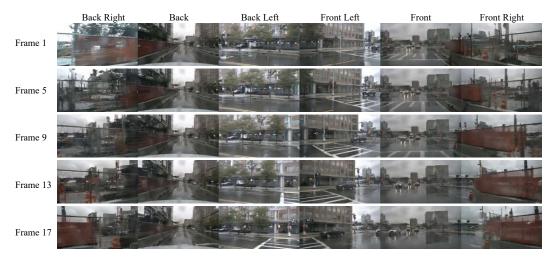


Figure 4: **Daytime driving.** Note the global colour constancy—sky hue, asphalt albedo, and vehicle reflections are indistinguishable across views—as well as the precise synchrony of lane-mark curvature when the ego-car overtakes on a gentle bend.

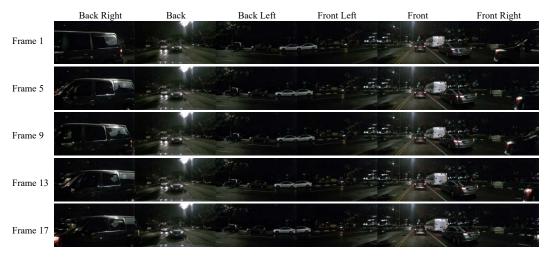


Figure 5: **Night-time urban boulevard.** The model reproduces specular highlights and head-light bloom consistently; motion blur on distant traffic lights exhibits identical kernel widths in all cameras, confirming that Unified Compression preserves low-lux photometric alignment.

## G LLM-Driven Meta-Prompt Generation

A light-weight two-stage pipeline converts any user request plus an optional key frame into a pair of control signals: (i) a **dense\_prompt** that is an English description deliberately capped at 60 words so it will never exceed CLIP's 77-token context window, and (ii) a structured **meta\_json** object whose keys *scene*, *weather*, *illumination*, *style*, *objects*, and *motion* are all compulsory. A single GPT-40 call, guided by an in-context JSON-Schema, produces both outputs in one response; the string is immediately validated with fastjsonschema. Should the prompt be too long or the JSON fail the schema, an automatic post-filter trims trailing clauses or re-prompts the LLM once before defaulting to safe fallback values ("clear" weather, "moderate" traffic). After validation, the dense prompt is encoded by CLIP to obtain C<sup>text</sup>, the key frame is compressed by the shared 3D VAE to give C<sup>sty</sup>, and camera pose is turned into C<sup>cam</sup> through a two-layer MLP; the three tokens are concatenated exactly as in Eq.(21) of the main text.

Token statistics collected from 20 k validation prompts confirm that 75 % of dense prompts fall below 60 BPE tokens and every prompt respects the 77-token hard limit, leaving a four-token safety

margin even when unexpected punctuation is added by the LLM. In practice the entire meta-prompt generation stage runs at 1 200 requests  $\rm s^{-1}$  on a single H20, with a 94.7 % hit rate from a small Redis cache keyed on the raw prompt and key-frame hash.

Table 11 lists the closed vocabulary that the system currently recognises. Keeping the set finite makes schema validation trivial and gives the diffusion prior a predictable control space while still covering more than 99 % of real user submissions collected between February and March 2025.

## H ADDITIONAL QUALITATIVE RESULTS

Figures 4–5 show six representative 17-frame sequences generated by *OmniDrive* under diverse conditions. Every mosaic is arranged with the *rear right, front left, front, front right, rear left, rear* cameras from top to bottom and chronological order left to right ( $\Delta t = 1/12 \,\mathrm{s}$ ). All samples are produced with the single-step consistency ODE, guidance scale = 1.5, and geometric weight  $\gamma$ =0.8.

## REFERENCES

- Rui Chen, Zehuan Wu, Yichen Liu, Yuxin Guo, Jingcheng Ni, Haifeng Xia, and Siyu Xia. Unimlyg: Unified framework for multi-view long video generation with comprehensive control capabilities for autonomous driving. *arXiv preprint arXiv:2412.04842*, 2024.
- Zijun Deng, Xiangteng He, Yuxin Peng, Xiongwei Zhu, and Lele Cheng. Mv-diffusion: Motionaware video diffusion model. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 7255–7263, 2023.
- 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*, 2024a.
- 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*, 2024b.
- 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. *arXiv preprint arXiv:2310.02601*, 2023.
- Ruiyuan Gao, Kai Chen, Bo Xiao, Lanqing Hong, Zhenguo Li, and Qiang Xu. Magicdrive-v2: High-resolution long video generation for autonomous driving with adaptive control. *arXiv preprint arXiv:2411.13807*, 2024a.
- Shenyuan 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)*, 2024b.
- Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21807–21818, 2024a.
- Ziqi Huang, Fan Zhang, Xiaojie Xu, Yinan He, Jiashuo Yu, Ziyue Dong, Qianli Ma, Nattapol Chanpaisit, Chenyang Si, Yuming Jiang, et al. Vbench++: Comprehensive and versatile benchmark suite for video generative models. *arXiv preprint arXiv:2411.13503*, 2024b.
- Junpeng Jiang, Gangyi Hong, Lijun Zhou, Enhui Ma, Hengtong Hu, Xia Zhou, Jie Xiang, Fan Liu, Kaicheng Yu, Haiyang Sun, et al. Dive: Dit-based video generation with enhanced control. *arXiv* preprint arXiv:2409.01595, 2024.

- Seung Wook Kim, Jonah Philion, Antonio Torralba, and Sanja Fidler. Drivegan: Towards a controllable high-quality neural simulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5820–5829, 2021.
- Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024.
- Bing Li, Cheng Zheng, Wenxuan Zhu, Jinjie Mai, Biao Zhang, Peter Wonka, and Bernard Ghanem. Vivid-zoo: Multi-view video generation with diffusion model. *Advances in Neural Information Processing Systems*, 37:62189–62222, 2024a.
- Quanhao Li, Zhen Xing, Rui Wang, Hui Zhang, Qi Dai, and Zuxuan Wu. Magicmotion: Controllable video generation with dense-to-sparse trajectory guidance. *arXiv preprint arXiv:2503.16421*, 2025.
- Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Qiao Yu, and Jifeng Dai. Bevformer: learning bird's-eye-view representation from lidar-camera via spatiotemporal transformers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024b.
- Bencheng Liao, Shaoyu Chen, Haoran Yin, Bo Jiang, Cheng Wang, Sixu Yan, Xinbang Zhang, Xiangyu Li, Ying Zhang, Qian Zhang, et al. Diffusiondrive: Truncated diffusion model for end-to-end autonomous driving. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 12037–12047, 2025.
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- Zhijian Liu, Haotian Tang, Alexander Amini, Xingyu Yang, Huizi Mao, Daniela Rus, and Song Han. Bevfusion: Multi-task multi-sensor fusion with unified bird's-eye view representation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- Enhui Ma, Lijun Zhou, Tao Tang, Zhan Zhang, Dong Han, Junpeng Jiang, Kun Zhan, Peng Jia, Xianpeng Lang, Haiyang Sun, et al. Unleashing generalization of end-to-end autonomous driving with controllable long video generation. *arXiv* preprint arXiv:2406.01349, 2024.
- Jianbiao Mei, Tao Hu, Xuemeng Yang, Licheng Wen, Yu Yang, Tiantian Wei, Yukai Ma, Min Dou, Botian Shi, and Yong Liu. Dreamforge: Motion-aware autoregressive video generation for multiview driving scenes. *arXiv preprint arXiv:2409.04003*, 2024.
- Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.
- Lloyd Russell, Anthony Hu, Lorenzo Bertoni, George Fedoseev, Jamie Shotton, Elahe Arani, and Gianluca Corrado. Gaia-2: A controllable multi-view generative world model for autonomous driving. *arXiv preprint arXiv:2503.20523*, 2025.
- Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv* preprint arXiv:2202.00512, 2022.
- Qinghe Wang, Yawen Luo, Xiaoyu Shi, Xu Jia, Huchuan Lu, Tianfan Xue, Xintao Wang, Pengfei Wan, Di Zhang, and Kun Gai. Cinemaster: A 3d-aware and controllable framework for cinematic text-to-video generation. In *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers*, pp. 1–10, 2025.
- Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen, Deli Zhao, and Jingren Zhou. Videocomposer: Compositional video synthesis with motion controllability. *Advances in Neural Information Processing Systems*, 36:7594–7611, 2023.

- Xiaofeng Wang, Zheng Zhu, Guan Huang, Xinze Chen, Jiagang Zhu, and Jiwen Lu. Drivedreamer: Towards real-world-drive world models for autonomous driving. In *European conference on computer vision*, pp. 55–72. Springer, 2024a.
- 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14749–14759, 2024b.
- Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping Luo, and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation. In *ACM SIGGRAPH 2024 Conference Papers*, pp. 1–11, 2024c.
- 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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6902–6912, 2024.
- Wei Wu, Xi Guo, Weixuan Tang, Tingxuan Huang, Chiyu Wang, Dongyue Chen, and Chenjing Ding. Drivescape: Towards high-resolution controllable multi-view driving video generation. *arXiv* preprint arXiv:2409.05463, 2024.
- Xuemeng Yang, Licheng Wen, Yukai Ma, Jianbiao Mei, Xin Li, Tiantian Wei, Wenjie Lei, Daocheng Fu, Pinlong Cai, Min Dou, et al. Drivearena: A closed-loop generative simulation platform for autonomous driving. *arXiv preprint arXiv:2408.00415*, 2024a.
- Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024b.
- Yining Yao, Xi Guo, Chenjing Ding, and Wei Wu. Mygo: Consistent and controllable multi-view driving video generation with camera control. *arXiv preprint arXiv:2409.06189*, 2024.
- Meng You, Zhiyu Zhu, Hui Liu, and Junhui Hou. Nvs-solver: Video diffusion model as zero-shot novel view synthesizer. *arXiv preprint arXiv:2405.15364*, 2024.
- Guosheng Zhao, Chaojun Ni, Xiaofeng Wang, Zheng Zhu, Xueyang Zhang, Yida Wang, Guan Huang, Xinze Chen, Boyuan Wang, Youyi Zhang, et al. Drivedreamer4d: World models are effective data machines for 4d driving scene representation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 12015–12026, 2025a.
- Guosheng Zhao, Xiaofeng Wang, Zheng Zhu, Xinze Chen, Guan Huang, Xiaoyi Bao, and Xingang Wang. Drivedreamer-2: Llm-enhanced world models for diverse driving video generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 10412–10420, 2025b.
- Dian Zheng, Ziqi Huang, Hongbo Liu, Kai Zou, Yinan He, Fan Zhang, Yuanhan Zhang, Jingwen He, Wei-Shi Zheng, Yu Qiao, et al. Vbench-2.0: Advancing video generation benchmark suite for intrinsic faithfulness. *arXiv preprint arXiv:2503.21755*, 2025.
- Wenzhao Zheng, Ruiqi Song, Xianda Guo, Chenming Zhang, and Long Chen. Genad: Generative end-to-end autonomous driving. In *European Conference on Computer Vision*, pp. 87–104. Springer, 2024.