

Risk-Controllable Multi-View Diffusion for Driving Scenario Generation

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Generating safety-critical driving scenarios is crucial for*
002 *evaluating and improving autonomous driving systems, but*
003 *long-tail risky situations are rarely observed in real-world*
004 *data and difficult to specify through manual scenario de-*
005 *sign. Existing generative approaches typically treat risk*
006 *as an after-the-fact label and struggle to maintain ge-*
007 *ometric consistency in multi-view driving scenes. We*
008 *present RiskMV-DPO, a general and systematic pipeline for*
009 *physically-informed, risk-controllable multi-view scenario*
010 *generation. By integrating target risk levels with physically-*
011 *grounded risk modeling, we autonomously synthesize di-*
012 *verse and high-stakes dynamic trajectories that serve as ex-*
013 *PLICIT geometric anchors for a diffusion-based video gen-*
014 *erator. To ensure spatial-temporal coherence and geomet-*
015 *ric fidelity, we introduce a geometry-appearance alignment*
016 *module and a region-aware direct preference optimization*
017 *(RA-DPO) strategy with motion-aware masking to focus*
018 *learning on localized dynamic regions. Experiments on the*
019 *nuScenes dataset show that RiskMV-DPO can freely gen-*
020 *erate a wide spectrum of diverse long-tail scenarios while*
021 *maintaining state-of-the-art visual quality, improving 3D*
022 *detection mAP from 18.17 to 30.50 and reducing FID to*
023 *15.70. Our work shifts the role of world models from pas-*
024 *sive environment prediction to proactive, risk-controllable*
025 *synthesis, providing a scalable toolchain for the safety-*
026 *oriented development of embodied intelligence.*

027 1. Introduction

028 Autonomous driving systems are ultimately judged by how
029 safely and reliably they behave in the open world, includ-
030 ing rare but safety-critical events [5, 10, 39]. Since real-
031 world logs under-sample these cases, driving scenario gen-
032 eration has become essential for data augmentation, stress
033 testing, and systematic safety evaluation [17, 22, 32]. With
034 production vehicles relying heavily on camera suites, gen-
035 erating temporally coherent, multi-view consistent driving
036 sequences is particularly valuable [19, 41].

037 Despite rapid progress, generating long-tail risky sce-

narios remains a significant bottleneck [11]. A common
practice is to generate a large pool of driving scenes and
then mine rare events, but truly high-risk situations occur
with extremely low probability [40]. Another practice is to
handcraft scenarios by specifying parameters such as near-
collisions or road intrusions [37]. While offering direct con-
trol, these "sampling-and-filtering" or manual design ap-
proaches inevitably reflect human priors and lack the effi-
ciency to synthesize the vast, non-intuitive space of haz-
ardous interactions. What is missing is a general pipeline to
transform "risk" from an after-the-fact label into a proac-
tive, time-resolved control signal that can autonomously
guide the synthesis of diverse long-tail hazards.

At the same time, the multi-view, outdoor, dynamic
nature of driving scenes introduces a second bottleneck:
geometry and robustness [2]. While diffusion models
have advanced visual realism, driving applications demand
more. They require cross-view metric-level consistency,
motion continuity, and stability under occlusion and light-
ing changes. For consistency, many attempts rely on large-
scale question-answer datasets to teach VLMs spatial con-
cepts, but this supervision rarely provides a solid geomet-
ric prior that enforces physically grounded alignment across
cameras [24, 30]. This motivates mechanisms that inject 3D
priors into diffusion training and align the model's capacity
with localized, motion-dominant regions that matter most
for driving risk.

In short, these challenges are twofold: *making risk*
controllable for long-tail scenario generation, and *making*
multi-view geometry reliable under outdoor dynamic con-
ditions. Therefore, we propose risk-controllable multi-view
diffusion for driving scenario generation, a framework that
makes risk a first-class control knob while improving geo-
metric fidelity and localized dynamic realism. Our key idea
is to turn risk into a time-resolved conditioning signal, gen-
erate trajectories and 3D boxes at a target risk level, and use
them as structured controls for multi-view diffusion. The
contributions can be summarized as follows:

- We propose a general pipeline for risk-controllable syn-
thesis that reformulates driving risk as a proactive control
signal, enabling the autonomous generation of dynamic

- 079 trajectories and 3D boxes at target risk levels.
 080 • We introduce a geometry-appearance alignment module
 081 that injects compact 3D priors into multi-view diffusion to
 082 ensure metric-level cross-view consistency and geometric
 083 plausibility.
 084 • We develop multi-view consistent RA-DPO with motion-
 085 aware masking to align scene synthesis with localized
 086 preferences in dynamic regions, significantly enhancing
 087 temporal and motion realism.

088 2. Related Work

089 This section reviews driving scenario generation, risk-aware
 090 modeling and simulation, geometric priors for multi-view
 091 diffusion, and preference alignment for generative models.

092 2.1. Driving scene and multi-view video generation

093 Driving scenario generation has evolved from neural sim-
 094 ulators to large-scale driving world models [18]. Drive-
 095 GAN [14] learns an end-to-end simulator conditioned on
 096 ego actions. GAIA-1 [12] models driving as token se-
 097 quences with video/text/action conditioning. Diffusion-
 098 based world models further improve realism and controlla-
 099 bility, such as DriveDreamer [31] and DriveDreamer-2 [43].

100 Multi-view driving video synthesis introduces explicit
 101 structure controls to improve cross-view consistency. Mag-
 102 icDrive [7] conditions on camera poses, maps, and 3D
 103 boxes with cross-view attention. DrivingDiffusion [16]
 104 generates multi-view videos from 3D layouts and enforces
 105 cross-view and cross-frame consistency. MagicDrive-
 106 V2 [9] scales to longer and higher-resolution synthesis, and
 107 recent work explores stronger spatiotemporal-view atten-
 108 tion [21]. Despite these advances, long-tail risk is still hard
 109 to target explicitly and multi-view geometry remains fragile
 110 in dynamic outdoor scenes, motivating our risk-controllable
 111 and geometry-aware design.

112 2.2. Risk modeling and risk-aware simulation

113 Risk modeling often relies on trajectory-based criticality
 114 metrics (e.g., TTC variants) [34, 38] and formal safety en-
 115 velopes such as RSS [25]. Recent work moves toward risk
 116 forecasting under interaction and uncertainty. RiskNet [20]
 117 combines field-theoretic interactions with multi-modal tra-
 118 jectory prediction for time-evolving risk estimation, and
 119 SafeDrive [45] quantifies perceived risk via risk fields to
 120 guide an LLM driving agent.

121 Risk-aware testing uses risk signals to search for rare
 122 failures efficiently. Scenic [6] specifies scenarios as distri-
 123 butions with constraints, VerifAI [4] supports simulation-
 124 based falsification and parameter synthesis, and Adaptive
 125 Stress Testing [15] uses RL to discover failure trajectories.
 126 Most toolchains still operationalize risk via thresholds or
 127 specification violations rather than enabling generation at a

128 desired, time-resolved risk level, which motivates our risk-
 129 as-control formulation.

2.3. 3D priors for multi-view generation 130

131 Multi-view diffusion increasingly incorporates explicit 3D
 132 signals to stabilize generation and reduce cross-view
 133 drift [27]. In driving, MagicDrive [7] and DrivingDiffu-
 134 sion [16] condition on 3D controls and add cross-view mod-
 135 eling, while DriveScape [35] scales to high-resolution syn-
 136 thesis via multi-view feature fusion with 3D guidance. Even
 137 so, models can “understand the language” but still can-
 138 not “see through the geometry”: existing conditioning of-
 139 ten provides surface-level reinforcement and remains brit-
 140 tle under occlusion and fast motion. This motivates our
 141 geometry-appearance alignment that injects compact 3D
 142 priors into diffusion training.

2.4. DPO for generative models 143

144 DPO is a practical tool for aligning generative mod-
 145 els with preference data. Diffusion-DPO [28] adapts
 146 DPO to diffusion models, and DSPO [46] improves ob-
 147 jective matching under score-based training. For video,
 148 DenseDPO [36] motivates finer temporal preference sig-
 149 nals, and LocalDPO [13] introduces region-level mask-
 150 guided preference optimization. Existing preference align-
 151 ment is rarely designed for multi-view driving, where pref-
 152 erences should be consistent across cameras and concen-
 153 trated on motion-dominant regions, motivating our multi-
 154 view motion-aware masking and region-aware DPO with
 155 progressive corruption fusion.

3. Problem Formulation 156

157 We consider the task of generating multi-view driving sce-
 158 narios at a user-specified risk level. Let $\mathcal{V} = \{1, \dots, V\}$
 159 denote the camera set and $\mathbf{I}_{1:T} = \{I_{1:T}^v\}_{v \in \mathcal{V}}$ denote the ob-
 160 served multi-view history. Given structured context \mathbf{C} (in-
 161 cluding HD maps \mathbf{M} and text \mathbf{y}), the user provides a target
 162 risk control \mathbf{r}^* , which can be a scalar r^* or a time-resolved
 163 profile $\mathbf{r}_{1:H}^*$ over a future horizon H .

164 Unlike passive prediction models, our framework reformu-
 165 lates driving risk as a proactive control signal used for
 166 structured conditioning. Let \mathcal{A} denote the set of traffic par-
 167 ticipants. For each agent $a \in \mathcal{A}$, we define the motion con-
 168 trol signal \mathbf{U} as the set of future trajectories and 3D bound-
 169 ing boxes:

$$\mathbf{U} = \left\{ \left(\mathbf{x}_{1:H}^a, \mathbf{b}_{1:H}^a \right) \right\}_{a \in \mathcal{A}}. \quad (1) \quad 170$$

171 To achieve risk-controllable synthesis, we decouple the pro-
 172 cess into two stages. First, a risk control module $g(\cdot)$ gener-
 173 ates physically-grounded motion \mathbf{U} conditioned on the tar-
 174 get risk:

$$\mathbf{U} = g(\mathbf{I}_{1:T}, \mathbf{M}, \mathbf{y}, \mathbf{r}^*). \quad (2) \quad 175$$

176 Second, a multi-view diffusion generator $f(\cdot)$ renders the
177 future frames $\hat{\mathbf{I}}_{T+1:T+H}$ conditioned on the structured motion
178 control:

$$179 \quad \hat{\mathbf{I}}_{T+1:T+H} = f(\mathbf{I}_{1:T}, \mathbf{M}, \mathbf{y}, \mathbf{U}). \quad (3)$$

180 Our objective is to generate diverse long-tail scenarios
181 where the induced risk matches \mathbf{r}^* while preserving cross-
182 view consistency and temporal plausibility.

183 4. Risk-Controlled Motion Generation

184 This section presents how we model driving risk, reveals
185 the risk distribution in typical scenarios, and shows how to
186 control the risk level to generate corresponding trajectories
187 and 3D boxes.

188 4.1. Per-frame Risk Computation

189 At each time step t , we assume the ego vehicle state
190 $(\mathbf{p}_e^t, \mathbf{v}_e^t)$ and a set of surrounding agents $\mathcal{A}_t = \{1, \dots, N_t\}$
191 with states $(\mathbf{p}_i^t, \mathbf{v}_i^t)$ are available from logs or annotations.
192 We compute a per-agent risk contribution R_i^t and optionally
193 aggregate them into a per-frame scalar risk r_t . Our formu-
194 lation follows the field-inspired, direction-aware risk mod-
195 eling principle in RiskNet [20], and adapts it to an efficient
196 per-frame computation.

197 Let the relative displacement and its unit direction be

$$198 \quad \mathbf{r}_i^t = \mathbf{p}_i^t - \mathbf{p}_e^t, \quad \hat{\mathbf{r}}_i^t = \frac{\mathbf{r}_i^t}{\|\mathbf{r}_i^t\|_2 + \epsilon}. \quad (4)$$

199 We use two dot products to indicate whether ego and agent
200 i are moving towards each other:

$$201 \quad d_{e \rightarrow i}^t = (\mathbf{v}_e^t)^\top \mathbf{r}_i^t, \quad d_{i \rightarrow e}^t = (\mathbf{v}_i^t)^\top (-\mathbf{r}_i^t). \quad (5)$$

202 A larger positive $d_{e \rightarrow i}^t$ implies ego is moving towards agent
203 i , and a larger positive $d_{i \rightarrow e}^t$ implies agent i is moving to-
204 wards ego.

205 Moreover, we map the approaching cues to a discrete
206 interaction weight ω_i^t that reflects asymmetric driving risk:

$$207 \quad \omega_i^t = \begin{cases} \omega_{\text{bi}}, & d_{e \rightarrow i}^t > 0 \wedge d_{i \rightarrow e}^t > 0, \\ \omega_{\text{agent}}, & d_{e \rightarrow i}^t \leq 0 \wedge d_{i \rightarrow e}^t > 0, \\ \omega_{\text{ego}}, & d_{e \rightarrow i}^t > 0 \wedge d_{i \rightarrow e}^t \leq 0, \\ \omega_{\text{away}}, & \text{otherwise.} \end{cases} \quad (6)$$

208 Typically ω_{bi} is the largest since both parties are closing,
209 while $\omega_{\text{agent}} > \omega_{\text{ego}}$ reflects that an approaching agent can
210 be riskier when ego may not react in time.

211 Also, we introduce an agent-type coefficient μ_i to reflect
212 severity differences across categories (e.g., heavy trucks vs.
213 pedestrians):

$$214 \quad \mu_i = \text{TypeCoeff}(\text{cls}_i), \quad (7)$$

where cls_i is the category label.

To calibrate longitudinal and lateral factors, we define a
closing speed along the line of sight

$$s_i^t = \max\left(0, (\mathbf{v}_e^t - \mathbf{v}_i^t)^\top \hat{\mathbf{r}}_i^t\right), \quad (8)$$

and convert it to a longitudinal amplification factor

$$\alpha_i^t = \exp(\kappa s_i^t), \quad (9)$$

where κ controls sensitivity. To penalize lateral (non-
colliding) motion, we compute the sine term using the rela-
tive velocity $\mathbf{v}_{\text{rel}}^t = \mathbf{v}_e^t - \mathbf{v}_i^t$:

$$\sin^2 \theta_i^t = \frac{\|\mathbf{v}_{\text{rel}}^t \times \hat{\mathbf{r}}_i^t\|_2^2}{\|\mathbf{v}_{\text{rel}}^t\|_2^2 + \epsilon}, \quad \beta_i^t = \exp(-\lambda \sin^2 \theta_i^t), \quad (10)$$

where λ controls the lateral attenuation.

Therefore, the final per-agent risk contribution is

$$R_i^t = K C \cdot \frac{\omega_i^t \mu_i \alpha_i^t \beta_i^t}{\|\mathbf{r}_i^t\|_2 + \epsilon}, \quad (11)$$

where K and C are calibration constants. We optionally
aggregate per-agent risks into a per-frame scalar risk

$$r_t = \sum_{i \in \mathcal{A}_t} R_i^t \quad \text{or} \quad r_t = \max_{i \in \mathcal{A}_t} R_i^t, \quad (12)$$

depending on whether we emphasize cumulative risk or the
most critical interaction.

4.2. Risk Mining in Typical Scenarios

Driving risk is often not visually obvious. Many safety-
critical situations arise from subtle conflicts among the ego
vehicle, surrounding traffic participants, and road geome-
try, where small timing differences, limited sight distance,
or tight lateral clearance can quickly increase crash prob-
ability. To reveal such “hidden” hazards, we quantify risk
per frame and use the resulting signals to analyze when and
why risk accumulates in typical driving scenarios.

Large-scale crash statistics further motivate this analysis.
NHTSA introduced a pre-crash scenario typology that cate-
gorizes hazardous situations into 37 distinct scenarios using
the General Estimates System (GES) crash database [23].
A later national analysis reorganized these into nine major
groups (rear-end, crossing paths, lane change, road depart-
ure, control loss, opposite direction, animal, pedestrian,
and pedalcyclist) using 2011–2015 FARS and GES data
[26]. Among police-reported crashes where the light vehi-
cle performs the critical action, rear-end and crossing-path
conflicts dominate crash exposure (about 34% and 23%),
followed by lane-change conflicts (about 13%) and road-
departure events (about 11%). Opposite-direction crashes
occur less frequently (about 2%) but show a much higher

256 fatal-to-all crash ratio (about 3.23%), indicating that rare
 257 scenarios can be disproportionately severe. These findings
 258 suggest that risk-aware generation should target not only
 259 frequent crash types but also rare, high-severity scenarios
 260 often overlooked in manual scenario design.

261 We mine latent risks in diverse traffic scenarios by com-
 262 puting a per-frame risk signal and analyzing when risk ac-
 263 cumulates over time. Figure 1 presents three representative
 264 left-turn cases. In Fig. 1(a), the ego vehicle performs an un-
 265 protected left turn at an unsignalized intersection while an
 266 oncoming vehicle proceeds straight, creating a lateral col-
 267 lision threat. The risk rises sharply around the mid-turn
 268 phase when the ego vehicle becomes laterally exposed to
 269 the oncoming stream. In Fig. 1(b), two vehicles execute
 270 parallel left turns with small lateral separation, where slight
 271 asynchrony in turning trajectories can induce a sideswipe
 272 risk that is hard to anticipate from high-level intent alone.
 273 In Fig. 1(c), a parking-lot exit appears immediately after
 274 a left turn, forming a localized blind spot, and the risk spikes
 275 as the ego vehicle approaches the exit region where cross-
 276 traffic may emerge. These examples show that non-intuitive
 277 risks are concentrated in specific moments and regions, moti-
 278 vating the use of time-resolved risk as a controllable signal
 279 for long-tail motion synthesis and scenario generation.

280 4.3. Controllable Motion Generation

281 We treat risk as a conditioning variable for motion synthe-
 282 sis. Given a target risk \mathbf{r}^* (either a scalar r^* or a pro-
 283 file $\mathbf{r}_{1:H}^*$), our goal is to generate future motion $\mathbf{U} =$
 284 $\{(\mathbf{x}_{1:H}^a, \mathbf{b}_{1:H}^a)\}_{a \in \mathcal{A}}$ such that the induced risk matches the
 285 target. We build on a GENAD-style conditional motion
 286 generator [44] and introduce an explicit risk-conditioning
 287 path.

288 We then encode the target risk with an MLP $h(\cdot)$ and
 289 apply feature-wise linear modulation (FiLM) on the BEV
 290 (or latent) feature map \mathbf{F} :

$$291 \mathbf{e}_r = h(\mathbf{r}^*), \quad (\gamma_r, \beta_r) = \psi(\mathbf{e}_r), \quad \tilde{\mathbf{F}} = \gamma_r \odot \mathbf{F} + \beta_r, \quad (13)$$

292 where $\psi(\cdot)$ is a linear head and \odot denotes element-wise
 293 multiplication. The modulated feature $\tilde{\mathbf{F}}$ is then used to
 294 generate multi-modal future trajectories and 3D boxes.

295 Let the generator produce M motion hypotheses
 296 $\{\hat{\mathbf{U}}^{(m)}\}_{m=1}^M$. We compute a scalar risk for each mode us-
 297 ing the per-frame risk functional in Eqs. (11)–(12), denoted
 298 by $\mathcal{R}(\hat{\mathbf{U}}^{(m)})$. We enforce a lower-bound matching and an
 299 upper-bound capability constraint:

$$300 R_{\min} = \min_m \mathcal{R}(\hat{\mathbf{U}}^{(m)}), \quad R_{\max} = \max_m \mathcal{R}(\hat{\mathbf{U}}^{(m)}),$$

$$301 \mathcal{L}_{\min} = |R_{\min} - r^*|, \quad \mathcal{L}_{\max} = |R_{\max} - \tau r^*|.$$

(14)

302 where τ is a scaling factor (e.g., $\tau = 2$) that encourages the
 model to retain the capacity to generate higher-risk modes.

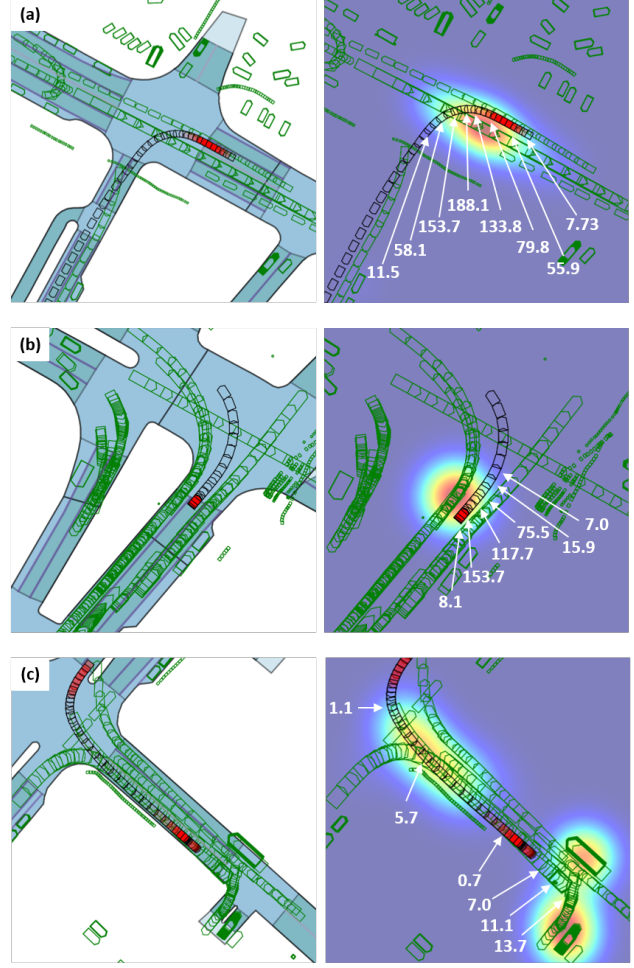


Figure 1. Examples of mined potential risks in typical left-turn scenarios using the proposed per-frame risk quantification. The numbers indicate the risk coefficient of the ego vehicle at that location. (a) Unprotected left turn at an unsignalized intersection with an oncoming straight-moving vehicle. (b) Parallel left turns with tight lateral clearance. (c) Left turn followed by a nearby parking-lot exit that creates a local blind spot.

The final training loss combines the base GENAD objective $\mathcal{L}_{\text{GENAD}}$ with risk guidance:

$$\mathcal{L}_{\text{risk}} = \mathcal{L}_{\text{GENAD}} + \lambda_{\min} \mathcal{L}_{\min} + \lambda_{\max} \mathcal{L}_{\max}. \quad (15)$$

Since r_t is computed per frame, we can localize when risk emerges within a scenario by inspecting peaks or sustained segments of $\{r_t\}_{t=1}^H$. This provides interpretable evidence of non-intuitive risk patterns and supports targeted conditioning of motion generation at a desired risk level.

Therefore, our GENAD-based module produces risk-conditioned trajectories and 3D boxes for ego and other agents, which are used as the primary motion controls for the subsequent multi-view diffusion model to generate driving videos consistent with the specified risk level.

316 5. Multi-View Diffusion for Driving Scenarios

317 Figure 2 illustrates the architecture of the proposed frame-
318 work, termed *RiskMV-DPO*, which integrates diffusion-
319 based generation with multi-view localized preference
320 alignment for risk-aware driving scenario synthesis. Given
321 risk-conditioned motion controls, including trajectories and
322 3D bounding boxes, the model generates temporally coher-
323 ent multi-view driving videos. A multimodal encoder fuses
324 view tokens, motion signals, and scenario context, while
325 geometry-aware conditioning and localized DPO guide
326 the generation toward spatially consistent and dynamically
327 plausible driving scenarios.

328 5.1. Geometry-Appearance Alignment

329 Multi-view diffusion models can produce visually plausi-
330 ble frames, but they often lack a strong geometric prior and
331 may break metric-level cross-view consistency. We inject
332 compact 3D priors by extracting geometry-aware features
333 from DGGT and its VGGT backbone [3, 29], and aligning
334 them with the diffusion model’s appearance features during
335 training.

336 **Geometric feature extraction.** Given the first-frame
337 multi-view images $I \in \mathbb{R}^{B \times N_c \times C \times H \times W}$, where B is the
338 batch size, N_c is the number of cameras, and C is the num-
339 ber of channels, we extract patch-wise VGGT features

$$340 F^{\text{VGGT}} = \Phi_{\text{VGGT}}(I), \quad F^{\text{VGGT}} \in \mathbb{R}^{(BN_c) \times P \times D_{\text{VGGT}}}, \quad (16)$$

341 where P is the number of patches and D_{VGGT} is the VGGT
342 feature dimension (e.g., 3072). A learnable projection maps
343 VGGT features to the diffusion latent dimension D (e.g.,
344 1152):

$$345 F = \text{MLP}_{\text{proj}}(F^{\text{VGGT}}), \quad F \in \mathbb{R}^{(BN_c) \times P \times D}. \quad (17)$$

346 **Learnable geometric token compression.** To obtain a
347 fixed number of compact geometric tokens, we introduce
348 learnable queries $Q \in \mathbb{R}^{N_{\text{tok}} \times D}$ (e.g., $N_{\text{tok}}=16$) and com-
349 press F by cross-attention:

$$350 G = \text{Attn}(Q, F, F), \quad G \in \mathbb{R}^{(BN_c) \times N_{\text{tok}} \times D}. \quad (18)$$

351 For classifier-free guidance, we apply conditional dropout
352 by replacing G with a learned null token set G_{\emptyset} using a
353 Bernoulli mask $m_{\text{drop}} \sim \text{Bernoulli}(1 - p_{\text{drop}})$:

$$354 \tilde{G} = m_{\text{drop}} G + (1 - m_{\text{drop}}) G_{\emptyset}. \quad (19)$$

355 **Appearance feature extraction.** During diffusion train-
356 ing, we extract intermediate token features from the last K
357 backbone layers (e.g., $K=8$). Let $R^{(\ell)} \in \mathbb{R}^{(BN_c) \times N_s \times D}$

be the token features from layer ℓ , where N_s is the num-
ber of spatial tokens. We spatially pool tokens within each
selected layer and average across layers:

$$r^{(\ell)} = \frac{1}{N_s} \sum_{i=1}^{N_s} R_{\cdot, i, \cdot}^{(\ell)}, \quad r = \frac{1}{K} \sum_{\ell} r^{(\ell)}, \quad (20)$$

where $r \in \mathbb{R}^{(BN_c) \times D}$ is the final appearance feature.

Alignment loss. We pool geometric tokens into a single
vector $g \in \mathbb{R}^{(BN_c) \times D}$, normalize both g and r , and align
them with cosine similarity:

$$366 g = \frac{1}{N_{\text{tok}}} \sum_{k=1}^{N_{\text{tok}}} \tilde{G}_{\cdot, k, \cdot}, \quad \hat{g} = \frac{g}{\|g\|_2}, \quad \hat{r} = \frac{r}{\|r\|_2}, \quad (21)$$

$$368 \mathcal{L}_{\text{align}} = 1 - \frac{1}{BN_c} \sum_{i=1}^{BN_c} \hat{g}_i^{\top} \hat{r}_i. \quad (22)$$

This alignment injects compact DGGT/VGGT-derived ge-
ometric priors into diffusion training and improves cross-
view consistency.

372 5.2. Multi-view Consistent Motion-Aware Masking

373 In driving videos, most pixels belong to static background,
374 while safety-critical cues are concentrated in dynamic re-
375 gions such as moving agents and their interactions. Uni-
376 form corruption wastes model capacity and weakens learn-
377 ing on motion-dominant areas. We therefore construct a
378 motion-aware mask and enforce multi-view consistency so
379 that corruption focuses on dynamic regions while remaining
380 geometrically aligned across cameras.

381 **Motion-aware mask.** Let $x \in \mathbb{R}^{B \times N_c \times C \times T \times H \times W}$ de-
382 note the video latent (or feature) tensor, where B is batch
383 size, N_c is the number of cameras, C is the channel dimen-
384 sion, and T is the number of frames. We estimate per-frame
385 motion magnitude by temporal differencing and channel av-
386 eraging:

$$387 D_{b,c,t} = |x_{b,c,\cdot,t+1} - x_{b,c,\cdot,t}|, \quad (23)$$

$$M_{b,c,t}^{\text{mot}} = \frac{1}{C} \sum_{j=1}^C D_{b,c,t,j}.$$

We normalize M^{mot} to $[0, 1]$ and optionally apply a soft
threshold to suppress weak motion.

390 **Geometric consistency across views.** Motion cues alone
391 do not guarantee cross-view consistency. We addition-
392 ally render a geometry-consistent mask by sampling 3D
393 points along agent trajectories (or Bézier curves) and pro-
394 jecting them into each camera view using known intrin-
395 sics and extrinsics. Let $\mathbf{p} = (x, y)$ be a pixel location

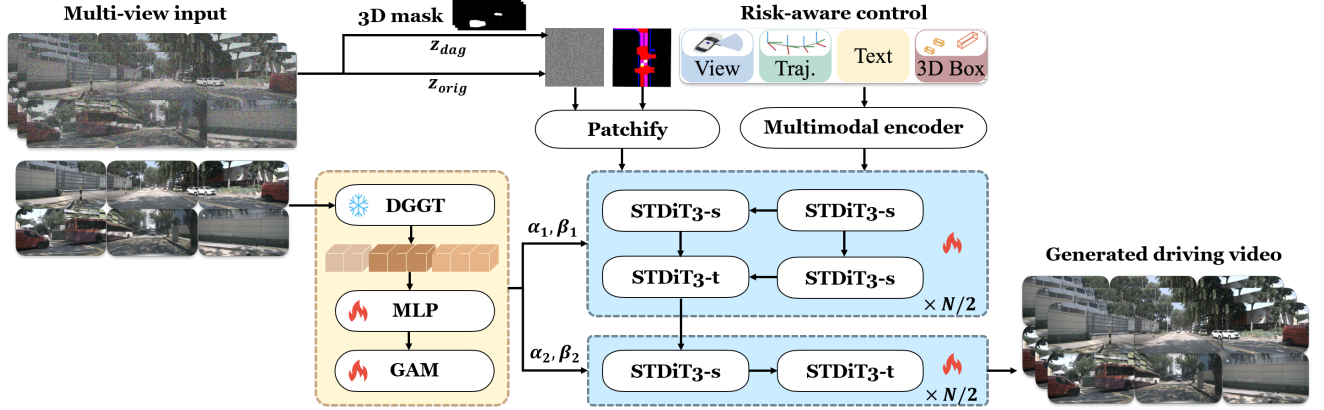


Figure 2. Overview of the proposed *RiskMV-DPO* framework. Given multi-view observations and scenario context, trajectories and 3D bounding boxes generated by the risk control module at a specified risk level are used as structured motion conditions. A multimodal encoder embeds view tokens and motion cues, which are then injected into a diffusion backbone composed of spatial and temporal STDiT3 blocks. Region-aware DPO further aligns the generation toward localized dynamic regions, producing temporally coherent multi-view driving videos consistent with the specified risk-conditioned motion.

396 and $\mathbf{u}_{bck} = (u_{bck}^x, u_{bck}^y)$ be the projected center of the k -
 397 th point for batch b and camera c , with visibility indicator
 398 $V_{bck} \in \{0, 1\}$. We rasterize a Gaussian blob:

$$M_{bck}^{\text{geo}}(\mathbf{p}) = V_{bck} \exp\left(-\frac{(x - u_{bck}^x)^2}{2\sigma_w^2} - \frac{(y - u_{bck}^y)^2}{2\sigma_h^2}\right). \quad (24)$$

399 We aggregate over sampled points to obtain a per-frame,
 400 per-view geometric mask $M_{b,c,t}^{\text{geo}}(\mathbf{p})$.
 401

402 **Mask fusion.** We fuse motion and geometric cues to form
 403 the final corruption mask:

$$M = \text{clip}(M^{\text{geo}} \odot M^{\text{mot}}, 0, 1), \quad (25)$$

405 so that corruption focuses on motion-dominant regions
 406 while remaining consistent across cameras.

407 5.3. Region-Aware Direct Preference Optimization

408 Global preference signals can be ambiguous for driving
 409 videos and may underweight dynamic regions that domi-
 410 nate risk. We align the model using localized preference
 411 pairs constructed on the masked regions, while keeping the
 412 unmasked background unchanged.

413 **Progressive corruption fusion.** Given a noise level t and
 414 clean latent z_0 , we corrupt only the masked region and keep
 415 the unmasked region clean:

$$z_t^{\text{mask}} = M \odot z_t + (1 - M) \odot z_0, \quad (26)$$

417 where z_t is the noised latent at level t . To construct lo-
 418 calized preference pairs, we use two corruption strengths

$t_w < t_l$ on the masked region:

$$\begin{aligned} z_{t_w}^{\text{mask}} &= M \odot z_{t_w} + (1 - M) \odot z_0, \\ z_{t_l}^{\text{mask}} &= M \odot z_{t_l} + (1 - M) \odot z_0. \end{aligned} \quad (27)$$

t_w yields an easier (preferred) target on the masked region,
 421 while t_l yields a harder (dispreferred) one, providing local-
 422 ized supervision without altering global context.
 423

424 **Masked flow-matching proxy.** We maintain an EMA
 425 reference model v_{ref} . Using the negative masked flow-
 426 matching residual as a proxy for log-probability, we define

$$\begin{aligned} \log p_{\theta}(y | x) &\approx -\text{FM}_{\theta}(y), \\ \delta_{\theta} &= v_{\theta}(z_t^{\text{mask}}, t, \mathbf{c}) - (\epsilon - z_0), \\ \text{FM}_{\theta}(y) &= \frac{1}{\|M\|_1} \|M \odot \delta_{\theta}\|_2^2, \end{aligned} \quad (28)$$

428 where \mathbf{c} denotes the conditioning (history, map, optional
 429 text, and motion tokens), and $\|M\|_1$ normalizes by masked
 430 area.

431 **Region-aware DPO objective.** We apply DPO on the lo-
 432 calized pair (y_w, y_l) :

$$\mathcal{L}_{\text{RA-DPO}} = -w(t) \log \sigma\left(\beta [\Delta_{\theta}(y_w) - \Delta_{\theta}(y_l)]\right), \quad (29)$$

$$\Delta_{\theta}(y) = \log p_{\theta}(y | x) - \log p_{\text{ref}}(y | x),$$

434 where $\sigma(\cdot)$ is the sigmoid, β is the DPO temperature, and
 435 $w(t)$ is a noise-adaptive weight. This objective aligns gen-
 436 eration with localized preferences in dynamic regions.

Method	FID↓	FVD↓	mAP↑
DriveDreamer-2	25.00	105.10	–
Panacea	16.96	139.00	–
MagicDriveV2	20.91	94.84	18.17
RiskMV-DPO (ours)	15.70	87.65	30.50

Table 1. Comparison with representative driving video generators on nuScenes.

Supervised fine-tuning loss. We apply supervised training on the masked region using the flow-matching residual.

$$\mathcal{L}_{\text{SFT}} = \frac{1}{\|M\|_1} \|M \odot \delta_\theta\|_2^2, \quad (30)$$

$$\delta_\theta = v_\theta(z_t^{\text{mask}}, t, \mathbf{c}) - (\epsilon - z_0).$$

Total objective. The complete training objective is

$$\mathcal{L}_{\text{total}} = \lambda_{\text{SFT}} \mathcal{L}_{\text{SFT}} + \lambda_{\text{RA}} \mathcal{L}_{\text{RA-DPO}} + \lambda_{\text{align}} \mathcal{L}_{\text{align}}. \quad (31)$$

6. Experiments

6.1. Experimental Setup

We conduct experiments on the nuScenes dataset [1] using the MagicDrive-STDiT3 backbone. We use 700 scenes for training and 150 scenes for evaluation. All videos are processed at 448×840 resolution per camera, and each clip contains $T=16$ frames. We compare RiskMV-DPO with MagicDriveV2 [8], DriveDreamer-2 [42], and Panacea [33], and also report ablations of our components.

We report eight metrics that cover visual fidelity, temporal quality, and multi-view geometry consistency. FID measures overall visual quality, and FVD measures temporal coherence. We evaluate 3D realism using mAP from a pretrained 3D detector applied to generated videos. To assess multi-view consistency, we report MV-SSIM computed on overlapping regions between adjacent cameras. For geometry accuracy, we report Depth AbsRel using VGGT-predicted depth against reference depth. We further report Foreground FID (Fg-FID) on cropped foreground objects and DINO-FID in DINO feature space for semantic fidelity.

We initialize from the MagicDriveV2 checkpoint [8] and train for 20,000 steps with AdamW. We use a learning rate of 2×10^{-5} and batch size 2 per GPU on 8 GPUs. We adopt BF16 training and an EMA model with decay $\gamma=0.9999$. The trainable parameters include the VGGT geometry adapter ($\sim 45\text{M}$), the GAM module ($\sim 12\text{M}$), and the motion-mask generator ($< 1\text{M}$), totaling $\sim 58\text{M}$ trainable parameters on top of a $\sim 2.5\text{B}$ diffusion backbone. All experiments are run on NVIDIA H800-90GB GPUs.

6.2. Main Results

After extensive experiments, we find that the trajectories and 3D bounding boxes generated under a specified risk

control lead to synthesized scenarios whose measured risk closely matches the target value. Due to space limitations, we omit the detailed analysis here.

Table 1 reports quantitative results on nuScenes. RiskMV-DPO achieves consistent improvements over MagicDriveV2. It reduces FID from 20.91 to 15.70 and improves FVD from 94.84 to 87.65, indicating better visual fidelity and temporal coherence. It also substantially improves 3D detection mAP from 18.17 to 30.50, suggesting that the generated videos contain more realistic and geometrically consistent 3D cues. Compared with other recent methods, RiskMV-DPO improves over DriveDreamer-2 (FID 25.00 \rightarrow 15.70) and Panacea (FID 16.96 \rightarrow 15.70).

6.3. Ablation Study

Table 2 evaluates the contribution of the key components in our framework. Starting from the baseline (A), introducing Random Mask DPO (B) slightly improves temporal quality and geometric consistency, but leaves FID unchanged. Replacing random masks with motion-aware masks (C) yields a clearer gain, improving FID (20.91 \rightarrow 19.42) and MV-SSIM (0.812 \rightarrow 0.825), indicating that focusing training on dynamic regions provides more informative localized preference signals. Using a 3D multi-view mask (D) further strengthens cross-view consistency, improving MV-SSIM to 0.841 and reducing Depth AbsRel to 0.228.

Next, adding VGGT-GAM without alignment (E) substantially improves geometric accuracy, reducing Depth AbsRel from 0.228 to 0.205, which validates the benefit of injecting geometry priors. Finally, enabling the geometry-appearance alignment loss (F) achieves the best overall performance, with FID 15.70, MV-SSIM 0.856, and Depth AbsRel 0.204.

6.4. Qualitative Visualization Across Risk Levels

Figure 3 shows qualitative examples generated by RiskMV-DPO at different target risk levels. Due to space limitations, we show three representative views (front-left, front, and front-right) and visualize a 7-frame clip by sampling one frame every 5 frames. As shown in Fig. 3(a), the higher-quantile scenario depicts the ego vehicle going straight through an unsignalized intersection while encountering a fast left-turning heavy truck at very close distance under partial occlusion, resulting in an imminent collision risk. In the same scene, the medium-quantile example in Fig. 3(b) shows the truck turning more slowly and keeping a larger clearance to the ego vehicle, which reduces the near-collision risk. The lower-quantile case in Fig. 3(c) corresponds to a common straight-driving pattern with no immediate conflict, leading to a low estimated risk.

We note that *high*, *medium*, and *low* risk do not have a strict boundary in our setting. Instead, the reported quantiles indicate the relative position of the induced risk coef-

Configuration	FID↓	FVD↓	MV-SSIM↑	Depth AbsRel↓	Fg-FID↓	DINO-FID↓
(A) Baseline	20.91	94.84	0.812	0.250	35.0	28.4
(B) + Random Mask DPO	20.91	93.01	0.818	0.243	33.6	27.9
(C) + Motion-Aware Mask	19.42	91.67	0.825	0.236	32.1	27.4
(D) + 3D MV Mask	19.33	89.51	0.841	0.228	31.0	26.8
(E) + VGGT-GAM (no Align)	16.79	88.70	0.848	0.205	30.2	26.3
(F) RiskMV-DPO (ours)	15.70	87.65	0.856	0.204	29.5	25.8

Table 2. Ablation study on nuScenes.

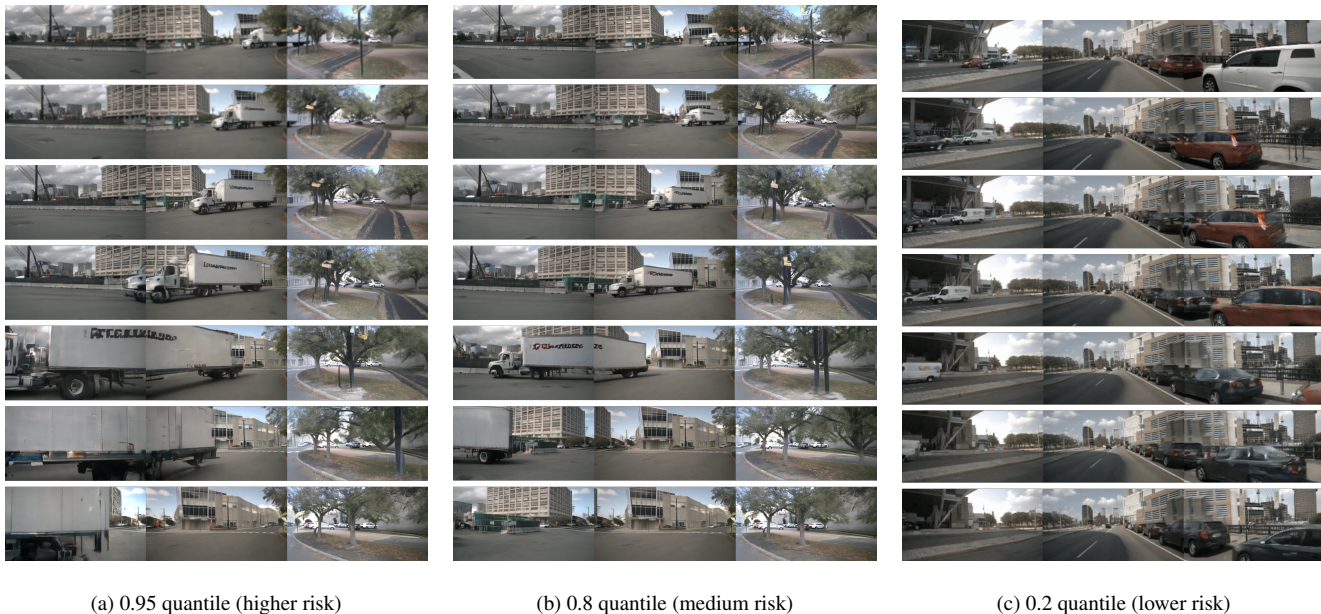


Figure 3. Qualitative examples generated by RiskMV-DPO at different target risk quantiles. The quantiles indicate the relative position of the induced risk coefficient among all generated scenarios. (a) Higher risk: an unsignalized intersection with a fast left-turning heavy truck at very close distance under partial occlusion. (b) Medium risk: the same interaction with slower turning and larger clearance, reducing near-collision risk. (c) Lower risk: a common straight-driving scenario with no immediate conflict.

525 efficient among all generated scenarios, where 0.95, 0.8, and
 526 0.2 correspond to higher, medium, and lower portions of
 527 the risk distribution, respectively, which is apparent on the
 528 provided examples. Code and additional qualitative results
 529 (including videos) are available at <https://github.com/venshow-w/RiskMV-DPO>.
 530

531 7. Conclusion and Limitations

532 In this paper, we present RiskMV-DPO, a general and sys-
 533 tematic pipeline for physically-informed, risk-controllable
 534 multi-view driving scenario generation. It transforms driv-
 535 ing risk from an after-the-fact label into an actionable con-
 536 trol signal that guides scene generation. By coupling risk-
 537 conditioned motion synthesis with diffusion-based multi-
 538 view video generation, the framework produces temporally
 539 coherent driving scenes consistent with specified risk lev-
 540 els. Experimental results on the nuScenes dataset demon-

strate that RiskMV-DPO achieves state-of-the-art perfor-
 mance, notably improving 3D detection mAP to 30.50 and
 reducing FID to 15.70. More broadly, our work highlights
 the potential of integrating risk-aware modeling with gen-
 erative world simulation for safety-oriented development of
 autonomous driving systems.

547 Despite these results, several limitations remain. While
 548 the proposed risk modeling framework captures a range
 549 of typical traffic interactions, many potential risk patterns
 550 in complex real-world environments remain underexplored.
 551 Besides, the overall performance depends on the quality of
 552 trajectories and 3D bounding boxes produced by the risk
 553 control module. Improving the robustness and generality of
 554 risk-conditioned motion generation is an important direc-
 555 tion for future research.

556

References

557

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

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

- [1] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 7
- [2] Xuan Cai, Xuesong Bai, Zhiyong Cui, Danmu Xie, Daocheng Fu, Haiyang Yu, and Yilong Ren. Text2scenario: Text-driven scenario generation for autonomous driving test. *Automotive Innovation*, pages 1–26, 2026. 1
- [3] Xiaoxue Chen, Ziyi Xiong, Yuantao Chen, Gen Li, Nan Wang, Hongcheng Luo, Long Chen, Haiyang Sun, Bing Wang, Guang Chen, Hangjun Ye, Hongyang Li, Ya-Qin Zhang, and Hao Zhao. Dggt: Feedforward 4d reconstruction of dynamic driving scenes using unposed images. *arXiv preprint arXiv:2512.03004*, 2025. 5
- [4] Tommaso Dreossi, Daniel J Fremont, Shromona Ghosh, Edward Kim, Hadi Ravanbakhsh, Marcell Vazquez-Chanlatte, and Sanjit A Seshia. Verifai: A toolkit for the formal design and analysis of artificial intelligence-based systems. In *International Conference on Computer Aided Verification*, pages 432–442. Springer, 2019. 2
- [5] Yang Fei, Peng Shi, Yang Liu, and Liang Wang. Critical roles of control engineering in the development of intelligent and connected vehicles. *Journal of Intelligent and Connected Vehicles*, 7(2):79–85, 2024. 1
- [6] Daniel J Fremont, Tommaso Dreossi, Shromona Ghosh, Xiangyu Yue, Alberto L Sangiovanni-Vincentelli, and Sanjit A Seshia. Scenic: a language for scenario specification and scene generation. In *Proceedings of the 40th ACM SIGPLAN conference on programming language design and implementation*, pages 63–78, 2019. 2
- [7] 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. 2
- [8] Ruiyuan Gao, Kai Chen, et al. Magicdrive-v2: High-resolution long video generation for autonomous driving with adaptive control. *arXiv preprint arXiv:2411.13807*, 2024. 7
- [9] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 28135–28144, 2025. 2
- [10] Xu He, Ji Li, Chuan Hu, Mingming Liu, and Hongming Xu. Dual-discriminator generative adversarial network with long-tail feature capture for extreme scenarios in human-machine shared driving. *IEEE Transactions on Intelligent Transportation Systems*, 2025. 1
- [11] Yixu He, Hongyi Lin, Lan Yang, and Yang Liu. Generative models for the evolution of transportation systems. *Journal of Traffic and Transportation Engineering (English Edition)*, 2025. 1
- [12] Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shotton, and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving. *arXiv preprint arXiv:2309.17080*, 2023. 2
- [13] Zitong Huang, Kaidong Zhang, Yukang Ding, Chao Gao, Rui Ding, Ying Chen, and Wangmeng Zuo. Mind the generative details: Direct localized detail preference optimization for video diffusion models. *arXiv preprint arXiv:2601.04068*, 2026. 2
- [14] 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*, pages 5820–5829, 2021. 2
- [15] Mark Koren, Saud Alsaif, Ritchie Lee, and Mykel J Kochenderfer. Adaptive stress testing for autonomous vehicles. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, pages 1–7. IEEE, 2018. 2
- [16] Xiaofan Li, Yifu Zhang, and Xiaoqing Ye. Drivingdiffusion: Layout-guided multi-view driving scenarios video generation with latent diffusion model. In *European Conference on Computer Vision*, pages 469–485. Springer, 2024. 2
- [17] Hongyi Lin, Yang Liu, Shen Li, and Xiaobo Qu. How generative adversarial networks promote the development of intelligent transportation systems: A survey. *IEEE/CAA journal of automatica sinica*, 10(9):1781–1796, 2023. 1
- [18] Hongyi Lin, Yang Liu, Liang Wang, and Xiaobo Qu. Big data-driven advancements and future directions in vehicle perception technologies: From autonomous driving to modular buses. *IEEE Transactions on Big Data*, 11(3):1568–1587, 2025. 2
- [19] Hongyi Lin, Yang Liu, Liang Wang, and Xiaobo Qu. A high-precision calibration and evaluation method based on binocular cameras and lidar for intelligent vehicles. *IEEE Transactions on Vehicular Technology*, 74(5):7404–7415, 2025. 1
- [20] Qichao Liu, Heye Huang, Shiyue Zhao, Lei Shi, Soyoung Ahn, and Xiaopeng Li. Risknet: interaction-aware risk forecasting for autonomous driving in long-tail scenarios. *Transportation Research Part E: Logistics and Transportation Review*, 205:104478, 2026. 2, 3
- [21] Hannan Lu, Xiaohe Wu, Shudong Wang, Xiameng Qin, Xinyu Zhang, Junyu Han, Wangmeng Zuo, and Ji Tao. Seeing beyond views: Multi-view driving scene video generation with holistic attention. *arXiv preprint arXiv:2412.03520*, 2024. 2
- [22] Yuewen Mei, Tong Nie, Jian Sun, and Ye Tian. Llm-attacker: Enhancing closed-loop adversarial scenario generation for autonomous driving with large language models. *IEEE Transactions on Intelligent Transportation Systems*, 2025. 1
- [23] Wassim G. Najm, David L. Smith, and Mikio Yanagisawa. Pre-crash scenario typology for crash avoidance research. Technical Report DOT HS 810 767, National Highway Traffic Safety Administration, 2007. Based on 2004 NASS GES crash database. 3
- [24] Xiaobo Qu, Hongyi Lin, and Yang Liu. Envisioning the future of transportation: Inspiration of chatgpt and large models, 2023. 1

- 669 [25] Shai Shalev-Shwartz, Shaked Shammah, and Amnon
670 Shashua. On a formal model of safe and scalable self-driving
671 cars. *arXiv preprint arXiv:1708.06374*, 2017. 2 727
- 672 [26] Elizabeth D. Swanson, Frank Foderaro, Mikio Yanagisawa,
673 Wassim G. Najm, and Philip Azeredo. Statistics of light-
674 vehicle pre-crash scenarios based on 2011–2015 national
675 crash data. Technical Report DOT HS 812 745, National
676 Highway Traffic Safety Administration, 2019. Scenario
677 groups and statistics based on 2011–2015 FARS and NASS
678 GES. 3 728
- 679 [27] Xuewei Tang, Mengmeng Yang, Tuopu Wen, Peijin Jia, Le
680 Cui, Mingshan Luo, Kehua Sheng, Bo Zhang, Kun Jiang,
681 and Diange Yang. Priorfusion: Unified integration of priors
682 for robust road perception in autonomous driving. *Communi-
683 cations in Transportation Research*, 5:100229, 2025. 2 729
- 684 [28] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou,
685 Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming
686 Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model align-
687 ment using direct preference optimization. In *Proceedings of
688 the IEEE/CVF Conference on Computer Vision and Pattern
689 Recognition*, pages 8228–8238, 2024. 2 730
- 690 [29] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea
691 Vedaldi, Christian Rupprecht, and David Novotny. Vggt:
692 Visual geometry ruppert transformer. In *Proceedings of
693 the IEEE/CVF Conference on Computer Vision and Pattern
694 Recognition (CVPR)*, 2025. 5 731
- 695 [30] Jie Wang, Guang Li, Zhijian Huang, Chenxu Dang, Hangjun
696 Ye, Yahong Han, and Long Chen. Vggdrive: Em-
697 powering vision-language models with cross-view geomet-
698 ric grounding for autonomous driving. *arXiv preprint
699 arXiv:2602.20794*, 2026. 1 732
- 700 [31] Xiaofeng Wang, Zheng Zhu, Guan Huang, Xinze Chen, Jia-
701 gang Zhu, and Jiwen Lu. Drivedreamer: Towards real-world-
702 drive world models for autonomous driving. In *European
703 conference on computer vision*, pages 55–72. Springer, 2024.
704 2 733
- 705 [32] Yong Wang, Daifeng Zhang, Yanqiang Li, Liguu Shuai,
706 Zhicheng Tang, and Yuxiang Hou. Safety-critical scenario
707 test for intelligent vehicles via hybrid participation of natural
708 and adversarial agents. *Journal of Intelligent and Connected
709 Vehicles*, 8(3):9210066–1, 2025. 1 734
- 710 [33] Yu Wen et al. Panacea: Panoramic and controllable video
711 generation for autonomous driving. In *Proceedings of the
712 IEEE/CVF Conference on Computer Vision and Pattern
713 Recognition (CVPR)*, 2024. 7 735
- 714 [34] Lukas Westhofen, Christian Neurohr, Tjark Koopmann,
715 Martin Butz, Barbara Schütt, Fabian Utesch, Birte Neurohr,
716 Christian Gutenkunst, and Eckard Böde. Criticality metrics
717 for automated driving: A review and suitability analysis of
718 the state of the art: L. westhofen et al. *Archives of Computa-
719 tional Methods in Engineering*, 30(1):1–35, 2023. 2 736
- 720 [35] Wei Wu, Xi Guo, Weixuan Tang, Tingxuan Huang, Chiyu
721 Wang, and Chenjing Ding. Drivescape: High-resolution
722 driving video generation by multi-view feature fusion. In
723 *Proceedings of the Computer Vision and Pattern Recognition
724 Conference*, pages 17187–17196, 2025. 2 737
- 725 [36] Ziyi Wu, Anil Kag, Ivan Skorokhodov, Willi Menapace,
726 Ashkan Mirzaei, Igor Gilitschenski, Sergey Tulyakov, and
Aliaksandr Siarohin. Densedpo: Fine-grained temporal
preference optimization for video diffusion models. *arXiv
preprint arXiv:2506.03517*, 2025. 2 738
- [37] Chejian Xu, Aleksandr Petiushko, Ding Zhao, and Bo Li.
Diffscene: Diffusion-based safety-critical scenario gener-
ation for autonomous vehicles. In *Proceedings of the
AAAI conference on artificial intelligence*, pages 8797–8805,
2025. 1 739
- [38] Mingxing Xu, Hongyi Lin, and Yang Liu. A deep learn-
ing approach for vehicle velocity prediction considering the
influence factors of multiple lanes. *Electronic Research
Archive*, 31(1):401, 2023. 2 740
- [39] Cong Zhang, Bangyang Wei, Yang Liu, and Samuel Labi.
World model-based long-tail and scenario-specific genera-
tion for autonomous driving. *Journal of Intelligent and Con-
nected Vehicles*, 2026. 1 741
- [40] Jiawei Zhang, Chejian Xu, and Bo Li. Chatscene:
Knowledge-enabled safety-critical scenario generation for
autonomous vehicles. In *Proceedings of the IEEE/CVF Con-
ference on Computer Vision and Pattern Recognition*, pages
15459–15469, 2024. 1 742
- [41] Meng Zhang, Jiatong Xu, Ying Gao, Dandan Shen, and Zhi-
gang Xu. Can combined virtual-real testing speed up au-
tonomous vehicle testing? findings from aeb field exper-
iments. *Communications in Transportation Research*, 5:
100216, 2025. 1 743
- [42] 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. *arXiv preprint arXiv:2403.06845*,
2024. 7 744
- [43] 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 Con-
ference on Artificial Intelligence*, pages 10412–10420, 2025.
2 745
- [44] Wenzhao Zheng, Ruiqi Song, Xianda Guo, Chenming
Zhang, and Long Chen. Genad: Generative end-to-end au-
tonomous driving. In *European Conference on Computer
Vision*, pages 87–104. Springer, 2024. 4 746
- [45] Zhiyuan Zhou, Heye Huang, Boqi Li, Shiyue Zhao, Yao Mu,
and Jianqiang Wang. Safedrive: Knowledge-and data-driven
risk-sensitive decision-making for autonomous vehicles with
large language models. *Accident Analysis & Prevention*,
224:108299, 2026. 2 747
- [46] Huaisheng Zhu, Teng Xiao, and Vasant G Honavar. Dspo:
Direct score preference optimization for diffusion model
alignment. International Conference on Learning Representa-
tions (ICLR 2025), 2025. 2 748