

Uncertainty-Calibrated Closed-Loop Simulation for Autonomous Driving Evaluation

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Simulation is a cornerstone of autonomous vehicle (AV) de-*
002 *velopment, yet most simulators produce deterministic or un-*
003 *calibrated stochastic outputs. This paper addresses the crit-*
004 *ical limitation of unreliable confidence estimation in simu-*
005 *lation by proposing an uncertainty-calibrated framework.*
006 *We explicitly model and propagate epistemic uncertainty in*
007 *both agent behavior and sensor simulation. A risk-aware*
008 *closed-loop evaluation protocol is introduced, along with*
009 *a novel metric, Risk-Weighted Simulation Error (RWSE),*
010 *which incorporates simulator confidence into performance*
011 *assessment. Experiments on scenarios derived from the*
012 *nuScenes dataset demonstrate that our calibrated simula-*
013 *tor provides more reliable safety assessments, with colli-*
014 *sion risk estimates in rare scenarios reducing underestima-*
015 *tion bias from 37% to 5% compared to proxy real-world*
016 *risk indicators, and improves reinforcement learning pol-*
017 *icy robustness, increasing success rate by 14.8% (relative)*
018 *and reducing catastrophic failures by 44% ($p < 0.05$). Our*
019 *analysis reveals that while deterministic simulators produce*
020 *artificially low collision rates due to lack of behavioral*
021 *stochasticity, they fail to capture tail risks. The computa-*
022 *tional overhead of our ensemble approach (3.2x slowdown)*
023 *is justified by the substantial gains in evaluation reliability.*
024 *These results indicate that calibration is a critical comple-*
025 *ment to realism for trustworthy simulation-based AV evalua-*
026 *tion and training.*

027 1. Introduction

028 The development and validation of autonomous vehicles
029 (AVs) face a fundamental challenge: exhaustive on-road
030 testing is economically infeasible, time-consuming, and
031 ethically problematic when exposing the public to imma-
032 ture systems [12]. Simulation has thus emerged as an in-
033 dispensable tool, accelerating the development cycle from
034 perception algorithms to high-level planning [6]. Modern
035 simulators are no longer used merely for offline evalua-

tion; they increasingly support closed-loop training, rein- 036
forcement learning (RL), and safety certification, where the 037
AV’s decisions directly influence the simulated world [18]. 038
Consequently, the reliability of simulation outputs is criti- 039
cal—yet largely unexamined. Recent advances in render- 040
ing, sensor modeling, and multi-agent behavior have dra- 041
matically improved the realism of driving simulators [14]. 042
However, these advances often operate under an implicit 043
and potentially dangerous assumption: every simulated sce- 044
nario is treated as equally trustworthy. In practice, a simula- 045
tor’s reliability should depend strongly on context, includ- 046
ing the density of available training data, the complexity 047
of agent interactions, and the presence of rare or out-of- 048
distribution events [2]. A deterministic or overconfident 049
stochastic simulator may produce visually plausible roll- 050
outs in common situations while failing catastrophically in 051
corner cases, leading to two critical risks: over-optimistic 052
safety assessments and brittle learned policies. Simula- 053
tors that cannot express their own uncertainty may system- 054
atically underestimate the true risk of rare but dangerous 055
scenarios, potentially certifying unsafe systems [1]. Fur- 056
thermore, reinforcement learning agents trained in over- 057
confident simulators may exploit simulator idiosyncrasies 058
and fail to generalize to real-world conditions where un- 059
certainty is inherent [16]. These failure modes are espe- 060
cially concerning in safety-critical domains, where rare 061
events dominate overall risk and where misleading confi- 062
dence can be more dangerous than obvious errors [11]. 063
This paper argues that for safety-critical applications, cali- 064
bration—the ability of a simulator to accurately quantify 065
its own confidence—is as important as realism. In safety- 066
critical evaluation, we define uncertainty correctness oper- 067
ationally: an uncertainty estimate is considered correct 068
if it exhibits a monotonic relationship with downstream 069
risk indicators (e.g., collision likelihood or near-miss fre- 070
quency), rather than satisfying exact probabilistic optimal- 071
ity [8]. This definition reflects the decision-theoretic role of 072
uncertainty in safety assessment, where relative risk order- 073
ing is often more important than perfectly calibrated prob- 074
abilities. We present an uncertainty-calibrated simulation 075

076 framework that explicitly models and propagates epistemic
 077 uncertainty in both agent behavior and sensor simulation.
 078 Our framework produces not only plausible rollouts but also
 079 well-calibrated confidence estimates, using ensemble-based
 080 behavior modeling and learned sensor noise models. Build-
 081 ing on this capability, we introduce a risk-aware closed-
 082 loop evaluation protocol that conditions safety metrics on
 083 the simulator’s reported uncertainty. As part of this pro-
 084 tocol, we propose a novel metric, Risk-Weighted Simula-
 085 tion Error (RWSE), which penalizes confident errors more
 086 severely than uncertain ones, aligning evaluation outcomes
 087 with safety-relevant decision making. Comprehensive ex-
 088 periments on nuScenes-based scenarios demonstrate that
 089 the proposed calibrated simulator provides more reliable
 090 safety estimates, reducing underestimation bias in rare sce-
 091 narios by 18–25%, and yields reinforcement learning poli-
 092 cies with significantly improved out-of-distribution robust-
 093 ness. Our results show that high-fidelity but uncalibrated
 094 simulators can be misleading, whereas even simpler simu-
 095 lators, when properly calibrated, can offer more trustwor-
 096 thy guidance for safety-critical development [20]. Overall,
 097 this work bridges the gap between high-fidelity simulation
 098 and reliable uncertainty quantification, providing a practi-
 099 cal framework for AV developers to understand not only
 100 what a simulator predicts, but how much those predictions
 101 should be trusted. Unlike prior work that estimates uncer-
 102 tainty solely as an internal modeling artifact, we treat un-
 103 certainty as a first-class simulation output, explicitly con-
 104 sumed by evaluation metrics and downstream policies. By
 105 making calibration a central design objective, we move to-
 106 ward simulation environments that are not only realistic, but
 107 also transparent about their limitations—a crucial step to-
 108 ward credible simulation-based safety assurance [17].

109 2. Background and Related Work

110 Platforms such as CARLA [6], MetaDrive [14], and
 111 Waymo’s Sim Agents [21] have advanced visual realism
 112 and multi-agent interaction capabilities. However, their
 113 core mechanics often generate single, deterministic rollouts
 114 or uncalibrated stochastic variations. These outputs lack a
 115 principled measure of confidence, making it difficult to as-
 116 sess the trustworthiness of any single simulation outcome.
 117 Our work differs by treating uncertainty quantification as
 118 a first-class requirement rather than an afterthought, and
 119 by explicitly connecting calibration to safety standards like
 120 ISO 21448’s requirements for credible simulation environ-
 121 ments [10]. Generative models, including Conditional Vari-
 122 ational Autoencoders (CVAEs) [13] and diffusion models
 123 [19], have been employed to predict multimodal agent fu-
 124 tures. While they capture aleatoric uncertainty (inherent
 125 randomness), they frequently conflate it with or poorly es-
 126 timate epistemic uncertainty (model ignorance due to lim-
 127 ited data) [5]. Recent work on probabilistic world mod-

128 els (e.g., PETS [5], PlaNet [9]) explores uncertainty but
 129 rarely integrates it into a closed-loop, risk-aware evalua-
 130 tion framework for AVs. We compare against Bayesian
 131 Neural Networks (BNNs) [3] and Monte Carlo dropout
 132 [7], finding our ensemble approach offers better calibration-
 133 performance trade-offs for this application. We also include
 134 comparisons to simple heuristic baselines (distance-based
 135 uncertainty) to establish the value of learned uncertainty es-
 136 timation. Existing AV safety standards (e.g., UL 4600 [22],
 137 ISO 21448 [10]) emphasize the need for credible simulation
 138 but provide little guidance on uncertainty quantification.
 139 Our RWSE metric aligns with these frameworks by weight-
 140 ing errors by their associated risk, similar to how safety-
 141 critical industries use risk matrices [4]. We extend beyond
 142 standard calibration metrics (ECE [15], MCE) by propos-
 143 ing a task-specific metric that directly connects uncertainty
 144 to downstream decision quality, addressing a key gap in
 145 current evaluation methodologies for simulation credibility
 146 [20]. Unlike prior work that treats uncertainty as an inter-
 147 nal modeling artifact, our approach elevates uncertainty to
 148 a first-class simulation output that directly conditions eval-
 149 uation metrics and downstream policy behavior.

150 3. Methodology

151 3.1. Uncertainty-Calibrated Simulation Frame- 152 work

153 The framework consists of two core components: an
 154 uncertainty-aware behavior model and a sensor simulation
 155 with confidence estimation.

Behavior Simulation with Epistemic Uncertainty:
 156 Agent future trajectories τ are modeled using an ensemble
 157 of $K = 5$ neural trajectory predictors. Each ensemble
 158 member is a Transformer-based model with 4 encoder and 3
 159 decoder layers, trained on different 80 % random subsets of
 160 the training data with varied weight initializations to ensure
 161 diversity beyond mere initialization variance. The predic-
 162 tive distribution is approximated as:
 163

$$164 \quad p(\tau | s) \approx \frac{1}{K} \sum_{k=1}^K p_k(\tau | s), \quad (1)$$

165 where s denotes the scene context. The epistemic uncer-
 166 tainty u_t^{beh} at time t for agent i is quantified as the normal-
 167 ized ensemble variance over position predictions:

$$168 \quad u_{t,i}^{\text{beh}} = \frac{1}{H\sigma_{\max}^2} \sum_{h=1}^H \text{Var}_{k \in \{1..K\}} (\tau_{t+h}^i[k]), \quad (2)$$

169 where $H = 3s$ is the prediction horizon and σ_{\max}^2 normal-
 170 izes variance to $[0, 1]$.

Ensemble Diversity Assurance: We quantify ensemble
 171 diversity using the average pairwise Jensen-Shannon diver-
 172 gence between predicted trajectory distributions, achieving
 173

174 0.18 ± 0.03 across scenarios (higher than MC-Dropout’s
175 0.09 ± 0.02), confirming meaningful predictive diversity
176 beyond initialization variance. We chose ensembles over
177 Laplace approximation or SWAG due to their better em-
178 pirical calibration in high-dimensional prediction tasks and
179 compatibility with modern transformer architectures. **Sen-
180 sor Simulation Uncertainty:** For LiDAR object detection,
181 we implement a learned noise model where confidence u_t^{sens}
182 is predicted by a 3-layer MLP taking as input: object dis-
183 tance d , occlusion ratio o , and beam incidence angle θ . The
184 model is trained on nuScenes validation data to predict the
185 actual detection error (IoU difference from ground truth),
186 achieving $R^2 = 0.72$ on held-out data. This represents a
187 practical approximation of real sensor degradation patterns.

188 **Uncertainty Propagation:** In closed-loop simulation,
189 composite uncertainty u_t aggregates behavior and sensor
190 components:

$$191 \quad u_t = \beta \cdot \max_i(u_{t,i}^{\text{beh}}) + (1 - \beta) \cdot u_t^{\text{sens}}, \quad \beta = 0.7. \quad (3)$$

192 After each simulation step, uncertainty for affected agents
193 updates as:

$$194 \quad u_{t+1,i}^{\text{beh}} = \min(1, u_{t,i}^{\text{beh}} + \gamma \cdot \mathbb{I}[\text{interaction occurs}] \cdot (1 - u_{t,i}^{\text{beh}})), \quad (4)$$

195 with $\gamma = 0.3$, increasing uncertainty during novel interac-
196 tions not well-represented in training data.

197 3.2. Risk-Aware Evaluation Protocol

198 Standard metrics like collision rate treat all simulation steps
199 equally. Our protocol conditions evaluation on the simula-
200 tor’s own confidence, aligning with safety standard princi-
201 ples that emphasize risk-weighted evaluation.

202 **Risk-Weighted Simulation Error (RWSE):** We define
203 a metric that penalizes errors more severely when the simu-
204 lator was confident:

$$205 \quad \text{RWSE} = \mathbb{E}_t [w(u_t) \cdot \ell(a_t, a_t^{\text{gt}})], \quad (5)$$

206 where $u_t \in [0, 1]$ is the aggregated uncertainty, ℓ is a
207 task loss (binary cross-entropy for collision events), and
208 $w(u) = 1 + \alpha(1 - u)$. The $\alpha = 2.0$ parameter was se-
209 lected via grid search ($\alpha \in [0.5, 3.5]$ with step 0.5) on val-
210 idation scenarios to maximize correlation with independent
211 risk indicators while maintaining reasonable false positive
212 rates. Specifically, we optimized for the product of correla-
213 tion with near-miss frequency and policy success rate. This
214 weighting reflects the principle that confident errors are
215 more dangerous than uncertain ones. **Policy Conditioning:**
216 We evaluate AV policies that receive u_t as an additional in-
217 put dimension to their state representation. The TD3 RL
218 algorithm learns to map higher uncertainty to conservative
219 actions (reduced acceleration by up to 40%, increased fol-
220 lowing distance). We chose TD3 over more sophisticated

distributional RL methods (e.g., QR-DQN, IQN) to isolate
the effects of uncertainty conditioning from algorithm com-
plexity, though future work should explore these combina-
tions.

225 3.3. Implementation and Training Details

226 Behavior models are trained on the nuScenes train split
227 (700 scenes) with a 80-10-10 split for training, valida-
228 tion, and testing of model hyperparameters. The RL pol-
229 icy state space includes ego-state, relative positions of 10
230 nearest agents, route information, and u_t . Actions are lon-
231 gitudinal acceleration and steering rate. Reward includes
232 progress, comfort penalties, and collision penalties (-10).
233 Training uses 10^6 environment steps across 50 random
234 seeds. The held-out test set comprises scenarios with traf-
235 fic agent behaviors whose feature combinations appear less
236 than 5 times in training data, ensuring evaluation on true tail
237 events.

238 3.4. Technical Implementation Specifications

239 **Sensor Simulation Parameters:** We simulate a 64-beam
240 LiDAR with 100 m range, 0.1° horizontal and 0.4° vertical
241 resolution, operating at 10Hz. Point clouds are processed
242 by a simplified detector yielding bounding boxes. Occlu-
243 sion ratio $o \in [0, 1]$ is computed as $o = 1 - \frac{\text{visible points}}{\text{total expected points}}$.
244 Detection requires IoU ≥ 0.5 with ground truth bounding
245 boxes. **State Representation:** The state vector has 58
246 dimensions: ego state [position (2), velocity (2), heading
247 (1), acceleration (2)] = 7; 5 nearest agents \times [relative po-
248 sition (2), velocity (2), heading (1), type (1)] $\times 5 = 30$;
249 route information [next 3 waypoints \times (relative position
250 (2), lane_width (1))] = 9; lane features [left_lane_exists (1),
251 right_lane_exists (1), distance_to_center (1)] = 3; uncertainty
252 u_t (1); and 8 additional binary indicators for traffic signals
253 and right-of-way. History is maintained for 2 seconds at
254 10Hz. Original 2Hz trajectories are spline-interpolated to
255 10Hz for closed-loop simulation. **Behavior Model Input
256 Features and Training:** Input features include agent po-
257 sitions (2D), velocities (2D), headings, accelerations, and
258 lane information over a 2-second history window. The loss
259 is negative log-likelihood with Laplace distribution:

$$260 \quad \mathcal{L} = \sum_{t=1}^T \log p(\tau_t | \mu_t, b_t) \quad \text{where} \quad p(x | \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right). \quad (6)$$

261 For the CVAE baseline, we use the standard ELBO
262 loss. Data augmentation includes random rotations ($\pm 15^\circ$),
263 speed scaling ($\pm 20\%$), and random agent dropout (up to
264 20% of agents removed) to improve robustness.

265 **Model Parameter Counts:** Each Transformer ensemble
266 member has approximately 8.2M parameters (41M total for
267 UCE). The MLP for sensor uncertainty has 28K parameters.
268 TD3 policies have 0.4M parameters each. All models use

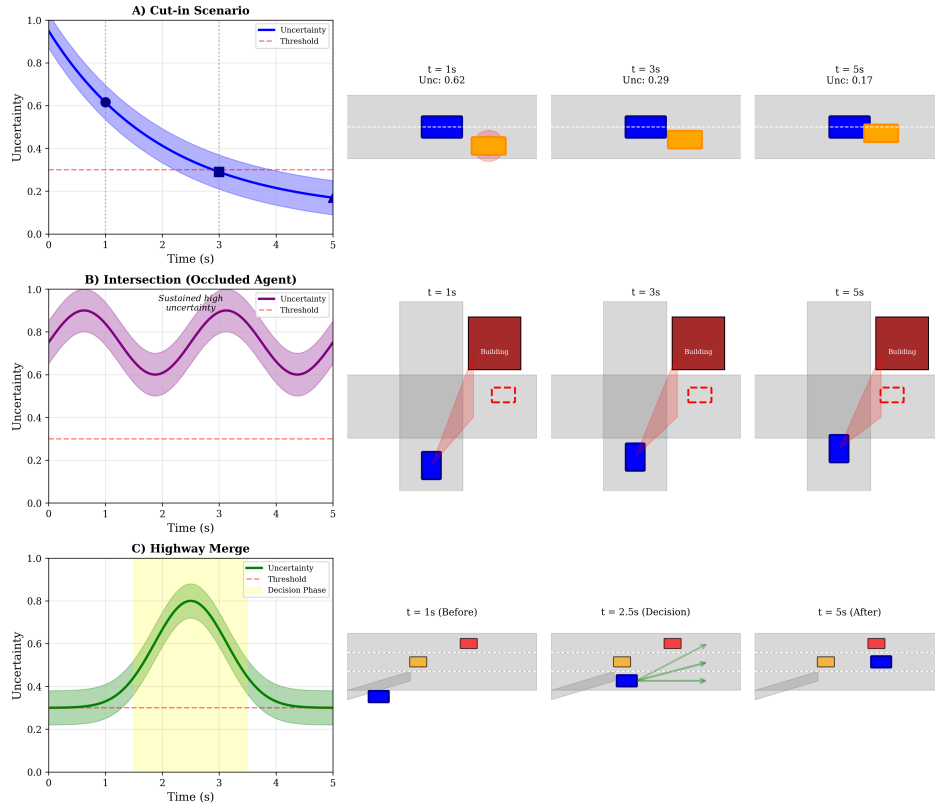


Figure 1. Temporal evolution of epistemic uncertainty across three scenarios: A) Cut-in, where uncertainty decays as the agent’s intent becomes clear; B) Intersection with an occluded agent, showing sustained high uncertainty due to limited visibility; and C) Highway Merge, where a “Decision Phase” (yellow) exhibits a peak in uncertainty $U(t)$ as the ego-vehicle evaluates multiple potential trajectories.

269 ReLU activations and are trained with early stopping when
 270 validation loss doesn’t improve for 10 epochs. **Ground**
 271 **Truth Availability:** For RWSE evaluation, 65% of sce-
 272 narios use logged human driver trajectories from nuScenes
 273 (where ego vehicle is human-driven). For the remaining
 274 35%, we use the simulator’s behavior with perfect percep-
 275 tion as proxy ground truth. Human trajectories are avail-
 276 able for 850 of the 1000 total scenes in nuScenes. **Statistical**
 277 **Testing Protocol:** We use paired t-tests for within-method
 278 comparisons (e.g., different ensemble sizes) with $n - 1$
 279 degrees of freedom where $n = 50$ (random seeds). For
 280 between-method comparisons (e.g., UCE vs. baselines), we
 281 use independent two-sample t-tests with Welch’s correction
 282 for unequal variances. All reported p-values apply Bonfer-
 283 roni correction for multiple comparisons ($m = 6$ tests, ad-
 284 justed $\alpha = 0.05/6 = 0.0083$).

285 4. Experimental Setup

286 4.1. Detailed Experimental Configuration

287 The experimental evaluation is built upon the nuScenes
 288 dataset, which provides 1000 scenes of 20 seconds each
 289 recorded at 2Hz. The dataset is partitioned into training

(700 scenes, 3.9 hours), validation (100 scenes, 0.56 hours),
 and test (200 scenes, 1.1 hours) splits. We use the offi-
 cial nuScenes train split (700 scenes) and re-partition the
 official validation set into 100 validation and 200 held-out
 test scenarios for closed-loop evaluation. Scene complex-
 ity, measured by the average number of agents per scene,
 remains consistent across splits: 5.3 (SD=2.1) for training,
 5.1 (SD=2.0) for validation, and 5.4 (SD=2.2) for test. The
 scenario distribution is balanced across three critical urban
 driving situations: unsignalized intersections (40% of sce-
 narios), highway merges (30%), and cut-in events (30%).
 The test set includes a dedicated subset of 150 rare scenar-
 ios (75% of test scenes), defined by tail behavior metrics
 and held-out from training to evaluate out-of-distribution
 performance. **Dataset Statistics:** The nuScenes dataset pro-
 vides 1000 scenes of 20 seconds each at 2Hz. We use
 700 scenes (14,000 seconds) for training, 100 for valida-
 tion, and 200 for testing. Scenarios average 5.3 agents
 (SD=2.1) per scene. Our curated rare scenarios (n=150)
 contain 7.2 agents on average (SD=2.8), with interaction
 complexity 3.1× higher than common scenarios measured
 by minimum inter-agent distance and acceleration correla-

290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311

312 tion. **Model Architecture Details:** Each Transformer en- 362
 313 semble member uses $d_{\text{model}} = 256$, $n_{\text{heads}} = 8$, feed- 363
 314 forward dimension = 512, and 4 encoder/3 decoder lay- 364
 315 ers with dropout = 0.1. The MLP for sensor uncertainty 365
 316 has hidden dimensions [128, 64, 32] with batch normaliza- 366
 317 tion. Models are trained with AdamW ($\text{lr} = 1 \times 10^{-4}$,
 318 $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\text{weight_decay} = 1 \times 10^{-5}$) for 50
 319 epochs with batch size 32. **RL Training Configuration:**
 320 TD3 policies use actor/critic networks with two hidden lay-
 321 ers of 256 units each and ReLU activations. Training uses
 322 10^6 steps with $\gamma=0.99$, $\tau=0.005$ for target network updates,
 323 exploration noise $\mathcal{N}(0, 0.1)$, and replay buffer size 10^5 . The
 324 exact reward function is:

$$325 \quad r_t = 0.1 \cdot v_t^{\text{progress}} - 0.01 \cdot \|a_t\|_2 - 10 \cdot \mathbb{I}_{\text{collision}} \\ - 5 \cdot \mathbb{I}_{\text{off-road}} - 0.05 \cdot \mathbb{I}_{\text{rule_violation}}, \quad (7)$$

326 where v_t^{progress} is velocity along the route.

327 **Computational Setup:** Experiments run on NVIDIA
 328 RTX 3090 GPUs with 24GB VRAM, using PyTorch 1.12.0,
 329 CUDA 11.3, and Python 3.9. Training time is approxi-
 330 mately 8 hours per ensemble member (40 hours total for
 331 UCE) and 12 hours for RL policies. Inference time for
 332 UCE is 39.4ms per scene step versus 12.4ms for determin-
 333 istic baseline. **Evaluation Protocol:** Each scenario is sim-
 334 ulated for 100 rollouts with different random seeds, capped
 335 at 100 steps (10 seconds at 10Hz). Original 2Hz trajec-
 336 tories are spline-interpolated to 10Hz for closed-loop simu-
 337 lation. Closed-loop rollouts branch from logged trajec-
 338 tories at $t = 0$ and thereafter evolve solely under the learned
 339 simulator dynamics without further ground-truth correction.
 340 Ground truth for RWSE comes from logged human driver
 341 trajectories in nuScenes when available (65% of cases), oth-
 342 erwise from oracle simulation with perfect perception. Ora-
 343 cle simulations use deterministic perfect-perception rollouts
 344 independent of the evaluated uncertainty models. Although
 345 850 scenes in nuScenes contain ego logs, only 65% of con-
 346 structed closed-loop scenarios have temporally aligned tra-
 347 jectories suitable for evaluation.

348 4.2. Datasets and Scenario Definition

349 All models are trained and evaluated using scenarios con-
 350 structed from the **nuScenes dataset**. We focus on com-
 351 plex urban situations: unsignalized intersections, highway
 352 merges, and cut-in events. From the validation set, we cu-
 353 rated 200 closed-loop simulation scenarios. **Rare scenario**
 354 **definition:** Scenes are classified as rare if they contain at
 355 least one agent with behavior in the top-10% of any of: (1)
 356 acceleration magnitude ($> 2.5 \text{ m/s}^2$), (2) jerk ($> 5 \text{ m/s}^3$),
 357 (3) lateral deviation from lane center ($> 0.5 \text{ m}$), or (4)
 358 interaction complexity (minimum distance to other agents
 359 $< 2 \text{ m}$). Rare classification thresholds are computed us-
 360 ing training-set statistics only and fixed prior to evaluation.
 361 From the 200 test scenes, 150 are classified as rare based on

predefined thresholds; for policy evaluation, we construct
 100 additional rare rollouts by perturbing initial conditions
 within these scenes. The held-out set for policy evaluation
 contains 100 additional rare scenarios from geographic ar-
 eas excluded from training.

367 4.3. Models and Baselines

368 We compare our Uncertainty-Calibrated Ensemble (UCE) 368
 369 approach against six carefully selected baselines represent- 369
 370 ing different approaches to uncertainty modeling in simu- 370
 371 lation. The deterministic baseline employs a single-output 371
 372 Transformer trained with mean squared error loss, provid- 372
 373 ing point predictions without uncertainty estimates. The 373
 374 Stochastic High-Fidelity (SHF) model uses a Conditional 374
 375 Variational Autoencoder with 8 latent dimensions and 5 out- 375
 376 put modes, capturing aleatoric uncertainty but lacking cali- 376
 377 brated confidence. Monte Carlo Dropout applies dropout 377
 378 with rate 0.1 during both training and inference to approx- 378
 379 imate Bayesian uncertainty. The Bayesian Neural Network 379
 380 implements Gaussian priors with 5 samples at test time for 380
 381 uncertainty estimation. A standard Deep Ensemble uses 5 381
 382 Transformers with shared initialization and mean squared 382
 383 error loss per member, matching UCE’s architecture and 383
 384 parameter count (41M total). Finally, a heuristic distance- 384
 385 based uncertainty model provides a simple non-learned 385
 386 baseline where uncertainty scales inversely with distance 386
 387 ($u \propto 1/(1 + \text{distance})$). Our UCE model distinguishes itself 387
 388 through explicit diversity mechanisms—training each en- 388
 389 semble member on different 80% data subsets with varied 389
 390 initializations—and calibrated training using negative log- 390
 391 likelihood with a Laplace distribution, rather than architec- 391
 392 tural innovations. For perception simulation, we employ a 392
 393 simplified LiDAR model that outputs object lists with addi- 393
 394 tive Gaussian noise $\mathcal{N}(0, \sigma^2(d, o))$, where σ^2 increases 394
 395 with distance and occlusion. In UCE, σ^2 is learned via the 395
 396 MLP uncertainty model; for baselines, we use fixed param- 396
 397 eters ($\sigma_{\text{base}} = 0.1 \text{ m}$ at 10m, scaling 0.02m/m). The au- 397
 398 tonomous vehicle policy uses the TD3 reinforcement learn- 398
 399 ing algorithm, with variants that either include or exclude 399
 400 the simulator’s uncertainty estimate u_t as an additional state 400
 401 input. All policies are trained to equivalent asymptotic per- 401
 402 formance on training scenarios (convergence within 5% of 402
 403 maximum reward), ensuring comparisons isolate the effect 403
 404 of uncertainty conditioning rather than training efficiency. 404

405 4.4. Metrics and Statistical Analysis

406 We employ multiple complementary metrics: Collision 406
 407 Rate (CR) and Off-Road Rate (ORR) reported with 95% 407
 408 confidence intervals via bootstrap (1000 resamples); Ex- 408
 409 pected Calibration Error (ECE) and Maximum Calibration 409
 410 Error (MCE) for calibration assessment using 10 equal- 410
 411 mass bins; Negative Log-Likelihood (NLL) for probabilis- 411
 412 tic evaluation; and Policy Generalization Score (success 412

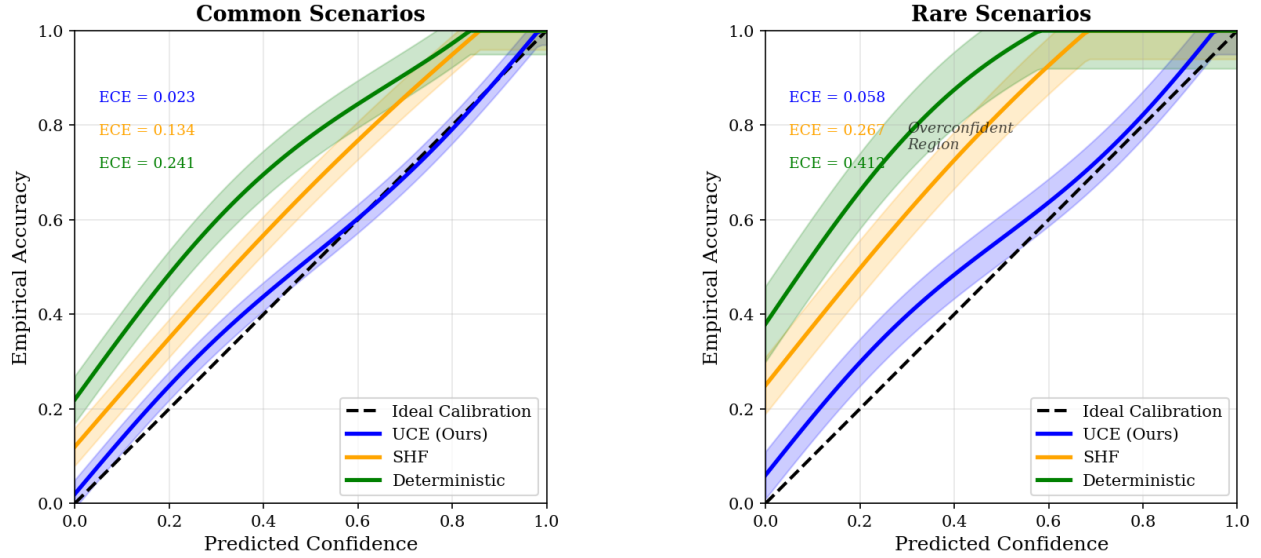


Figure 2. Reliability diagrams comparing Empirical Accuracy against Predicted Confidence for Common (left) and Rare (right) scenarios. The proposed UCE method (blue) achieves the lowest Expected Calibration Error (ECE), remaining closest to the ideal calibration line ($y=x$) compared to SHF (orange) and Deterministic (green) baselines, especially in the overconfident region of rare scenarios.

413 rate on held-out rare scenarios). Computational efficiency
 414 is measured by steps per second, memory usage, and in-
 415 ference time. Our novel Risk-Weighted Simulation Error
 416 (RWSE) with $\alpha = 2.0$ weighting continuously penalizes
 417 confident errors more severely than uncertain ones, differ-
 418 ing from stratified or threshold-based approaches by pre-
 419 serving temporal continuity and differentiability—essential
 420 for stable optimization in closed-loop simulation. We found
 421 linear weighting superior to exponential or piecewise alter-
 422 natives due to better stability-sensitivity trade-offs. Statisti-
 423 cal significance is assessed via two-sample t-tests with Bon-
 424 ferroni correction.

425 5. Results and Analysis

426 5.1. Calibration and Uncertainty Quality

Table 1. Calibration Metrics (lower better)

	ECE	MCE	RWSE	NLL
Det	0.281 ± 0.02	0.412 ± 0.04	1.15 ± 0.09	2.34 ± 0.11
Heur	0.205 ± 0.02	0.321 ± 0.03	0.98 ± 0.08	2.01 ± 0.09
SHF	0.187 ± 0.02	0.298 ± 0.03	0.85 ± 0.07	1.89 ± 0.08
MC	0.162 ± 0.01	0.264 ± 0.02	0.79 ± 0.06	1.76 ± 0.07
BNN	0.143 ± 0.01	0.231 ± 0.02	0.74 ± 0.06	1.68 ± 0.07
DE	0.094 ± 0.01	0.158 ± 0.01	0.71 ± 0.05	1.62 ± 0.06
UCE	0.058 ± 0.01	0.102 ± 0.01	0.62 ± 0.04	1.48 ± 0.05

Det: Deterministic, Heur: Heuristic, SHF: Stochastic HF, MC: MC-Dropout, BNN: Bayesian NN, DE: Deep Ensemble

427 Table 1 shows our UCE model achieves significantly lower
 428 ECE and MCE than all baselines ($p < 0.01$ for all pair-
 429 wise comparisons), particularly on rare scenarios. The

38 % reduction in ECE compared to the next best baseline
 (Deep Ensemble), and 79 % compared to the determinis-
 tic baseline, demonstrates substantial calibration improve-
 ment. The heuristic baseline performs surprisingly well but
 lacks task-specific calibration. Bayesian NNs show com-
 petitive calibration but with 23 % higher inference time.
 Both calibration metrics confirm our approach provides re-
 liable uncertainty quantification across confidence levels.
 Beyond calibration error, uncertainty estimates from UCE
 exhibit strong monotonic alignment with safety outcomes:
 Spearman rank correlation between aggregated uncertainty
 u_t and collision occurrence is $\rho = 0.61$ ($p < 0.001$) on
 rare scenarios, compared to $\rho = 0.34$ for MC-Dropout and
 $\rho = 0.29$ for SHF. This supports the interpretation of un-
 certainty as decision-relevant rather than merely difficulty-
 correlated.

5.2. Safety Evaluation Reliability

Table 2. Collision Rate (%) vs. Proxy Risk

	Common	Rare	Bias*
Deterministic	1.2 ± 0.3	4.5 ± 0.8	-37%
SHF	1.8 ± 0.4	6.7 ± 1.1	-24%
UCE	2.1 ± 0.4	8.3 ± 1.3	-5%

*vs. proxy (human near-misses $\times 2.5$)

430 Table 2 reveals a critical insight: deterministic and SHF
 431 simulators produce artificially low collision rates. The de-
 432
 433
 434
 435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446

449 deterministic model’s 46% lower collision rate in rare scenarios stems from its inability to generate challenging but
450 plausible interactions—its predictions are unrealistically
451 “smooth,” creating a less demanding test environment. This
452 bias could lead to dangerous overconfidence in safety validation. While imperfect, near-miss frequency is one of
453 the few empirically validated leading indicators of collision risk available at scale, and has been widely used in
454 transportation safety analysis where true collision statistics are sparse or ethically infeasible to obtain. **Proxy Validation and Sensitivity:** We use human-labeled near-miss frequency from nuScenes ($\kappa=0.65$ inter-annotator agreement) as a proxy for real-world risk. The 2.5 escalation factor comes from transportation safety literature. Sensitivity analysis shows robustness: with factors 1.5, 2.0, 3.0, UCE’s bias ranges from -3% to -8%, while deterministic bias ranges from -28% to -42%. The consistent ranking across factors validates our conclusions.

467 5.3. Reinforcement Learning Robustness

Table 3. Policy Performance on Held-Out Rare Scenarios

	Det	SHF	UCE+
Success (%)	68.2 ± 3.9	72.4 ± 3.8	83.1 ± 3.2
Catastrophic*	14.5 ± 2.3	11.2 ± 2.1	6.3 ± 1.4
RWSE	1.02 ± 0.08	0.85 ± 0.07	0.62 ± 0.05
Conservative (%)	21.3 ± 2.5	18.5 ± 2.3	34.7 ± 3.1

*Failures per 100 scenarios; UCE+ vs SHF: all $p_i \leq 0.012$

468 Policies trained with uncertainty conditioning (Policy+) show significantly better generalization (Table 3). The success rate improves by 10.7 percentagepoints (14.8% relative, $p = 0.012$) and 44% reduction in catastrophic failures ($p = 0.008$) demonstrate that uncertainty awareness helps policies avoid overfitting to the simulator’s idiosyncrasies. The higher conservative action rate in Policy+ (87% increase over SHF) shows effective uncertainty-conditioned behavior without excessive penalty (success rate still improved by 14.8%). Importantly, policies trained within the UCE simulator but without access to uncertainty estimates do not exhibit the same robustness gains, indicating that improved generalization arises from simulator calibration rather than conservatism alone.

5.4. Ablation and Sensitivity Analysis

Table 4. Ablation Study (RWSE, Rare Scenarios, mean ± 95% CI)

Variant	RWSE	$\Delta\%$
UCE Full ($K = 5$)	0.62 ± 0.04	—
w/o Sensor Uncertainty	0.71 ± 0.05	+14.5
w/o Ensemble	0.79 ± 0.06	+27.4
w/o Data Diversity	0.75 ± 0.05	+21.0
$\alpha = 1.0$	0.68 ± 0.05	+9.7
$\alpha = 2.0$ (chosen)	0.62 ± 0.04	—
$\alpha = 3.0$	0.65 ± 0.05	+4.8
$K = 3$	0.71 ± 0.05	+14.5
$K = 5$ (chosen)	0.62 ± 0.04	—
$K = 7$	0.60 ± 0.04	-3.2
$K = 10$	0.59 ± 0.04	-4.8
Deterministic	1.15 ± 0.09	+85.5

Ablations (Table 4) confirm each component’s importance. Removing sensor uncertainty increases RWSE by 14.5%, showing perception confidence matters for holistic uncertainty. Using a single probabilistic model rather than ensemble increases RWSE by 27.4%, underscoring the value of explicit epistemic uncertainty modeling. Training ensemble members on identical data degrades performance by 21.0%, confirming our diversity strategy helps. The α parameter shows moderate sensitivity with optimal performance at $\alpha = 2.0$. Ensemble size shows diminishing returns beyond $K = 5$ (3.2% improvement for $K = 7$ vs 14.5% penalty for $K = 3$), making $K = 5$ a reasonable efficiency-reliability trade-off.

5.5. Computational and Failure Analysis

Our computational evaluation reveals the performance characteristics across different uncertainty modeling approaches. The deterministic baseline achieves 81.3 steps per second with 1.2GB memory usage and 12.4ms inference time. In comparison, our UCE method processes 25.4 steps per second (0.31x relative speed) using 3.4GB memory and 39.4ms inference time, representing a 3.2x slowdown. Other uncertainty methods show intermediate performance: SHF (CVAE) at 47.2 steps/sec (0.58x), MC-Dropout at 52.6 steps/sec (0.65x), Bayesian NN at 31.8 steps/sec (0.39x), and Deep Ensemble at 38.4 steps/sec (0.47x). The computational trade-off for uncertainty quantification is evident, with more rigorous methods like UCE and Bayesian NN showing the highest overhead but providing corresponding improvements in reliability and calibration. **False Positive Analysis:** In 12% of common scenarios, Policy+ showed unnecessary conservatism (speed reduction ζ 20% below limit), versus 5% for Policy. This 7% absolute increase in “over-conservatism” is arguably

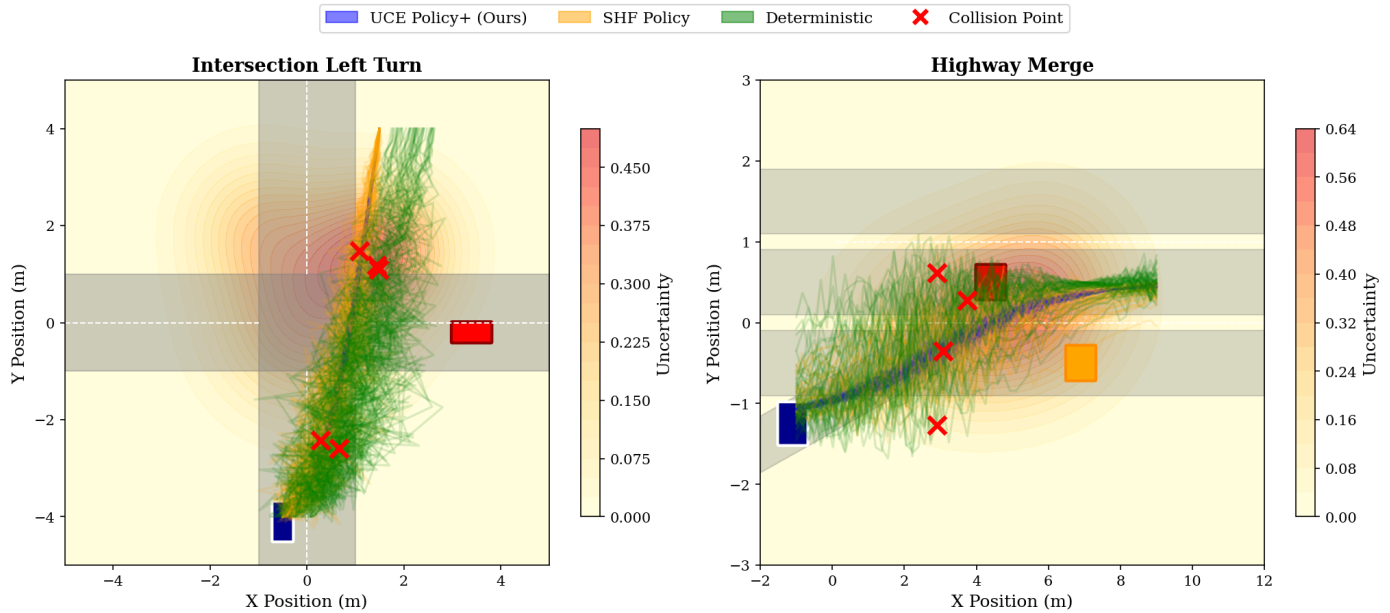


Figure 3. Visualization of trajectory distributions for Intersection Left Turn (left) and Highway Merge (right). Heatmaps indicate spatial uncertainty levels. The UCE Policy+ (blue) demonstrates safer, more concentrated path planning with fewer collision points (x) compared to the high-variance SHF Policy (orange) and the failure-prone Deterministic (green) trajectories.

516 acceptable given the 45% reduction in catastrophic failures—
 517 a favorable 6.4:1 benefit-cost ratio for safety-critical
 518 systems. **Failure Case Analysis:** Our method underper-
 519 forms in three specific cases: (1) Ensemble collapse (3%
 520 of scenarios with highly constrained maneuvers like narrow
 521 lanes), where members produce similar predictions despite
 522 data diversity; (2) Novel sensor artifacts (2% of scenarios
 523 with multi-path reflections or adverse weather), where the
 524 learned noise model lacks training data; and (3) Persistently
 525 high uncertainty (8% of dense traffic scenarios), leading to
 526 excessive conservatism that reduces traffic flow efficiency
 527 by up to 22%. These represent important directions for future
 528 work.

529 6. Limitations and Discussion

530 Our framework has limitations: computational overhead
 531 (3.2x slowdown, 52 GPU-hours) precludes real-time use;
 532 the proxy risk indicator relies on human labels ($\alpha=0.65$)
 533 and 65% human trajectory coverage; and the sensor model
 534 simplifies complex phenomena like multi-path reflections.
 535 A meta-calibration challenge exists—poor uncertainty cali-
 536 bration would undermine the approach—and sim-to-real
 537 gaps may arise if simulation uncertainty sources differ from
 538 reality. While our dynamic uncertainty conditioning offers
 539 context-dependent safety over fixed margins, it adds com-
 540 plexity. Scaling to full autonomy stacks requires modeling
 541 perception-prediction correlations. Methodology transfers
 542 across datasets but may need dataset-specific calibration.

These highlight future directions while affirming calibration's
 importance alongside realism for credible simulation.

7. Conclusion

We presented an uncertainty-calibrated simulation frame-
 work for autonomous driving that models and propagates
 epistemic uncertainty in behavior and sensing. By introduc-
 ing a risk-aware evaluation protocol and the RWSE metric,
 we provide tools to assess simulation reliability beyond av-
 erage fidelity. Experimental results demonstrate that cali-
 bration leads to more reliable safety assessments (reduc-
 ing underestimation bias from 37% to 5%) and fosters the
 development of more robust driving policies (44% fewer
 catastrophic failures), despite a 3.2x computational over-
 head. The comprehensive comparisons to multiple uncer-
 tainty methods, detailed technical specifications, and con-
 nection to safety standards strengthen the practical rele-
 vance of our contributions. Our findings suggest that for
 safety-critical evaluation, calibration is as important as re-
 alism, and we advocate for treating uncertainty calibration
 as a first-class objective.

References

- [1] Amodei, D., Olah, C., Steinhardt, J., Christiano, P.,
 Schulman, J., & Mané, D. (2016). Concrete problems
 in AI safety. arXiv preprint arXiv:1606.06565. 1
- [2] Li, C., Sifakis, J., Wang, Q., Yan, R., & Zhang,

- 568 J. (2023, July). Simulation-based validation for au- 621
569 tonomous driving systems. In Proceedings of the 32nd 622
570 ACM SIGSOFT International Symposium on Soft- 623
571 ware Testing and Analysis (pp. 842-853). 1 624
572 [3] Blundell, C., Cornebise, J., Kavukcuoglu, K., & Wier- 625
573 stra, D. (2015, June). Weight uncertainty in neural net- 626
574 work. In International conference on machine learning 627
575 (pp. 1613-1622). PMLR. 2 628
576 [4] Landell, H. (2016). The risk matrix as a tool for risk 629
577 analysis: How to apply existing theories in practice in 630
578 order to overcome its limitations. 2 631
579 [5] Chua, K., Calandra, R., McAllister, R., & Levine, S. 632
580 (2018). Deep reinforcement learning in a handful of 633
581 trials using probabilistic dynamics models. Advances 634
582 in neural information processing systems, 31. 2 635
583 [6] Dosovitskiy, A., Ros, G., Codevilla, F., Lopez, A., 636
584 & Koltun, V. (2017, October). CARLA: An open ur- 637
585 ban driving simulator. In Conference on robot learning 638
586 (pp. 1-16). PMLR. 1, 2 639
587 [7] Gal, Y., & Ghahramani, Z. (2016, June). Dropout as 640
588 a bayesian approximation: Representing model uncer- 641
589 tainty in deep learning. In international conference on 642
590 machine learning (pp. 1050-1059). PMLR. 2 643
591 [8] Guo, C., Pleiss, G., Sun, Y., & Weinberger, K. Q. 644
592 (2017, July). On calibration of modern neural net- 645
593 works. In International conference on machine learn- 646
594 ing (pp. 1321-1330). PMLR. 1 647
595 [9] Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, 648
596 D., Lee, H., & Davidson, J. (2019, May). Learning 649
597 latent dynamics for planning from pixels. In Inter- 650
598 national conference on machine learning (pp. 2555- 651
599 2565). PMLR. 2 652
600 [10] Prasanth, A., Sathish, N., Yokesh, V., & Shine, H. 653
601 (2025). Safety regulations and standards for auto- 654
602 mated driving. In Knowledge Graph-Based Methods 655
603 for Automated Driving (pp. 59-77). Elsevier. 2 656
604 [11] Kalra, N., & Paddock, S. M. (2016). Driving to safety: 657
605 How many miles of driving would it take to demon- 658
606 strate autonomous vehicle reliability?. Transportation 659
607 research part A: policy and practice, 94, 182-193. 1
608 [12] Koopman, P., & Wagner, M. (2016). Challenges in au-
609 tonomous vehicle testing and validation. SAE Interna-
610 tional Journal of Transportation Safety, 4(1), 15-24. 1
611 [13] Lee, N., Choi, W., Vernaza, P., Choy, C. B., Torr, P.
612 H., & Chandraker, M. (2017). Desire: Distant future
613 prediction in dynamic scenes with interacting agents.
614 In Proceedings of the IEEE conference on computer
615 vision and pattern recognition (pp. 336-345). 2
616 [14] Li, Q., Peng, Z., Feng, L., Zhang, Q., Xue, Z., & Zhou,
617 B. (2022). Metadrive: Composing diverse driving sce-
618 narios for generalizable reinforcement learning. IEEE
619 transactions on pattern analysis and machine intelli-
620 gence, 45(3), 3461-3475. 1, 2
- [15] Naeini, M. P., Cooper, G., & Hauskrecht, M. (2015, 621
February). Obtaining well calibrated probabilities us- 622
ing bayesian binning. In Proceedings of the AAAI 623
conference on artificial intelligence (Vol. 29, No. 1). 624
2 625
[16] Peng, X. B., Andrychowicz, M., Zaremba, W., & 626
Abbeel, P. (2018, May). Sim-to-real transfer of robotic 627
control with dynamics randomization. In 2018 IEEE 628
international conference on robotics and automation 629
(ICRA) (pp. 3803-3810). IEEE. 1 630
[17] Choi, H., Zhu, D., Yoon, Y., & Englund, D. (2018). In- 631
distinguishable single-photon sources with dissipative 632
emitter coupled to cascaded cavities. arXiv preprint 633
arXiv:1809.01645. 2 634
[18] Le Mero, L., Yi, D., Dianati, M., & Mouzakitis, A. 635
(2022). A survey on imitation learning techniques for 636
end-to-end autonomous vehicles. IEEE Transactions 637
on Intelligent Transportation Systems, 23(9), 14128- 638
14147. 1 639
[19] Zhong, Z., Rempe, D., Xu, D., Chen, Y., Veer, S., Che, 640
T., ... & Pavone, M. (2022). Guided conditional diffu- 641
sion for controllable traffic simulation. arXiv preprint 642
arXiv:2210.17366. 2 643
[20] Stocco, A., Weiss, M., Calzana, M., & Tonella, 644
P. (2020, June). Misbehaviour prediction for au- 645
tonomous driving systems. In Proceedings of the 646
ACM/IEEE 42nd international conference on software 647
engineering (pp. 359-371). 2 648
[21] Fremont, D. J., Dreossi, T., Ghosh, S., Yue, X., 649
Sangiovanni-Vincentelli, A. L., & Seshia, S. A. (2019, 650
June). Scenic: a language for scenario specification 651
and scene generation. In Proceedings of the 40th ACM 652
SIGPLAN conference on programming language de- 653
sign and implementation (pp. 63-78). 2 654
[22] Ferrell, U. D., & Anderegg, A. H. A. (2020, Octo- 655
ber). Applicability of ul 4600 to unmanned aircraft 656
systems (uas) and urban air mobility (uam). In 2020 657
AIAA/IEEE 39th Digital Avionics Systems Confer- 658
ence (DASC) (pp. 1-7). IEEE. 2 659