

TinyHazardSynth: Industrial Grade Realistic Data Augmentation for Autonomous Driving Using 3D Modeling and Depth-Aware Occlusion Modeling

Anonymous CVPR submission

Paper ID ****

Abstract

001 We introduce *TinyHazardSynth*, a depth-aware synthetic
002 model-auditing pipeline that can be used to stress-testing
003 autonomous-driving detectors by controllably inserting
004 photorealistic small road obstacles—situations rarely cap-
005 tured by public datasets—into real dash-cam videos. NeRF-
006 reconstructed assets are rendered with exact camera in-
007 trinsics; precise occlusion handling arise from a two-stage
008 fitting fusion of sparse LiDAR and Depth-Anything depth
009 that converts relative estimates to metric scale. A Mask-
010 Former semantic prior prevents ground and road clipping,
011 and modular fog/shadow layers vary visibility to probe ro-
012 bustness. The fully-automated factory produces thousands
013 of labelled frames and can wrap around any perception
014 stack. On inserted tiny obstacles, the in-house detector
015 achieved a recall of only 29.4%, indicating a high miss rate
016 for rare, small-scale hazards. And, after targeted retraining
017 on our clips, we lifted accuracy by 1.4pp, demonstrating
018 the pipeline’s value for safety-critical model assessment.
019 We leave as future work an investigation of how the same
020 controllable-insertion pipeline adapts to other rare hazard
021 types (e.g., deformable debris or transparent objects) and
022 to public datasets such as nuScenes and KITTI.

023 1. Introduction

024 Small static hazards—such as fist-sized rocks, dropped
025 cargo, or road-surface debris—pose a disproportionate
026 safety risk for autonomous vehicles, yet they appear
027 only rarely in public driving datasets. Purely simulator-
028 generated clips like Airsim and Carla [3, 6] help to enlarge
029 coverage, but the gap in texture realism, motion blur and
030 camera ego-motion still limits their usefulness once a model
031 is deployed on real dash-cam streams, which is the data used
032 for real world driving task.

033 We explore a complementary route: *in-video* insertion
034 of realistic tiny obstacles directly into ordinary dash-cam
035 footage. Our system, **TinyHazardSynth**, (i) captures ob-

stacle assets with NeRF [5], (ii) renders them per frame using
exact camera intrinsics, and (iii) achieves pixel-accurate
occlusions by fusing sparse LiDAR with Depth-Anything
predictions [8] to fit the relationship between relative depth
and metric depth. A lightweight *MaskFormer* prior [2] pre-
serves road and vehicle boundaries to help refine further,
while optional fog and shadow layers vary visibility condi-
tions. The pipeline produces more than thousands labelled
frames and can feed either regular data augmentation or tar-
geted robustness checks of existing detectors.

Contributions

1. A controllable dash-cam video insertion workflow fo-
cused on small road hazards, bridging the realism gap
left by simulator-only data.
2. A two-stage fusion of LiDAR and monocular depth that
converts relative estimates into metric scale, giving sta-
ble occlusions under fast camera motion.
3. An industrial-scale implementation that can synthesize
thousands of frames; we intend to share a streamlined
code snapshot after acceptance so that others can inspect
our pipeline logic and replicate the core steps on their
own inputs.

2. Literature Review

Synthetic data for autonomous driving. Large-scale
simulators such as CARLA [3], AIRSIM [6] render traf-
fic scenes with controllable lighting, weather and sensor
rigs, enabling low-cost data generation and *in-sim* per-
turbations for robustness studies. Yet a non-trivial do-
main gap—texture realism, motion blur and long-tail oc-
clusions—still limits transferable performance on real
dash-cam footage. Consequently, there is a gap demand for
production pipelines that augment real videos in real world
dash camera captured scenarios.

Object insertion in dash-cam videos. A straight-
forward approach would be to composite 2-D cut-outs.
However, this breaks down under strong parallax and, to
the best of our knowledge, lacks any pipeline capable of
handling complex viewpoint geometry or batch-level in-

074 tegration across video frames. NeRF-based approaches
075 [5] allow high-fidelity asset capture. In practice, exist-
076 ing insertion tools assume either a fixed camera or pre-
077 computed dense depth, conditions rarely met by long dash-
078 cam sequences with sparse LiDAR. Few papers tackle *small*
079 ground obstacles, whose limited pixel footprint amplifies
080 any depth or mask error.

081 **Our position.** We bridge these gaps by (i) inserting
082 NeRF-reconstructed tiny hazards into real videos using
083 a LiDAR + Depth-Anything metric fusion, (ii) preserving
084 road semantics via *MaskFormer* thereby further refining
085 synthesis, and (iii) modulating visibility with controllable
086 shadows—with additional diffusion-based weather effects
087 left as future work—so that quantitative audits remain en-
088 tirely within the real-image domain. The resulting clips
089 both enrich training data *and* serve as a fine-grained, dash-
090 cam-authentic testbed for obstacle-detection auditing.

091 3. Methodology

092 3.1. NeRF-based 3D Reconstruction for Obstacle 093 Cutout Images Gathering

094 To generate accurate visual representations of small
095 ground obstacles, we first employ Neural Radiance Fields
096 (NeRF) [5] to reconstruct detailed 3D models in the .ply
097 format from multiple views. Utilizing intrinsic and extrinsic
098 camera parameters at various timestamps, we render obsta-
099 cle images from precise virtual camera positions. This ap-
100 proach allows us to generate continuous and realistic obsta-
101 cle representations for subsequent insertion into dash cam
102 videos. In practice, obstacle cutout images can be cut out
103 with other methods or even manually.

104 3.2. Accurate Depth Estimation and Occlusion 105 Handling

106 A critical challenge in synthesizing realistic dash cam data
107 arises when dynamic objects, such as animals, suddenly
108 emerge from behind obstacles, creating complex occlusion
109 relationships. Traditional depth estimation models, includ-
110 ing Depth Anything V2 [8], output relative depth maps (0–
111 255 values) rather than absolute distances, which is not
112 compatible with LiDAR data. To bridge this gap, we in-
113 corporate LiDAR point clouds as ground truth data to assist
114 depth conversion. However, projecting point clouds onto
115 images typically results in “edge outlining effects,” where
116 the projected cloud is larger than the actual object in the
117 image, producing numerous discrete, noisy values at object
118 boundaries.

119 To address this, we introduce an innovative two-stage
120 curve fitting approach using *Scipy* [7] optimization. Ini-
121 tially, we apply a preliminary fit to the LiDAR-based depth
122 data and Depth Anything V2’s relative depth output, defined
123 by a function of the form: This first fitting pass is specifi-

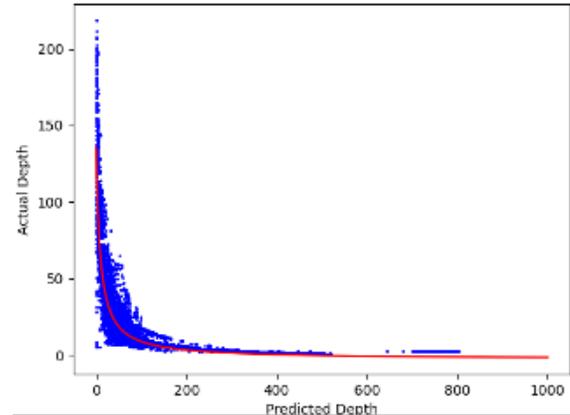


Figure 1. Fitting curve between predicted (relative) and actual depth. After outlier filtering, the final curve (in the form of $y = a / (x + b) + c$) aligns well with ground truth.

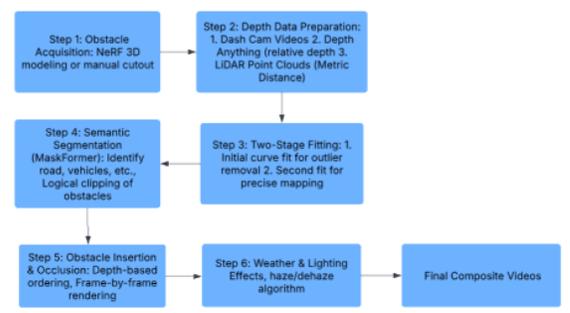


Figure 2. the whole flowchart of the pipeline

124 cally utilized to filter outliers by removing points exceed-
125 ing a defined margin threshold from the initial fit. Subse-
126 quently, we perform a second fit using the filtered dataset to
127 accurately determine the final mapping parameters between
128 relative and absolute depth, to reduce noise and enhancing
129 occlusion realism.

130 3.3. Semantic Segmentation for Clip Handling

131 Another practical challenge encountered is the “clipping”
132 effect due to camera perspective, where objects appear
133 partially embedded within road surfaces or vehicle struc-
134 tures, resulting in unrealistic visualization. This arises
135 from perspective distortion where parallel lines converge
136 in camera projection. To mitigate this issue, we employ
137 the *MaskFormer* [2] semantic segmentation model (trained
138 on the ADE20K dataset). *MaskFormer* enables the pre-
139 cise identification of road and vehicle hood areas, allowing
140 logic-based controls to manage object visibility effectively.
141 Specifically, we enforce rules to hide or show small objects
142 visually positioned beneath vehicle hoods or incorrectly in-
143 tersecting road surfaces, thus maintaining the semantic in-
144 tegrity and realism of synthesized videos.

	Accuracy	Recall
Baseline (without synthetic clips)	86.9 %	29.4 %
+ TinyHazardSynth clips	88.3 %	36.6 %

Table 1. Internal obstacle-detection results. Synthetic tiny-hazard clips lift accuracy by 1.4pp and recall by 7–8pp (24.5 % relative).

145 4. Experimental Validation

146 The goal of this section is to examine whether the clips
147 produced by our pipeline translate into measurable per-
148 formance gains on real-world autonomous-driving tasks.
149 While a full public benchmark is in preparation, we report
150 the first numbers from an internal test bed provided by an
151 industry partner.

152 The experiment uses proprietary dash-cam footage into
153 which we insert *metal plates*—a canonical small, rigid haz-
154 ard that is hard to detect because of its thin profile and
155 low reflectivity. After one round of fine-tuning on the aug-
156 mented set, the detector achieves the improvements listed
157 in Table 1. Although the study covers a single class
158 and a limited geography, the gain suggests that depth-
159 and shadow-aware synthesis reveals failure cases other-
160 wise under-represented in natural data and thus serves as
161 a lightweight stress test of model robustness.



Figure 3. Frame pair before (up) and after (bottom) inserting a partially obscured fallen barricade at the hood–road junction.

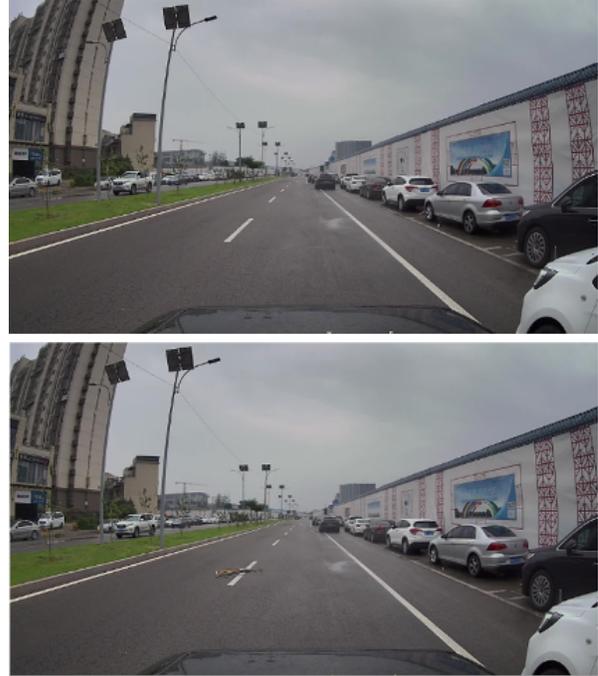


Figure 4. Another example with additional synthetic shadowing to match scene illumination.

Future work may include extending the pipeline to ad- 162
ditional object classes and investigating public release of a 163
small validation subset. 164

165 5. Conclusion and Future Work

166 We have presented a depth-aware pipeline that inserts pho-
167 torealistic tiny obstacles into dash-cam videos, combining
168 NeRF asset capture with a LiDAR + Depth-Anything fitting
169 scheme and a *MaskFormer* prior for clean semantic bound-
170 aries. On an internal test set the synthetic clips improved
171 iron-plate detection by 1.4 pp in accuracy and about 25% in
172 recall, indicating that such rare-event enrichment can bene-
173 fit production models.

174 Future directions

- 175 **1. Wider sensor settings.** The modular design makes it fea-
176 sible to plug in alternative depth sources—e.g. stereo or
177 monocular SfM—but the impact of domain shift remains
178 to be quantified.
- 179 **2. Richer environment effects.** Adding physics-based vol-
180 umetric fog or lightweight generative post-processing
181 could further close the realism gap under adverse
182 weather.
- 183 **3. Broader evaluation.** Extending the study to deformable
184 or transparent objects, and to public datasets such as
185 nuScenes [1] and KITTI [4], would clarify how well the
186 approach generalizes beyond metal plates.

All experiments in this paper rely on proprietary video due 187

188 to incorporate restrictions; public benchmarking is left to fu-
189 ture collaborative work.

190 References

- 191 [1] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora,
192 Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Gi-
193 ancarno Baldan, and Oscar Beijbom. nuscenes: A multimodal
194 dataset for autonomous driving. In *Conference on Computer
195 Vision and Pattern Recognition (CVPR)*, pages 11621–11631,
196 2020. 3
- 197 [2] Bowen Cheng, Alexander G. Schwing, and Alexander Kir-
198 illov. Per-pixel classification is not all you need for semantic
199 segmentation. *NeurIPS*, 2021. arXiv:2107.06278. 1, 2
- 200 [3] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio
201 López, and Vladlen Koltun. Carla: An open urban driving
202 simulator. In *CoRL*, pages 1–16, 2017. arXiv:1711.03938. 1
- 203 [4] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we
204 ready for autonomous driving? the kitti vision benchmark
205 suite. In *Conference on Computer Vision and Pattern Recog-
206 nition (CVPR)*, pages 3354–3361, 2012. 3
- 207 [5] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik,
208 Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf:
209 Representing scenes as neural radiance fields for view synthe-
210 sis. In *ECCV*, pages 405–421. Springer, Cham, 2020. 1, 2
- 211 [6] Shital Shah, Debadeepta Dey, Chris Lovett, and Ashish
212 Kapoor. Airsim: High-fidelity visual and physical simulation
213 for autonomous vehicles. In *Field and Service Robotics*, pages
214 621–635. Springer, Cham, 2018. 1
- 215 [7] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt
216 Haberland, Tyler Reddy, David Cournapeau, ..., and SciPy 1.0
217 Contributors. Scipy 1.0: Fundamental algorithms for scien-
218 tific computing in python. *Nature Methods*, 17(3):261–272,
219 2020. 2
- 220 [8] L. Yang, B. Kang, Z. Huang, Z. Zhao, X. Xu, J. Feng, and H.
221 Zhao. Depth anything v2. *NeurIPS*, 2024. arXiv:2406.09414.
222 1, 2