# RenderBender: A Survey on Adversarial Attacks **Using Differentiable Rendering**

Matthew Hull<sup>1</sup>, Haoran Wang<sup>1</sup>, Matthew Lau<sup>1</sup>, Alec Helbling<sup>1</sup>,

Mansi Phute<sup>1</sup>, Chao Zhang<sup>1</sup>, Zsolt Kira<sup>1</sup>, Willian Lunardi<sup>2</sup>, Martin Andreoni<sup>2</sup>, Wenke Lee<sup>1</sup> and Polo Chau<sup>1</sup>

<sup>1</sup>Georgia Tech, <sup>2</sup>Technology Innovation Institute

<sup>1</sup>{matthewhull,haoran.wang, mlau40, alechelbling, mphute6, chaozhang, zkira, wenke, polo}@gatech.edu, <sup>2</sup>{willian.lunardi, martin.andreoni}@tii.ae

### Abstract

Differentiable rendering techniques like 1 Gaussian Splatting and Neural Radiance 2 Fields have become powerful tools for 3 generating high-fidelity models of 3D ob-4 jects and scenes. Their ability to pro-5 duce both physically plausible and differ-6 entiable models of scenes are key ingredi-7 ent needed to produce physically plausible 8 adversarial attacks on DNNs. However, 9 the adversarial machine learning commu-10 nity has yet to fully explore these capabil-11 ities, partly due to differing attack goals 12 (e.g., misclassification, misdetection) and 13 a wide range of possible scene manipu-14 lations used to achieve them (e.g., alter 15 texture, mesh). This survey contributes a 16 framework that unifies diverse goals and 17 tasks, facilitating easy comparison of ex-18 isting work, identifying research gaps, and 19 highlighting future directions-ranging 20 from expanding attack goals and tasks to 21 account for new modalities, state-of-the-22 art models, tools, and pipelines, to under-23 scoring the importance of studying real-24 world threats in complex scenes. 25

#### 1 Introduction 26

Differentiable rendering has emerged as a power-27 ful tool for solving inverse problems in vision and 28 graphics by enabling gradient propagation through 29 the rendering process. Recent methods like Neural 30 Radiance Fields (NeRF) [Mildenhall et al., 2020] 31

and 3D Gaussian Splatting [Kerbl et al., 2023] en-32 able novel view synthesis from limited images to re-33 construct 3D models or scenes. These advancements 34 have spurred open-source tools, such as PyTorch3D<sup>1</sup> 35 and user-friendly platforms<sup>2</sup> that allow creating tex-36 tured 3D models from photos. 37

Differentiable rendering has also exposed vulnera-38 bilities in DNNs by enabling adversarial attacks. Ad-39 versaries exploit DNN gradients to optimize inputs, 40 training, or outputs for malicious purposes, leading 41 to misclassifications in systems such as stop signs in 42 cars, LiDAR [Cao et al., 2019], facial recognition, 43 and 3D models [Xiao et al., 2019]. Similarly, differ-44 entiable rendering allows attackers to optimize 3D 45 scene parameters (objects, materials, lighting) via 46 loss gradients. Research on differentiable rendering-47 based attacks is scattered across: 48

- 1. Attack goals: e.g., inducing misclassifications or 49 motion/depth errors; 50
- 2. Attackable components: e.g., preprocessing 51 steps or during inference; 52
- 3. Scene manipulation: e.g., targeting texture, ge-53 ometry, or combinations thereof. 54

In other words, progress in adversarial attacks us-55 ing differentiable rendering has been made, but sys-56 tematic comparisons, summaries of strengths, and 57 research gap identification remain challenging. Fig-58 ure 1 shows how our survey addresses this gap by 59 organizing tasks like texture manipulation, illumina-60 tion changes, and 3D mesh alterations, emphasizing 61 both techniques and their potential exploitation by 62 adversaries. 63

<sup>&</sup>lt;sup>1</sup>https://pytorch3d.org <sup>2</sup>https://poly.cam



### Unified Goals & Tasks Framework for Adversarial Attacks with Differentiable Rendering

Fig. 1: Visual overview of our unifying survey framework that, by unifying the diverse goals and tasks in identifying attackable components and manipulating scene representations, enables systematic summarization and comparison with existing differentiable rendering related adversarial attack research.

## 64 1.1 Related Survey and Methodology

This is the first survey to focus on task-based dif-65 ferentiable rendering capabilities for 3D adversarial 66 attacks. Existing work separates differentiable ren-67 dering and adversarial research. Kato et al. briefly 68 mention adversarial applications in their differen-69 tiable rendering survey but lack attack details [Kato 70 et al., 2020]. Since then, NeRF and 3D Gaus-71 sian Splatting have gained prominence, requiring 72 discussion. Surveys on NeRF [Xie et al., 2022; 73 Tewari et al., 2022; Gao et al., 2023; Mittal, 2024] 74 and 3D Gaussian Splatting [Chen and Wang., 2024; 75 Tosi et al., 2024] do not address adversarial use. Ex-76 isting adversarial attack surveys cover 2D/3D models 77 [Li et al., 2024b], robustness and defenses [Miller et 78 al., 2020], or image classification [Machado et al., 79 2023] but omit differentiable rendering. 80

We reviewed 28 works from top venues in com-81 puter vision, ML, and graphics, covering differen-82 83 tiable rendering methods (e.g., NeRF, 3D Gaussian Splatting) and their use in adversarial attacks. Us-84 ing a task-based framework, we categorized attacker 85 goals-texture, illumination, and mesh manipula-86 tion-to clarify methodologies and vulnerabilities. 87 88 As a newer field, differentiable rendering research 89 began in 2014, with adversarial applications emerging in 2019. 90

## **1.2** Contributions

C1. We present a comprehensive, attacker-task 92 guided survey on adversarial attacks using differ-93 entiable rendering, incorporating a use-inspired 94 approach (Fig. 1). Our framework positions each 95 work by attacker objectives and differentiable ren-96 dering techniques, defining the attack surface based 97 on feasible scene manipulations of the scene repre-98 sentations (Sec. 3). 99

91

- Our methodology links goals to tasks, providing a structured comparison of works and identifying research gaps.
   100
- Table 1 explains differentiable rendering's role in 103 attacks, relevant methods, and current strengths 104 and limitations. 105

C2. We provide comprehensive categorizations 106 of attack methods, highlighting their impact and 107 real-world implications (Sec. 3). We show a "Tar-108 get List" of attacked models, including object de-109 tection, image classification, and others along with 110 attacker access levels (Table 2). These resources en-111 able researchers to build on existing work, compare 112 outcomes, and develop new techniques to address 113 adversarial threats using differentiable rendering. 114

**C3. We identify key future research directions to address the growing threat of adversarial attacks** (Sec. 6). Priorities include developing robust defenses, exploring novel attack strategies, and in-**118** 

vestigating the physical plausibility of these attacks. 119

#### 2 **Attacker Goals** 120

To better understand how differentiable rendering 121 122 is used in adversarial attacks, we describe an attacker using the threat model concept to delineate 123 their goals, capabilities, and knowledge in the con-124 text of the attack they wish to carry out [Li et al., 125 2024b]. Using an attacker-task guided perspec-126 tive, we connect the attacker goals to required tasks 127 and sub-tasks (Sec. 3) that are used in differen-128 tiable rendering attacks. Attacker goals encompass 129 any threat affecting the integrity of a DNN's in-130 tended task Papernot et al., 2016b; Li et al., 2024b; 131 Wiyatno *et al.*, 2019]. In this survey, we identify 132 five attacker goals that are used in attacks on deep 133 learning models using differentiable rendering: 134

2.1. Misclassification - the model predicts an incor-135 rect class (untargeted) or a specified incorrect class 136 (targeted) [Papernot et al., 2016a]. 137

2.2. Misdetection - manifested as various errors: 138 nothing is detected (evasion), improper bounding 139 box localization, duplicate detections, or detecting 140 background as an object, or combinations thereof 141 [Bolya *et al.*, 2020]. 142

2.3. Reduce Confidence - the target class is not pre-143 dicted with high confidence [Papernot et al., 2016a]. 144 145

**2.4.** Misestimate Motion - the model misestimates 146 the motion of objects in the scene caused by adversar-147 ial movements or objects [Schmalfuss et al., 2023]. 148 149

**2.5. Misestimate Depth** - the model misestimates 150 151 depth, affecting the model's ability to perceive distances. [Zheng et al., 2024]. 152

153 In Table 1, we categorized our 28 survey papers 154 as S=Survey (3), M=Metrics (1), or A=Attack (24). Of our 24 Attack papers, we found that 18 works 155 chose goals of inducing misclassifications, 17 in-156 duced misdetections, and 10 induced reduction in 157 model confidence while only 1 work each pursued 158 attack goals of misestimation of motion or depth. 159

#### **Identify Attackable Components** 3 160

To achieve the attacker goals in Sec. 2, one must 161 162 identify which components can be manipulated by 163 analyzing the attack surface (Sec. 3) and understanding 3D scene representations. 164

### Attack Surface

The attack surface includes all data processing stages 166 [Papernot et al., 2016b] in Fig. 1, from sensor inputs 167 and pre-processing to model inference and output 168 actions. In differentiable rendering, this surface ex-169 tends to the renderer and scene representation, giving 170 adversaries multiple potential entry points. For in-171 stance, a robot scanning its 3D environment has: 172 • Sensor Inputs (e.g., camera, LiDAR). 173

• **Pre-processing** steps (e.g., generating 2D images 174 or point clouds). 175

• Inference by the DNN model.

• Predictions (e.g., labels, bounding boxes, segmen-177 tation). 178

• Actions or decisions based on model output. 179 An attack's effectiveness hinges on the adversary's 180 access level: white-box  $\circ$ , black-box  $\bullet$ , or combina-181 tion thereof •, classified in Table 2. In differentiable 182 rendering, attackers manipulate scene elements (Sec. 183 3), such as object textures or environmental factors 184 (e.g., adversarial weather [Schmalfuss et al., 2023]) 185 to deceive the DNN. 186

### **Analyze Scene Components**

In differentiable rendering attacks, the 3D scene rep-188 resentation is the main target. We categorize its 189 components under: **Geometry** (explicit or implicit 190 [Mildenhall et al., 2020; Kerbl et al., 2023]), Tex-191 ture (color and reflectance), **Position/Pose** (object 192 location/orientation), **Illumination** (light sources, 193 e.g., sun or lamps), and Sensors (camera/LiDAR 194 properties like resolution or field of view). Identify-195 ing these components helps attackers craft manipu-196 lations that produce realistic, adversarial inputs to 197 DNN models. 198

#### Manipulate 3D Scene 4

Differentiable rendering enables gradient-based ma-200 nipulation of any scene elements. With white-box 201 access to victim models, attackers can use loss gra-202 dients with respect to scene representations to guide 203 such manipulations. This section reviews common 204 manipulations on 3D scene representations, includ-205 ing geometry, texture, pose, illumination, and sen-206 sors. Among the surveyed works, texture attacks are 207 the most prevalent (15), followed by geometry (7), 208 pose (5), illumination (2), and sensors (1). 209

165

176

187

199

	ATTACKER GOALS					REQUIRED TASKS							DON	IAIN	WHERE	
						IDEN	ITIFY	MANIPULATE								
Work	2.1 Misclassification	2.2 Missed Detection	2.3 Reduce Model Confidence	2.4 Misestimation of Motion	2.5 Misestimation of Depth	3.1) Attack Surface	3.2 Analyze Scene Components	4.1 Attack Scene Geometry	4.2 Attack Scene Textures	4.3 Attack Scene Object Pose	4.4 Attack Scene Illumination	4.5 Attack Scene Sensors	5.1 Perform Digital Attack	5.2 Perform Physical Attack	Paper Type	Publication Venue
[Abdelfattah et al., 2021]															A	ICIP CVPB
Bolya et al. 2019															M N	FCCV
Byun et al. 2022															Δ	arXiv
[Cao et al. 2019]									-						Â	arXiv
[Dong et al. 2022]								_							A	NeurIPS
Huang et al., 2024										_					A	CVPR
Li et al., 2024b]															S	ACM CSUR
Leheng et al., 2023															Ă	arXiv
Li <i>et al.</i> , 2024a															Α	arXiv
Liu et al., 2019															Α	ICLR
Machado et al., 2023															S	ACM
[Maesumi et al., 2021]															Α	arXiv
[Meloni <i>et al.</i> , 2021]															Α	ICMLA
Papernot et al., 2016b															Α	EuroS&P
Papernot et al., 2016a															A	EuroS&P
Schmalfuss et al., 2023															A	ICCV
Shahreza and Marcel, 2023															A	TPAMI
Suryanto et al., 2022															A	CVPR
Suryanto et al., 2023															A	ICCV
[Tu <i>et al.</i> , 2021]	_													_	A	CoRL
[Wang et al., 2022]															A	AAAI
[wiyatho et al., 2019]								-							A	
[Ala0 <i>et al.</i> , 2019]															A	TNNLS
[Tuan <i>et al.</i> , 2019]										-	-			-	5	CVPP
[Zeng et al., 2019]															A	CVPR
[Zheng et al., 2024]	-		-												A	ICM
ZHOU <i>et al.</i> , 2024															A	ICIVIL

Table 1: Overview of representative works on adversarial attacks using differentiable rendering methods. Each row is one work; each column corresponds to a required attacker task or goal. A work's relevant goal or task is indicated by a colored cell. S = Survey, M = Metrics, A = Attack.

210 4.1 Attacks on Scene Geometry

Mesh. Attackers use differentiable render-211 ing to generate adversarial meshes by per-212 Multiply turbing vertex positions to minimize the 213 cross-entropy loss towards the target label. The ad-214 versarial meshes are re-rendered as inputs to victim 215 models. Beyond Pixel Norm-Balls [Liu et al., 2019] 216 introduced a differentiable rendering framework for 217 generating adversarial geometry V' by propagating 218 gradients through a rendering pipeline via chain rule: 219

$$V' \leftarrow V - \gamma \frac{\partial C}{\partial I} \frac{\partial I}{\partial N} \frac{\partial N}{\partial V}, \qquad (1)$$

where V are vertex positions, N per-face normals, and  $\gamma$  the attack strength. MeshAdv [Xiao *et al.*, 2019] used Neural Mesh Renderer to perturb vertices, attacking classifiers and object detectors like YOLO-v3. TT3D [Huang *et al.*, 2024] created adversarial geometry via NeRF and marching cubes but faced scalability challenges due to optimization 226 overhead [Tewari *et al.*, 2022]. Distracting Down-227 pour [Schmalfuss *et al.*, 2023] attacked optical flow 228 models by adding scene-specific spatiotemporally 229 consistent particulate geometry (e.g., rain or snow) 230 to create false motion signals in various datasets. 231

Point Cloud. LiDAR-ADV [Cao et al., 2019] 232 used a differentiable LiDAR simulator to per-233 turb point clouds, converting the initially non-234 differentiable features into differentiable ones with 235 smoothing. Two other works perturbed point cloud 236 objects and converted them to textured meshes to 237 target multi-modal systems [Abdelfattah et al., 2021; 238 Tu et al., 2021]. 239

Geometry Post-Processing and Stabilization. 240 Post-perturbation processing maintains realism and 241 avoids topological issues, such as self-intersections 242 or non-manifold meshes. Techniques like Lapla- 243



Table 2: DNNs attacked by differentiable rendering. Each column is one work; each row is a model.

cian smoothing, regularization loss, and Chamfer 244 distance loss ensure realistic and stable adversarial 245 geometry. For instance, Laplacian smoothing min-246 imizes deviations between original and adversarial 247 vertex, while Chamfer distance loss penalizes dis-248 similarities between point clouds P and Q. Depth 249 completion and lighting approximation from Tu et 250 al. [Tu et al., 2021] enhance realism by restricting 251 adversary scale within axis-aligned bounding boxes. 252

253

### 4.2 Attack Scene Texture

Texture adversarial attacks manipulate an 254 object's appearance by perturbing its color, 255 A pattern, or light reflection properties. Us-256 ing differentiable rendering, the model's loss gra-257 dient is used to perturb the texture mappings (e.g., 258 UV maps) via world-aligned methods that optimize 259 2D textures, UV map-based methods that directly 260 optimize 3D textures, and neural-rendered methods 261 that dynamically generate textures from 3D repre-262 sentations [Zhou et al., 2024]. 263

Multi-Object Texture Attacks. Two works ex-264 plore multi-object texture attacks to study transfer-265 ability. Meloni et al. [Meloni et al., 2021] create poi-266 soned data by perturbing texels using a saliency map 267 from a non-differentiable renderer. Byun et al. [Byun 268 et al., 2022] demonstrate transferability by applying 269 2D adversarial textures to various 3D objects, achiev-270 ing successful impersonation and dodging attacks 271 against facial recognition classifiers. 272

Adversarial Camouflage. Adversarial camou-273 flage targets vehicles and humans. For vehicles, FCA 274 [Wang et al., 2022] applied adversarial textures to 275 an Audi e-Tron in CARLA scenes using the Neural 276 Mesh Renderer (NMR), while DTA [Survanto et al., 277 2022] used EoT for texture projections for a Tesla 278 Model 3 and ACTIVE [Survanto *et al.*, 2023] made a 279 further improvement with tri-planar mapping, allow-280 ing complex shapes. Li et al. [Li et al., 2024a] con-281 ducted a flexible physical camouflage attack (FPA) 282 using diffusion models to generate UV-map-based 283 textures, improving the environmental adaptability 284 in neural rendering. RAUCA [Zhou et al., 2024] 285 extended this by incorporating environmental condi-286 tions via an encoder-decoder Environmental Feature 287 Extractor (EFE) for optimized textures. For humans, 288 Maesumi et al. [Maesumi et al., 2021] developed 289 adversarial clothing using UV maps and SMPL mod-290 els, using Blender's subdivision surface modifier to 291 improve texture resolution for more effective attacks. 292

Texture Attacks on Autonomous Driving Sys-293 294 tems. Abdelfattah et al. attacked object textures by treating vertex colors as learnable parameters, reduc-295 ing YOLOv3 detection in cascaded models used in 296 self-driving [Tu et al., 2021]. Adv3D [Leheng et al., 297 2023] used NeRF with semantic branch augmenta-298 tion along with EoT to enhance physical transfer-299 ability and reduce the confidence of LiDAR detec-300 tor, while 3D<sup>2</sup>Fool [Zheng et al., 2024] developed 301 object-agnostic adversarial patches via EoT and tex-302 ture conversion to attack monocular depth estimation 303 304 models.

305 Texture Post-Processing and Stabilization. To enhance the appearance and physical transferabil-306 ity of adversarial textures, many works incorporate 307 post-processing techniques such as hyperparameter 308 tuning, Total Variation (TV) loss, Smooth loss, and 309 Non-Printability Score (NPS). TV loss [Mahendran 310 and Vedaldi, 2015] penalizes differences between ad-311 jacent texture pixels, reducing noise and promoting 312 smoothness: 313

$$TV(\boldsymbol{x}) = \sum_{i,j} \left( (\boldsymbol{x}_{i,j+1} - \boldsymbol{x}_{ij})^2 + (\boldsymbol{x}_{i+1,j} - \boldsymbol{x}_{ij})^2 \right)^{\frac{1}{2}}.$$

NPS [Sharif *et al.*, 2016] assesses the physical printability of textures by evaluating pixel proximity to printable RGB triplets  $P \subset [0, 1]^3$ :

$$NPS(\hat{p}) = \prod_{p \in P} |\hat{p} - p|,$$

where a low score indicates higher printability.
These methods ensure adversarial textures are both
visually plausible and physically realizable.

320 4.3 Attack Scene Illumination

Illumination manipulation in differentiable 321 rendering attacks is underexplored due to 322 Its global impact, but Liu et al. [Liu et al., 323 2019] and Zeng et al. [Zeng et al., 2019] demonstrate 324 325 its potential. Liu et al. used spherical harmonic lighting for global adjustments, optimizing coefficients 326 via the chain rule (Eq. 1) to preserve realism while at-327 tacking DNNs. Zeng et al. manipulated point lights, 328 creating adversarial lighting to mislead DNNs. 329

330 4.4 Attack Scene Sensors

While many attacks test robustness to camera angle changes, Shahreza et al. [Shahreza and Marcel, 2023] directly manipulated camera parameters for attacks. Using 334 NeRF, they optimized camera rotations to find face 335 poses capable of impersonating target identities in fa-336 cial recognition models. Their method reconstructs 337 3D faces from 2D facial templates, enabling practi-338 cal presentation attacks, such as digital screen replay 339 or printed photographs, which can be further ex-340 tended to create wearable face masks for physical 341 impersonation. 342

## 4.5 Attack Scene Object Pose/Translation 343

DNNs are vulnerable to subtle pose or 344 position changes, with tools like 3DB 345 [Leclerc et al., 2022] available for vulnerability ex-346 ploration. Using differentiable rendering, attackers 347 can generate precise object poses or translations to 348 induce misclassification in arbitrary settings. Alcorn 349 et al. [Alcorn et al., 2019] demonstrated that Incep-350 tionv3 misclassifies 97% of the pose space for Im-351 ageNet objects recognized in their canonical poses, 352 with adversarial poses transferring at high rates to 353 AlexNet, ResNet-50, and YOLOv3. Similarly, Zeng 354 et al. [Zeng et al., 2019] demonstrated adversarial 355 poses misleading Visual Question Answering mod-356 els resulting in wrong scene descriptions. ViewFool 357 [Dong et al., 2022] trained NeRF models on 3D ob-358 jects from BlenderKit<sup>3</sup>, sampling 100 images per 359 model, and demonstrated that ViT-B/16 was more 360 robust to pose attacks than ResNet-50. 361

## 5 Digital and Physical Attack Domains 362

In this section, we focus on the challenges of creating and evaluating attacks in both digital and physical domains (see Table 1) and the tools used for attack. 365

366

### 5.1 Attacks in the Digital Domain

Digital attacks often rely on simulations for con-367 trolled testing. However, since differentiable ren-368 derers (e.g., Mitsuba, PyTorch3D) often lack simula-369 tion features, researchers use non-differentiable tools 370 with simulation features instead. RAUCA [Zhou et 371 al., 2024] and FPA [Li et al., 2024a] used Unreal 372 Engine and CARLA, non-differentiable tools that 373 support data capture, diverse scene setups, lighting 374 conditions and self-driving simulations. Two other 375 works produced adversarial textures and meshes us-376 ing PyTorch 3D and then evaluated their robustness 377 within scenes rendered by non-differentiable Blender 378

<sup>&</sup>lt;sup>3</sup>https://www.blenderkit.com/

and Unity tools [Zeng *et al.*, 2019; Meloni *et al.*,
2021]. TT3D [Huang *et al.*, 2024] created attacked
objects using NeRF and then used Blender and Meshlab for testing cross-render transferability.

### **383** 5.2 Attacks in the Physical Domain

Implementing real-world adversarial attacks poses 384 several challenges, especially when manufacturing 385 adversarial meshes and textures. Post-processing and 386 mesh stabilization may require advanced techniques 387 like Marching Cubes [Tu et al., 2021] to ensure a wa-388 tertight, non-degenerate mesh that can be 3D printed. 389 390 Researchers have also developed flexible "universal" attacks that can be 3D printed once and deployed in 391 multiple scenarios without retraining Abdelfattah 392 et al., 2021]. When applying adversarial textures, 393 high-resolution printing or color constraints (Sec. 394 395 4.2) can enhance feasibility; however, covering large 396 surfaces is costly, prompting the use of localized sticker-mode" approaches [Li et al., 2024a] that only 397 modify a small area (e.g., a vehicle door). 398

### **399 6 Future Directions**

400 To further expose DNN vulnerabilities, we propose401 four directions for differentiable rendering research:

Target Diversity. Many differentiable rendering at-402 tacks focus on targeting cars used for autonomous 403 driving. Meanwhile, use of NeRF and 3D Gaus-404 sian Splatting has recently expanded into other real-405 world applications for robotics and unmanned aerial 406 systems (UAS) but remains largely unconsidered in 407 adversarial ML research. Exploring more diverse 408 targets in these applications would expand Task 4.1 409 and Task 4.2. 410

SOTA Models and Other Modalities. Existing dif-411 ferentiable rendering attacks mainly target image 412 classifiers and object detectors, with limited work on 413 optical flow, depth estimation, point cloud classifiers, 414 and multi-modal or multi-task fusion models [Ab-415 delfattah et al., 2021]. Attacks on 3D scene under-416 standing and advanced tasks like tracking or video 417 recognition remain underexplored, despite the grow-418 ing use of robust models in robotics and AR. Future 419 research could also include newer architectures such 420 as EfficientNet, ViT, and DeiT, which could exhibit 421 different vulnerabilities from older models. Exploit-422 423 ing these emerging vulnerabilities would advance attacker goals 2.1–2.5. 424

Attacks Considering Real-World Phenomena. 425 Current methods use only basic lighting and cam-426 era adjustments, such as varying lighting intensity 427 and position or camera resolution. This overlooks 428 complex environmental factors (e.g., variable light 429 shapes, shadows, color) and camera parameters (lens-430 warping, field of view, focus distance, and exposure) 431 that create new attack surfaces in drones and other 432 camera-equipped systems. Other physical attacks 433 involving placement of lens covers and rolling shut-434 ter exploitation [Sayles *et al.*, 2021] are also under-435 studied. Broadening research on such real-world 436 phenomena would strengthen Tasks 4.3, 4.4, and 437 5.2. 438

**Tools and Pipelines.** While simulators like CARLA 439 are widely used for attack research, differentiable 440 rendering libraries often require specialized knowl-441 edge and manual scene configuration. Existing GUIs, 442 such as Blender plugins for Mitsuba<sup>4</sup>, help export 443 scenes but still demand significant expertise in 3D 444 modeling. More user-friendly interfaces and inte-445 grated pipelines for differentiable renderers would 446 streamline digital attacks (Task 5.1), ultimately facil-447 itating transfer to physical scenarios (Task 5.2). 448

## 7 Conclusion

Understanding the evolving capabilities of differen-450 tiable rendering is essential for safeguarding deep 451 neural networks. This survey presents a task-guided 452 review of adversarial attacks using differentiable ren-453 dering, covering manipulations of 3D objects and 454 scenes that compromise applications like image clas-455 sification and object detection. By categorizing at-456 tacker tasks and linking them to goals, we highlight 457 research gaps such as attacks targeting scene pa-458 rameters (lighting, camera configurations) and the 459 need for user-friendly resources. Future work should 460 explore novel attack methods and practical physical 461 evaluations, facilitating more resilient DNN defenses 462 in this rapidly advancing area. 463

## References

### 464

449

[Abdelfattah *et al.*, 2021] M. Abdelfattah, K. Yuan,
 Z. Wang, and R. Ward. Towards Universal Physical Attacks On Cascaded Camera-Lidar 3d Object
 Detection Models. *ICIP*, 2021.

<sup>4</sup>https://github.com/mitsuba-renderer/ mitsuba-blender

- 469 [Alcorn et al., 2019] M. Alcorn, Q. Li, Z. Gong,
- 470 C. Wang, L. Mai, W. Ku, and A. Nguyen. Strike
- 471 (With) a Pose: Neural Networks Are Easily

Fooled by Strange Poses of Familiar Objects.*CVPR*, 2019.

- 474 [Bolya et al., 2020] D. Bolya, S. Foley, J. Hays, and
- 475 J. Hoffman. TIDE: A General Toolbox for Identi-
- 476 fying Object Detection Errors. *ECCV*, 2020.
- 477 [Byun et al., 2022] J. Byun, S. Cho, M. Kwon,
- H. Kim, and C. Kim. Improving the Transferability of Targeted Adversarial Examples through
  Object-Based Diverse Input. *arXiv*, 2022.
- 481 [Cao et al., 2019] Y. Cao, C. Xiao, D. Yang, J. Fang,
- R. Yang, M. Liu, and B. Li. Adversarial Objects
  Against LiDAR-Based Autonomous Driving Systems. *arXiv*, 2019.
- [Chen and Wang., 2024] G. Chen and W. Wang. A
   Survey on 3D Gaussian Splatting. *arXiv*, 2024.
- [Dong *et al.*, 2022] Y. Dong, S. Ruan, H. Su,
  C. Kang, X. Wei, and J. Zhu. ViewFool: evaluating the robustness of visual recognition to adversarial viewpoints. *NeurIPS*, 2022.
- 491 [Gao *et al.*, 2023] K. Gao, Y. Gao, H. He, D. Lu,
  492 L. Xu, and J. Li. NeRF: Neural Radiance Field
  493 in 3D Vision, A Comprehensive Review. *arXiv*,
  494 2023.
- <sup>495</sup> [Hu *et al.*, 2023] Z. Hu, W. Chu, X. Zhu, H. Zhang,
  <sup>496</sup> B. Zhang, and X. Hu. Physically Realizable
  <sup>497</sup> Natural-Looking Clothing Textures Evade Person
- <sup>498</sup> Detectors via 3D Modeling. *CVPR*, 2023.
- [Huang *et al.*, 2024] Y. Huang, Y. Dong, S. Ruan,
  X. Yang, H. Su, and X. Wei. Towards Transferable
  Targeted 3D Adversarial Attack in the Physical
  World. *CVPR*, 2024.
- [Jiang *et al.*, 2024] W. Jiang, H. Zhang, X. Wang,
  Z. Guo, and H. Wang. NeRFail: Neural Radiance Fields-Based Multiview Adversarial Attack. *AAAI*, 2024.
- [Kato *et al.*, 2020] H. Kato, D. Beker, M. Morariu,
  T. Ando, T. Matsuoka, W.Kehl, and A. Gaidon.
  Differentiable Rendering: A Survey. *arXiv*, 2020.
- 510 [Kerbl *et al.*, 2023] B. Kerbl, G. Kopanas,
  511 T. Leimkuehler, and G. Drettakis. 3D Gaussian
  512 Splatting for Real-Time Radiance Field Ren513 dering. ACM Transactions on Graphics, 42(4),
- 514 2023.

- [Leclerc *et al.*, 2022] G. Leclerc, H. Salman, 515
  A. Ilyas, S. Vemprala, L. Engstrom, V. Vineet, 516
  K. Xiao, P. Zhang, S. Santurkar, G. Yang, 517
  A. Kapoor, and A. Madry. 3DB: A Framework for Debugging Computer Vision Models. 519 *NeurIPS*, 2022. 520
- [Leheng *et al.*, 2023] L. Leheng, Q. Lian, and 521 Y. Chen. Adv3D: Generating 3D Adversarial Examples in Driving Scenarios with NeRF. *arXiv*, 523 2023. 524
- [Li et al., 2024a] Y. Li, W. Tan, C. Zhao, S. Zhou,
   X. Liang, and Q. Pan. Flexible Physical Camouflage Generation Based on a Differential Approach. arXiv, 2024.
- [Li et al., 2024b] Y. Li, B. Xie, S. Guo, Y. Yang, 529 and B. Xiao. A Survey of Robustness and Safety 530 of 2D and 3D Deep Learning Models against Adversarial Attacks. ACM Computing Surveys, 532 56(6), 2024. 533
- [Liu et al., 2019] H. Liu, M. Tao, C. Li, 534
   D. Nowrouzezahrai, and A. Jacobson. Beyond Pixel Norm-Balls: Parametric Adversaries 536
   using an Analytically Differentiable Renderer. 537
   *ICLR*, 2019. 538
- [Machado et al., 2023] G. Machado, E. Silva, and539R. Goldschmidt. Adversarial Machine Learning540in Image Classification: A Survey Toward the541Defender's Perspective. ACM Computing Surveys,54255(1), 2023.543
- [Maesumi et al., 2021] A. Maesumi, M. Zhu, 544
   Y. Wang, T. Chen, Z. Wang, and C. Bajaj. 545
   Learning Transferable 3D Adversarial Cloaks for 546
   Deep Trained Detectors. arXiv, 2021. 547
- [Mahendran and Vedaldi, 2015] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. *CVPR*, 2015. 550
- [Meloni *et al.*, 2021] E. Meloni, M. Tiezzi, 551 L. Pasqualini, M. Gori, and S. Melacci. Messing Up 3D Virtual Environments: Transferable Adversarial 3D Objects. *ICMLA*, 2021. 554
- [Mildenhall *et al.*, 2020] B. Mildenhall, P. Srinivasan, M. Tancik, J. Barron, N. Ramamoorthi, and R. Ng. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. *ECCV*, 2020.
- [Miller *et al.*, 2020] D. Miller, Z. Xiang, and G. Kesidis. Adversarial Learning Targeting Deep Neu-561

- ral Network Classification: A Comprehensive Re-
- view of Defenses Against Attacks. *Proceedings* of the IEEE, 108(3), 2020.
- [Mittal, 2024] A. Mittal. Neural Radiance Fields:
  Past, Present, and Future. *arXiv*, 2024.
- 567 [Papernot et al., 2016a] N. Papernot, P. McDaniel,
- 568 S. Jha, M. Fredrikson, Z. Celik, and A. Swami.
   569 The Limitations of Deep Learning in Adversarial
- The Limitations of Deep Learning in Adv Settings. *EuroS&P*, 2016.
- 571 [Papernot et al., 2016b] N. Papernot, P. McDaniel,
- A. Sinha, and M. Wellman. Towards the Science
   of Security and Privacy in Machine Learning. *EuroS&P*, 2016.
- [Sayles *et al.*, 2021] A. Sayles, A. Hooda,
  M. Gupta, R. Chatterjee, and E. Fernandes.
  Invisible Perturbations: Physical Adversarial
  Examples Exploiting the Rolling Shutter Effect. *CVPR*, 2021.
- [Schmalfuss *et al.*, 2023] J. Schmalfuss, L. Mehl,
  and A. Bruhn. Distracting Downpour: Adversarial Weather Attacks for Motion Estimation. *ICCV*,
  2023.
- [Shahreza and Marcel, 2023] H. Shahreza and
  S. Marcel. Comprehensive Vulnerability Evaluation of Face Recognition Systems to Template
  Inversion Attacks via 3D Face Reconstruction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(12), 2023.
- [Sharif *et al.*, 2016] M. Sharif, S. Bhagavatula,
  L. Bauer, and M. K. Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. *CCS*, 2016.
- [Suryanto *et al.*, 2022] N. Suryanto, Y. Kim,
  H. Kang, H. Larasati, Y. Yun, T. Le, H. Yang,
  S. Oh, and H. Kim. DTA: Physical Camouflage
  Attacks using Differentiable Transformation
  Network. *CVPR*, 2022.
- [Suryanto *et al.*, 2023] N. Suryanto, Y. Kim,
  H. Larasati, H. Kang, T. Le, Y. Hong, H. Yang,
  S. Oh, and H. Kim. ACTIVE: Towards Highly
  Transferable 3D Physical Camouflage for
  Universal and Robust Vehicle Evasion. *ICCV*,
  2023.
- [Tewari *et al.*, 2022] A. Tewari, J. Thies, andB. Mildenhall. Advances in Neural Rendering.
- 607 *Computer Graphics Forum*, 41(22), 2022.

- [Tosi *et al.*, 2024] F. Tosi, Y. Zhang, Z. Gong, 608
  E. Sandström, S. Mattoccia, M.R. Oswald, and 609
  M. Poggi. How NeRFs and 3D Gaussian Splatting are Reshaping SLAM: a Survey. *arXiv*, 2024. 611
- [Tu *et al.*, 2021] J. Tu, H. Li, X. Yan, M. Ren, Y. Chen, M. Liang, E. Bitar, E. Yumer, and R. Urtasun. Exploring Adversarial Robustness of Multisensor Perception Systems in Self Driving. *CoRL*, 2021. 616
- [Wang et al., 2022] D. Wang, T. Jiang, J. Sun, 617
  W. Zhou, X. Zhang, Z. Gong, W. Yao, and 618
  X. Chen. FCA: Learning a 3D Full-coverage 619
  Vehicle Camouflage for Multi-view Physical Adversarial Attack. AAAI, 2022. 621
- [Wiyatno *et al.*, 2019] R. Wiyatno, A. Xu, O. Dia, and A. de Berker. Adversarial Examples in Modern Machine Learning: A Review. *arXiv*, 2019. 624
- [Xiao et al., 2019] C. Xiao, D. Yang, B. Li, J. Deng, and M. Liu. MeshAdv: Adversarial Meshes for Visual Recognition. CVPR, 2019.
- [Xie *et al.*, 2022] Y. Xie, T. Takikawa, S. Saito, 628
  O. Litany, S. Yan, N. Khan, F. Tombari, J. Tompkin, V. Sitzmann, and S. Sridhar. Neural Fields in Visual Computing and Beyond. *arXiv*, 2022. 631
- [Yuan et al., 2019] X. Yuan, P. He, Q. Q. Zhu, and
   X. Li. Adversarial Examples: Attacks and De fenses for Deep Learning. *IEEE Transactions on Neural Networks and Learning Systems*, 30(9),
   2019.
- [Zeng et al., 2019] X. Zeng, C. Liu, Y. Wang, 637
   W. Qiu, L. Xie, Y. Tai, C. Tang, and A. Yuille. Adversarial Attacks Beyond the Image Space. CVPR, 639
   2019. 640
- [Zheng et al., 2024] J. Zheng, C. Lin, J. Sun, 641
   Z. Zhao, Q. Li, and C. Shen. Physical 3D Adver-642
   sarial Attacks against Monocular Depth Estimation in Autonomous Driving. *CVPR*, 2024. 644
- [Zhou *et al.*, 2024] J. Zhou, L. Lyu, D. He, and
   Y. Li. RAUCA: A Novel Physical Adversarial
   Attack on Vehicle Detectors via Robust and Accurate Camouflage Generation. *ICML*, 2024.