

Video2Game: Real-time, Interactive, Realistic and Browser-Compatible Environment from a Single Video

Anonymous CVPR submission

Paper ID

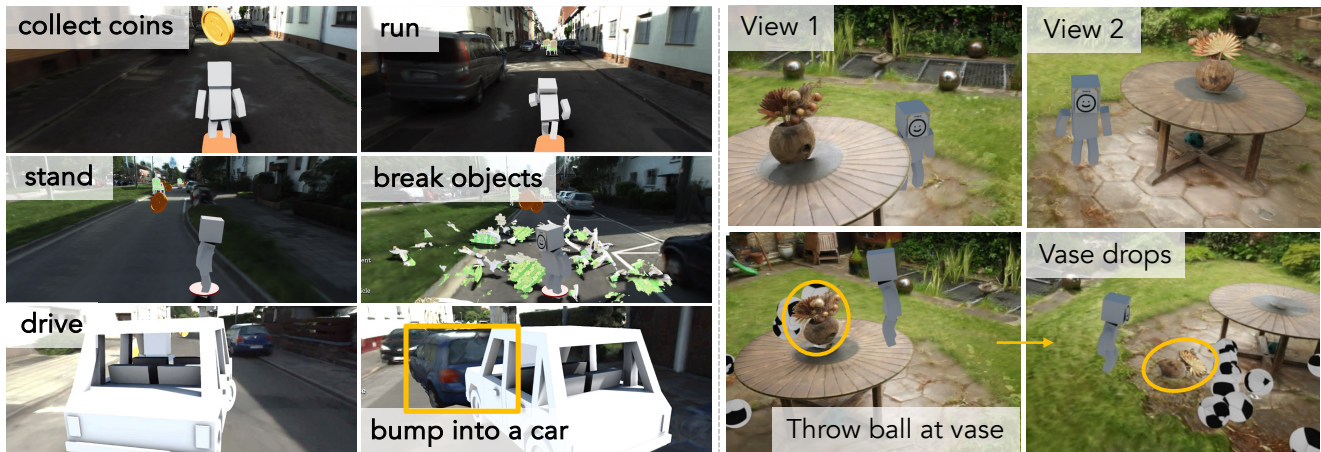


Figure 1. **Video2Game** takes an input video of an arbitrary scene and automatically transforms it into a real-time, interactive, realistic and browser-compatible environment. The users can freely explore the environment and interact with the objects in the scene.

Abstract

001 *Creating high-quality and interactive virtual environ-*
 002 *ments, such as games and simulators, often involves com-*
 003 *plex and costly manual modeling processes. In this paper,*
 004 *we present Video2Game, a novel approach that automati-*
 005 *cally converts videos of real-world scenes into realistic and*
 006 *interactive game environments. At the heart of our sys-*
 007 *tem are three core components: (i) a neural radiance fields*
 008 *(NeRF) module that effectively captures the geometry and*
 009 *visual appearance of the scene; (ii) a mesh module that dis-*
 010 *tills the knowledge from NeRF for faster rendering; and (iii)*
 011 *a physics module that models the interactions and physical*
 012 *dynamics among the objects. By following the carefully de-*
 013 *signed pipeline, one can construct an interactable and ac-*
 014 *tionable digital replica of the real world. We benchmark our*
 015 *system on both indoor and large-scale outdoor scenes. We*
 016 *show that we can not only produce highly-realistic render-*
 017 *ings in real-time, but also build interactive games on top.*

019 1. Introduction

020 Crafting a visually compelling and interactive environment
 021 is crucial for immersive experiences in various domains,

such as video games, virtual reality applications, and self-
 driving simulators. This process, however, is complex and
 expensive. It demands the skills of experts in the field and
 the use of professional software development tools [21, 24].
 For instance, Grand Theft Auto V [23], known for its in-
 tricately detailed environment, was one of the most expen-
 sive video games ever developed, with a budget over \$265
 million primarily for asset creation. Similarly, the develop-
 ment of the CARLA autonomous driving simulator [19] in-
 volves a multidisciplinary team of 3D artists, programmers,
 and engineers to meticulously craft and texture the virtual
 cityscapes, creating its lifelike environments.

An appealing alternative to extensive manual modelling
 is creating environments directly from the real world. For
 instance, photogrammetry, a technique for constructing dig-
 ital replicas of objects or scenes from overlapping real-
 world photographs, has been utilized for environment crea-
 tion [52, 53]. Success stories also span various games
 and simulators. However, most use cases are limited to
 creating object assets and necessitate significant post-
 processing, such as material creation, texturing, and geom-
 etry fixes [66]. People thus turns to neural radiance fields
 (NeRFs) [46], as it offers a more promising approach to

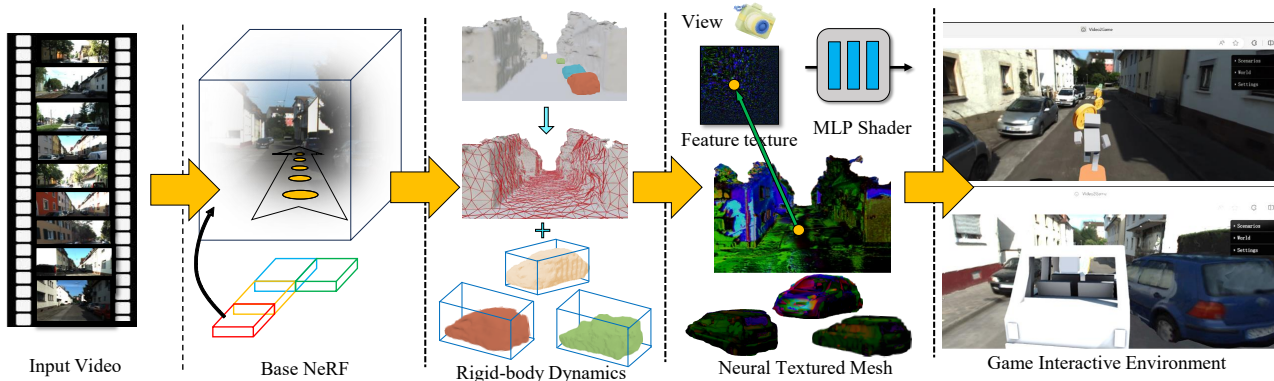


Figure 2. **Overview of Video2Game:** Given multiple posed images from a single video as input, we first construct a large-scale NeRF model that is realistic and possesses high-quality surface geometry. We then transform this NeRF model into a mesh representation with corresponding rigid-body dynamics to enable interactions. We utilize UV-mapped neural texture, which is both expressive and compatible with game engines. Finally, we obtain an interactive virtual environment that virtual actors can interact with, can respond to user control, and deliver high-resolution rendering from novel camera perspectives – all in real-time.

045 modeling large scenes. With careful design [13, 22, 39, 48,
046 61], NeRF is able to render free-viewpoint, photo-realistic
047 images efficiently. However, crafting an interactive environ-
048 ment entails more than just creating a visually high-fidelity
049 digital twin; it also involves building a **physically plau-**
050 **sible, immersive, real-time** and importantly, **interactive**
051 world tailored to user experiences. Furthermore, we expect
052 such a virtual world to be **compatible with real-time inter-**
053 **action interfaces** such as common game engines. Despite
054 its promise, the use of NeRF to create interactive environ-
055 ments from real-world videos remains largely unexplored.

056 In this paper, we introduce Video2Game, a novel ap-
057 proach to automatically converting a video of a scene into a
058 realistic and interactive virtual environment. Given a video
059 as input, we first construct a NeRF that can effectively cap-
060 ture the geometric and visual information of a large-scale,
061 (unbounded) scene. Then we distill the NeRF into a game
062 engine-compatible, neural textured mesh. This significantly
063 improves the rendering efficiency while maintains the over-
064 all quality. To model the interactions among the objects,
065 we further decompose the scene into individual actionable
066 entities and equip them with respective physics model. Fi-
067 nally, we import our automatically generated assets into a
068 WebGL-based game engine and create a playable game.
069 The resulting virtual environment is photo-realistic, inter-
070 active, and runs in real-time. See Fig. 1 for demonstration.
071 In summary, our key contributions are:

- 072 • A novel 3D modeling algorithm for real-time, free-
073 viewpoint rendering and physical simulation, surpassing
074 state-of-the-art NeRF baking methods with added rigid-
075 body physics for enhanced simulation.
- 076 • An automated game-creation framework to transform a
077 scene video into an interactive, realistic environment,
078 compatible with current game engines.

2. Related Works 079

080 Given a single video, we aim to create a real-time, interac-
081 tive game where the agents (*e.g.*, the character, the car) can
082 navigate and explore the reconstructed digital world, inter-
083 act with objects in the scene (*e.g.*, collision and manipulate
084 objects), and achieve their respective tasks (*e.g.*, collecting
085 coins, shooting targets). We draw inspirations from multi-
086 ple areas and combine the best of all. In this section, we
087 will briefly review those closely related areas which forms
088 the foundation of our work.

089 **Novel view synthesis (NVS):** Our work builds upon the
090 success of novel view synthesis [14, 25, 35, 62], which
091 is crucial for our game since it enables the agents to
092 move freely and view the reconstructed world seamlessly
093 from various perspectives. Among all these approaches
094 [26, 60, 68, 85, 86], we exploit neural radiance field (NeRF)
095 [46] as our underlying representation. NeRF has emerged as
096 one of the most promising tools in NVS since its introduc-
097 tion [49–51], and has great performance across a wide range
098 of scenarios [36, 56, 75, 81]. For instance, it can be easily
099 extended to handle various challenging real-world scenar-
100 ios such as learning from noisy camera poses [38, 70],
101 reflectance modeling for photo-realistic relighting [69, 83],
102 and real-time rendering [16, 39, 55, 65, 76]. In this work,
103 we combine recent advances in NeRF with *physics mod-*
104 *eling* to build an immersive digital world from one single
105 video, moving from *passive* NVS to our complete solution
106 for *embodied* world modeling where agents can *actively*
107 explore and interact with the scene.

108 **Controllable video generation:** Using different control
109 signals to manipulate the output of a visual model has gar-
110 nered great interest in the community. This has had a pro-
111 found impact on content creation [57, 58], digital editing
112 [11, 34], and simulation [30, 31, 40]. One could also lever-

113 age large foundation models to control video content using
 114 text [57, 58]. However, they lack fine-grained and real-time
 115 control over the generated content. Alternatively, training
 116 (conditional) generative models for each scene enables bet-
 117 ter disentanglement of dynamics (*e.g.*, foreground *vs.* back-
 118 ground) and supports better control. For instance, one can
 119 represent a self-driving scene [31] or a Pacman game [30] as
 120 latent codes and generate video frames based on control in-
 121 puts with a neural network. One can also learn to control the
 122 players within tennis games [44, 45, 79, 80]. Our work falls
 123 under the second line of research, where the model takes
 124 user control signals (*e.g.*, a keystroke from the keyboard)
 125 as input and responds by rendering a new scene. However,
 126 instead of focusing on a specific scene (*e.g.*, tennis games),
 127 we have developed a pipeline that allows the creation of
 128 a playable environment from a single video of a generic
 129 scene. Additionally, we model everything in 3D, which
 130 enables us to effectively capture not only view-dependent
 131 appearance but also physical interactions among rigid-body
 132 equipped objects. Importantly, we adopt a neural represen-
 133 tation compatible with graphics engines, enabling users to
 134 play the entire game in their browser at an interactive rate.

135 **Data-driven simulation:** Building a realistic simulation
 136 environment has been a longstanding challenge. [19, 29,
 137 67, 71]. While it’s promising, we come close to mirror the
 138 real world only in recent years [10, 15, 42, 43, 59, 74, 75].
 139 The key insight of these work is to build models by lever-
 140 aging real-world data. Our work closely aligns with this
 141 line of research on building high-fidelity simulators from
 142 real-world data, with a few key differences. First, exist-
 143 ing works mainly focus on *offline* training and evaluation
 144 [10, 43, 74, 75], whereas our system runs at an interactive
 145 rate and allows for *online, real-time* control. Second, some
 146 existing works[41, 43, 72, 87] need additional data modality
 147 like LiDAR point clouds for geometry reconstruction, but
 148 RGB video is all we need. Third, most photo-realistic sim-
 149 ulators don’t model physical interactions. However, we sup-
 150 ports various physics modeling and allows agents to interact
 151 with the environment. Last, existing simulators are typi-
 152 cally resource-intensive , while our system is lightweight
 153 and can be easily accessible in common engines.

154 3. Video2Game

155 Given a sequence of images or a video of a scene, our goal
 156 is to construct an *interactable* and *actionable* digital twin,
 157 upon which we can build real-time, interactive games or re-
 158 alistic (sensor) simulators. Based on the observations that
 159 prevalent approaches to constructing digital replica mainly
 160 focus on visual appearance and ignore the underlying phys-
 161 ical interactions, we carefully design our system such that
 162 it can not only produce high-quality rendering across view-
 163 points, but also support the modeling of physical actions
 164 (*e.g.*, navigation, collision, manipulation, etc). At the heart

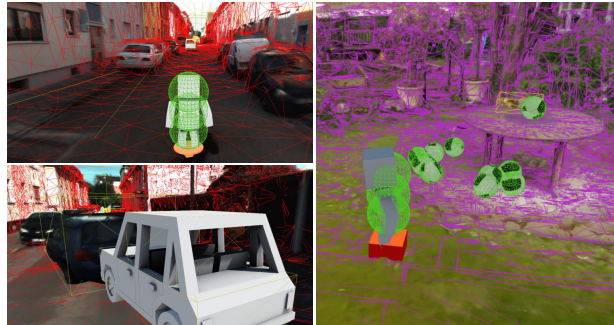


Figure 3. **Visualization of automatically computed collision geometry:** Sphere collider (green), box collider (yellow), convex polygon collider (purple) and trimesh collider (red).

of our systems is a compositional implicit-explicit 3D rep- 165
 resentation that is effective and efficient for both sensor and 166
 physics simulation. By decomposing the world into indi- 167
 vidual entities, we can better model and manipulate their 168
 physical properties (*e.g.*, specularity, mass, friction), and 169
 simulate the outcomes of interactions more effectively. 170

We start by introducing a NeRF model that can effective- 171
 ly capture the geometric and visual information of a 172
 large-scale, unbounded scene (Sec. 3.1). Next, we present 173
 an approach to convert the NeRF into a game-engine com- 174
 patible mesh with neural texture maps, significantly im- 175
 proving the rendering efficiency while maintaining the qual- 176
 ity (Sec. 3.2). To enable physical interactions, we further 177
 decompose the scene into individual actionable entities and 178
 equip them with respective physics models (Sec. 3.3). Fi- 179
 nally, we describe how we integrate our interactive environ- 180
 ment into a WebGL-based game engine, allowing users to 181
 play and interact with the virtual world in real time on their 182
 personal browser. Fig. 2 provides an overview of our pro- 183
 posed framework. 184

3.1. Large-scale NeRF 185

Preliminaries: Instant-NGP [48] is a notable variant of 186
 NeRF, which represents the radiance field with a combi- 187
 nation of spatial hash-based voxels and neural networks: 188
 $\mathbf{c}, \sigma = F_{\theta}(\mathbf{x}, \mathbf{d}; \Phi) = \text{MLP}_{\theta}(\text{It}(\mathbf{x}, \Phi), \mathbf{d})$. Given a 3D 189
 point $\mathbf{x} \in \mathbb{R}^3$ and a camera direction $\mathbf{d} \in \mathbb{R}^2$ as input, 190
 Instant-NGP first interpolate the point feature $\text{It}(\mathbf{x}, \Phi)$ 191
 from the adjacent voxel features Φ . Then the point feature 192
 and the camera direction are fed into a light-weight multi- 193
 layer perception (MLP) to predict the color $\mathbf{c} \in \mathbb{R}^3$ and 194
 density $\sigma \in \mathbb{R}^+$. To render the scene appearance, we first cast a 195
 ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ from the camera center \mathbf{o} through the pixel 196
 center in direction \mathbf{d} , and sample a set of 3D points $\{\mathbf{x}_i\}$ 197
 along the ray. We then query their respective color $\{\mathbf{c}_i\}$ 198
 and density $\{\sigma_i\}$ and obtain the color of the pixel through 199
 alpha-composition: $\mathbf{C}_{\text{NeRF}}(\mathbf{r}) = \sum_i w_i \mathbf{c}_i$. Similarly, the 200

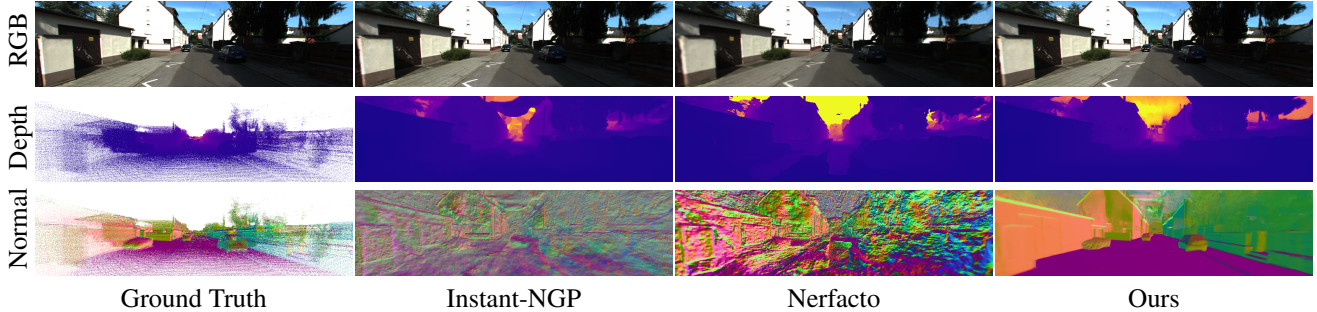


Figure 4. **Qualitative comparisons among NeRF models.** The rendering quality of our base NeRF is superior to baselines, and with leveraging monocular cues, we substantially improve rendered geometry compared to other baselines. This significantly facilitates NeRF baking in subsequent stages. Here we consider depths measured by LiDAR point cloud in KITTI-360 and compute normals based on it.

201 expected depth can be computed by: $\mathbf{D}_{\text{NeRF}}(\mathbf{r}) = \sum_i w_i t_i$.
 202 Here, w_i indicates the blending weight that is derived from
 203 the densities $\{\sigma_i\}$. We refer the readers to [46] for more de-
 204 tails. To learn the voxel features Φ and the MLP weights θ ,
 205 we compute the difference between the ground truth color
 206 and the rendered color: $\mathcal{L}_{\text{rgb}} = \sum_{\mathbf{r}} \|\mathbf{C}_{\text{GT}}(\mathbf{r}) - \mathbf{C}_{\text{NeRF}}(\mathbf{r})\|_2^2$.
 207 **Large-scale NeRF:** While Instant-NGP [48] has shown
 208 promising results on densely observed and bounded scenes,
 209 its performance starts to degrade when extending to
 210 sparsely-captured, large-scale, unbounded environments.
 211 To mitigate these issues, we propose several enhancements:
 212
 213
$$\mathbf{c}, \sigma, s, \mathbf{n} = F_{\theta}(\mathbf{x}, \mathbf{d}; \Phi) = \text{MLP}_{\theta}(\text{It}(\text{Ct}(\mathbf{x}), \Phi), \mathbf{d}). \quad (1)$$

214 First of all, we exploit the contraction function $\text{Ct}(\mathbf{x})$ [12]
 215 to map the unbounded coordinates into a bounded region.
 216 In addition to radiance and density, we predict the seman-
 217 tics s and the surface normal \mathbf{n} of the 3D points, guided with
 218 2D priors to better regularize the scene geometry. Further-
 219 more, we divide large-scale scenes into several blocks [63]
 220 to capture the fine-grained details. We now describe these
 221 enhancements in more details.

222 **Depth:** High-quality geometry is critical for model-
 223 ing physical interactions. Inspired by MonoSDF [78], we
 224 leverage off-the-shelf monocular depth estimators [20, 27]
 225 to guide and improve the underlying geometry. We first
 226 predict the depth of the scene from rendered RGB im-
 227 ages. Then we minimize the discrepancy between the
 228 rendered depth and the predicted depth via $\mathcal{L}_{\text{depth}} =$
 229 $\sum_{\mathbf{r}} \|(a\mathbf{D}_{\text{NeRF}}(\mathbf{r}) + b) - \mathbf{D}_{\text{mono}}(\mathbf{r})\|_2^2$, where a and b are the
 230 scale and shift that aligns the two distribution [54].

231 **Surface normals:** Similar to depth, we encourage the
 232 normal estimated from NeRF to be consistent with the nor-
 233 mal predicted by the off-the-shelf estimator [20, 27]. The
 234 normal of a 3D point \mathbf{x}_i can be either analytically derived
 235 from the estimated density $\mathbf{n}_i = (1 - \frac{\nabla_{\mathbf{x}} \sigma_i}{\|\nabla_{\mathbf{x}} \sigma_i\|})$, or pre-
 236 dicted by the MLP header as in Eq. 1. We could aggre-
 237 gate them via volume render: $\mathbf{N}(\mathbf{r}) = \sum_i \mathbf{w}_i \mathbf{n}_i$. Em-

238 pirically we find that adopting both normals and promot-
 239 ing their mutual consistency works the best, since the MLP
 240 header offers more flexibility. We thus employ $\mathcal{L}_{\text{normal}} =$
 241 $\|\mathbf{N}_{\text{mlp}}(\mathbf{r}) - \mathbf{N}_{\text{mono}}(\mathbf{r})\|_2^2 + \|\mathbf{N}_{\text{mlp}}(\mathbf{r}) - \mathbf{N}_{\text{density}}(\mathbf{r})\|_2^2$.

242 **Semantics:** We also predict semantic logits for each sam-
 243 pled 3D points with our MLP. This helps us capture the cor-
 244 relation between semantics and geometry [36, 84]. We ren-
 245 der the semantic map with volume rendering $\mathbf{S}_{\text{NeRF}}(\mathbf{r}) =$
 246 $\sum_i \mathbf{w}_i s_i$ and compute the cross-entropy with that of a 2D
 247 segmentation model $\mathcal{L}_{\text{semantics}} = \text{CE}(\mathbf{S}_{\text{mono}}, \mathbf{S}_{\text{NeRF}})$.

248 **Regularization:** We additionally adopt two regulariza-
 249 tion terms. To reduce floaters in the scene, for each ran-
 250 domly sampled 3D point \mathbf{x} , we penalize its density by
 251 $\mathcal{L}_{\text{sp}} = \sum 1 - \exp(-\alpha\sigma(\mathbf{x}))$, where $\alpha > 0$ [77]. For each
 252 sky pixel (which we derived from the semantic MLP), we
 253 encourage its depth $\mathbf{D}_{\text{NeRF}}(\mathbf{r}^{\text{sky}})$ to be as far as possible.
 254 The loss is defined as: $\mathcal{L}_{\text{sky}} = \sum_{\mathbf{r}^{\text{sky}}} \exp(-\mathbf{D}_{\text{NeRF}}(\mathbf{r}^{\text{sky}}))$.

255 **Blocking:** Capitalizing on a single Instant-NGP to cover
 256 an extraordinarily large scene such as KITTI-360 [37]
 257 would often lead to inferior results. We thus adopt a strat-
 258 egy akin to BlockNeRF [63] where we divided the whole
 259 scene into numerous blocks and model each region with a
 260 separate Instant-NGP. Adjacent regions maintain substan-
 261 tial overlaps to ensure smooth transition.

262 **Learning:** We jointly optimize the voxel feature Φ and
 263 the MLP weights θ by minimizing the following loss:

$$\mathcal{L}_{\text{total}}^{\text{NeRF}} = \mathcal{L}_{\text{rgb}} + \mathcal{L}_{\text{normal}} + \mathcal{L}_{\text{semantics}} + \mathcal{L}_{\text{depth}} + \mathcal{L}_{\text{sky}} + \mathcal{L}_{\text{sp}} \quad (2) \quad 264$$

3.2. NeRF Baking 265

266 We aim to create a digital replica that users (or agents)
 267 can freely explore and act upon in real time. Although
 268 our large-scale NeRF effectively renders high-quality im-
 269 ages and geometry, its efficiency is limited by the compu-
 270 tational costs associated with sampling 3D points. The un-
 271 derlying volume density representation further complicates

Method	Representation	KITTI-360			Gardenvase			Interactive Compatibility		
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Real time	Rigid-body physics	Scene decomposition
InstantNGP [48]	Volume	27.46	0.853	0.165	25.90	0.757	0.191	\times	\times	\times
Nerfacto [64]		23.20	0.763	0.238	22.16	0.517	0.283	\times	\times	\times
Video2Game		27.62	0.871	0.131	26.57	0.815	0.143	\times	\times	\times
Gauss. Spl. [28]	Points	17.85	0.615	0.428	27.50	0.858	0.099	\checkmark	\times	\times
MobileNeRF [16]	Mesh	19.67	0.627	0.452	22.80	0.505	0.365	\checkmark	\times	\times
BakedSDF* [76]		22.37	0.757	0.302	22.68	0.514	0.369	\checkmark	\checkmark	\times
Video2Game		23.35	0.765	0.246	22.81	0.508	0.363	\checkmark	\checkmark	\checkmark

Table 1. **Quantitative results on novel view synthesis and interactive compatibility analysis.** Video2Game produces better or comparable results across scenes, suggesting the effectiveness of our NeRF and mesh model. The performance improves the most when tackling unbounded, large-scale scenes in KITTI-360. We note that existing NeRFs cannot reach the interactive rate required for real-time games. While point-based rendering significantly improves the speed, it does not support rigid body physics. BakedSDF [76] represents the whole scene with one single mesh, thus does not support object-level interactions.

the problem. For instance, it’s unclear how to define physical interaction with such a representation (e.g., defining collision). The representation is also not compatible with common graphics engines. While recent software, such as the NeRFStudio Blender plugin and LumaAI Unreal add-on, has made some strides, their interaction capabilities and scene geometry quality are still not optimal for real-time user engagement, especially when the scene is large and the observations are relatively sparse. To overcome these challenges, we draw inspiration from recent NeRF meshing advancements and present a novel NeRF baking framework that efficiently transforms our NeRF representation into a game-engine compatible mesh. As we will show in Sec. 4, this conversion greatly enhances rendering efficiency while preserving quality and facilitates physical interactions.

Mesh representation: Our mesh $\mathcal{M} = (\mathbf{V}, \mathbf{F}, \mathbf{T})$ comprises vertices $\mathbf{V} \in \mathbb{R}^{|\mathbf{V}| \times 3}$, faces $\mathbf{F} \in \mathbb{N}^{|\mathbf{F}| \times 3}$ and a UV neural texture map $\mathbf{T} \in \mathbb{R}^{H \times W \times 6}$. Following [65], we store the base color in the first three dimension of \mathbf{T} , and encode the specular feature in the rest. The initial mesh topology are obtained by marching cubes in the NeRF density field. We further prune the invisible faces. conduct mesh decimation and edge length regularization. The UV coordinate of each vertex is calculated via xatlas [7].

Rendering: We leverage differentiable renderers [33] to render our mesh into RGB images \mathbf{C}_R and depth maps \mathbf{D}_R . Specifically, we first rasterize the mesh into screen space and obtain the UV coordinate for each pixel i . Then we sample the corresponding texture feature $\mathbf{T}_i = [\mathbf{B}_i; \mathbf{S}_i]$ and feed it into our customized shader. Finally, the shader computes the sum of the view-independent base color \mathbf{B}_i and the view-dependent MLP $\text{MLP}_\theta^{\text{shader}}(\mathbf{S}_i, \mathbf{d}_i)$:

$$\mathbf{C}_R = \mathbf{B}_i + \text{MLP}_\theta^{\text{shader}}(\mathbf{S}_i, \mathbf{d}_i). \quad (3)$$

The MLP is lightweight and can be baked in GLSL.

Learning: We train the shader MLP $\text{MLP}_\theta^{\text{shader}}$ and the neural texture map \mathbf{T} by minimizing the color difference

between the mesh and the ground truth, and the geometry difference between the mesh and the NeRF model:

$$\mathcal{L}_{\mathbf{T}, \theta}^{\text{mesh}} = \sum_{\mathbf{r}} \|\mathbf{C}_R(\mathbf{r}) - \mathbf{C}_{\text{GT}}(\mathbf{r})\| + \|\mathbf{D}_R(\mathbf{r}) - \mathbf{D}_{\text{NeRF}}(\mathbf{r})\|. \quad (4)$$

Anti-aliasing: Common differentiable rasterizers only take the center of each pixel into account. This may lead to aliasing in the learned texture map. To resolve this issue, we randomly perturb the optical center of the camera by 0.5 pixels along each axis at every training step. This ensure all the regions within a pixel get rasterized.

Relationship to existing work: Our approach is closely related to recent work on NeRF meshing [16, 55, 65, 76], but there exist key differences. While MobileNeRF [16] also adopts an explicit mesh with neural textures, they mainly capitalize on planar primitives. The quality of the reconstructed mesh is thus inferior. BakedSDF [76] offers excellent runtime and rendering quality, but their vertex coloring approach has limited resolution for large scenes. NeRF2Mesh [65] lacks depth distillation and doesn’t adopt contraction space for unbounded scenes. They also have a sophisticated multi-stage training and multi-resolution refinement process. Finally, MeRF [55], though efficient, still relies on volume-rendering.

3.3. Representation for Physical Interaction

Our mesh model facilitates efficient novel-view rendering in real time and allows for basic *rigid-body* physical interactions. For example, the explicit mesh structure permits an agent to “stand” on the ground. Nevertheless, beyond navigation, an agent should be capable of performing various actions including collision and manipulation. Furthermore, a scene comprises not only the background but also interactable foreground objects, each possessing unique physical properties. For instance, a street-bound car is much heavier than a flower vase. When struck by another object, a car may barely move but the vase may fall and shatter.

342 To enhance physical interaction realism, we decompose the
343 scene into discrete, actionable entities, each endowed with
344 specific physical characteristics (*e.g.*, mass, friction). This
345 approach, in conjunction with *rigid-body* physics, allows
346 for the effective simulation that adheres to physical laws.

347 **Scene decomposition:** Directly editing and decompos-
348 ing a mesh is extremely difficult due to topology change.
349 Fortunately, neural fields are inherently compositional in
350 3D. By identifying the objects each spatial region belongs
351 to, we can use neural fields to guide the decomposition of
352 the mesh. Specifically, we sample a 3D point \mathbf{x}_i within each
353 voxel i and determine its semantic category either through
354 the predicted semantic logits s_i or by verifying whether the
355 point is within a specified bounding box. The process is re-
356 peated for all voxels to segment the entire scene. Then, for
357 each object, we perform NeRF meshing individually, setting
358 the density of the remaining areas to zero. The intersec-
359 tions between objects are automatically resolved by march-
360 ing cube. Finally, we initialize the neural texture of these
361 new, individual meshes from the original mesh model. For
362 newly created faces, we employ nearest neighbor inpainting
363 on the neural texture map, which empirically yields satis-
364 factory results. Fig. 1 shows an example where a vase is
365 separated from a table. The middle of the table is original
366 occluded yet we are able to maintain high-quality rendering.

367 **Physical parameters reasoning:** The next step is to
368 equip decomposed individual meshes with various physics-
369 related attributes so that we can effectively model and sim-
370 ulate their physical dynamics. In this work, we focus on
371 rigid body physics, where each entity i is represented by
372 a collision geometry col_i , mass m_i , and friction parameters
373 f_i . We support four types of collision geometry with differ-
374 ent levels of complexity and efficiency: box, sphere, convex
375 polygon, and triangle mesh. Depending on the object and
376 the task of interest, one can select the most suitable collision
377 check for them. For other physical parameters (*e.g.* mass,
378 friction), one can either set them manually or query large
379 language models (LLMs) for an estimation.

380 **Physical interactions:** Rigid body dynamics, while sim-
381 ple, can support a variety of interactions. With the collision
382 check, an user/agent can easily navigate through the envi-
383 ronment while respecting the geometry of the scene. The
384 agents will no longer be stuck in a road or cut through a
385 wall. It also allows the agent to interact with the objects in
386 the scene. For instance, one can push the objects towards
387 the location of interest. The object movement will be deter-
388 mined by its mass and other physical properties such as the
389 friction. We can also manipulate the objects by adopting a
390 magnet grasper, following AI2-Thor [32]. This opens the
391 avenue towards automatic creation of realistic, interactive
392 virtual environment for robot learning.

3.4. Interactive Environment

393 We deploy our interactive environment within a real-time,
394 browser-based game engine. We manage the underlying
395 logic and assets using Sketchbook [3], a Game Engine
396 based on Three.js that leverages WebGL [4] for rendering.
397 This combination ensures high efficiency while offering the
398 flexibility and sophistication required for intricate render-
399 ing tasks. It also allows us to easily integrate content from
400 different scenes together. We have further extended Sketch-
401 book’s capabilities by implementing a GLSL-based shader
402 [2]. This enables real-time computation of our MLP-based
403 specular shader during deployment. For physics simula-
404 tion, we use Cannon.js [1], which assures realism and ef-
405 ficiency in the motion within our interactive environment.
406 It supports not only rigid body dynamics but also more so-
407 phisticated modeling techniques. For example, we can pre-
408 compute the fracturing effect for dynamic objects. Upon
409 experiencing a significant force, these objects are realisti-
410 cally simulated by the real-time physics engine, which han-
411 dles the interactions between the fractured pieces and the
412 rest of the scene, such as their falling and settling on the
413 ground. Besides browser-based engine, the virtual environ-
414 ments from Video2Game pipeline could be also integrated
415 into both **Blender** [17] and **Unreal** engines [21] (see Fig. 6).
416

4. Experiments

417 We begin by presenting our experimental setup, followed by
418 a comparison of our model with state-of-the-art approaches.
419 Next, we conduct an extensive analysis of our model’s dis-
420 tinctive features and design choices. Then we demonstrate
421 how we constructed a web browser-compatible game capa-
422 ble of delivering a smooth interactive experience exceed-
423 ing 100 frames per second (FPS), all derived from a single
424 video source. Finally, we showcase the capabilities of our
425 model in robot simulation through two demonstrations.
426

4.1. Setup

427 **Dataset:** We evaluate the effectiveness of Video2Game
428 across three distinct scenes in various scenarios, includ-
429 ing “Gardenvase” [12], an outdoor object-centric scene;
430 the KITTI-360 dataset [37], a large-scale self-driving scene
431 with a sequence that forms a closed loop, suitable for car-
432 racing and Temple Run-like games; and finally, an indoor
433 scene from the VR-NeRF [73] dataset, showcasing the po-
434 tential for robot simulations.
435

436 **Metrics:** To evaluate the quality of the rendered im-
437 ages, we adopt PSNR, SSIM, and LPIPS [82]. For geome-
438 try reconstruction, we evaluate with LiDAR point cloud in
439 KITTI-360 dataset. Root mean square deviation (RMSE),
440 mean absolute error (MAE), and outlier rate are applied to
441 measure the disparity existing in estimated geometry.
442

443 **Our model:** For NeRF, we adopt hashgrid encoding [47]
and two-layer MLP for each header. For textured mesh,

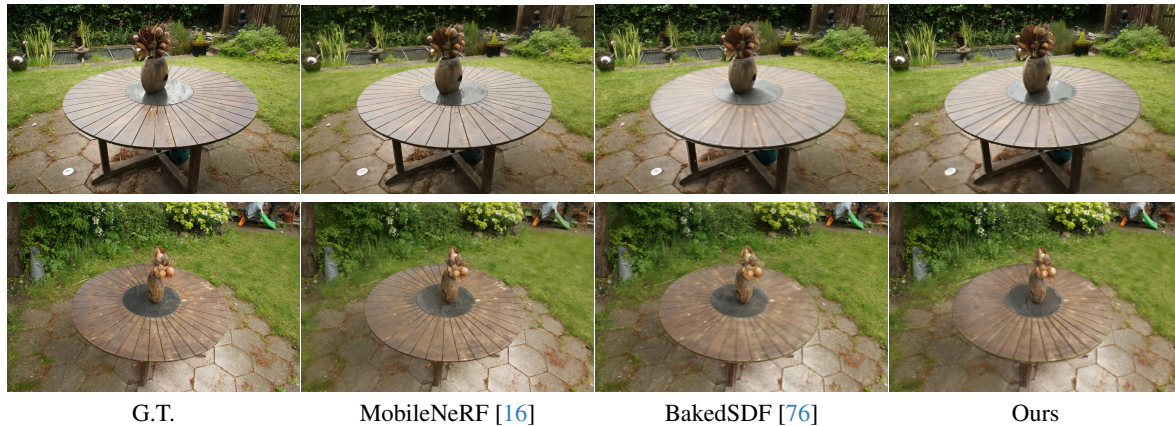


Figure 5. **Qualitative comparisons among mesh models.** We compare our mesh rendering method with others in Garden scene [12].



Figure 6. **Video2Game in Blender and Unreal Engine.**

Method	Outlier-% \downarrow	RMSE \downarrow	MAE \downarrow
Instant-NGP [48]	22.89	4.300	1.577
Nerfacto [64]	50.95	8.007	2.863
Gauss. Spl. [28]	91.08	11.768	8.797
BakedSDF* (offline) [76]	43.78	5.936	2.509
Video2Game (Our NeRF)	13.23	3.028	1.041

Table 2. **Quantitative evaluation on NeRF geometry.** Our NeRF renders significantly more accurate depth compared with the baselines. The unit is meter and the outlier threshold is 1.5 meters.

Method	Volume Rendering			Mesh Rasterization		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Vanilla NGP	27.46	0.853	0.165	22.54	0.716	0.350
+ Regularization terms	27.52	0.861	0.157	22.97	0.732	0.303
+ Monocular cues	27.62	0.871	0.131	23.35	0.765	0.246

Table 3. **Ablation studies on KITTI-360.**

444 we conduct marching cubes on the NeRF and post-process
 445 it to a fixed precision. We set the texture image size to
 446 4096x4096. For GLSL shader, we design a light-weight
 447 two-layer MLP, which enables efficient real-time rendering.
 448 For KITTI-360 (see Sec. 3.1), we divide the whole scene
 449 into 16 blocks and create a skydome mesh for the sky.

450 **Baselines:** To evaluate the visual and geometry quality
 451 of our model, we compare against SOTA approaches in
 452 neural rendering and neural reconstruction. **Instant-NGP**
 453 [48] is a NeRF-based method that exploits multi-resolution

hashing encoding. **Nerfacto** [64] extends the classic NeRF 454
 with learnable volumetric sampling and appearance embed- 455
 ding. **3D Gaussian Splatting** [28] leverages 3D Gaus- 456
 sians and achieves fast training and rendering. **MobileN-** 457
erf [16] adopts a hybrid NeRF-mesh representation. It can 458
 be baked into a texture map and enable real-time rendering. 459
BakedSDF [76] adopts a volume-surface scene representa- 460
 tion. It models view-dependent appearance efficiently by 461
 baking spherical Gaussians into the mesh. 462

4.2. Experimental results 463

Novel view synthesis: Tab. 1 shows the rendering perfor- 464
 mance and interactive compatibility of our model against 465
 the baselines on KITTI-360 [37] and Gardenvase [12]. Our 466
 NeRF achieves superior performance when compared to 467
 state-of-the-art neural volume render approaches across differ- 468
 ent scenes. Though [28] performs best in Gardenvase, 469
 it fails to handle the sparse camera settings in KITTI-360, 470
 where it learns bad 3D orientations of Gaussians. Our baked 471
 mesh outperforms other mesh rendering baselines signifi- 472
 cantly in KITTI-360 and performs similarly in Gardenvase 473
 as shown in Fig. 5. Additionally, our pipeline has the high- 474
 est interactive compatibility among all baselines. 475

Geometry reconstruction: Our model performs signifi- 476
 cantly better than the baseline regarding geometry accuracy 477
 (see Tab. 2). We provide some qualitative results in Fig. 478
 4, demonstrating that our model can generate high-quality 479
 depth maps and surface normals, whereas those produced 480
 by the baselines contain more noise. 481

Ablation study: To understand the contribution of each 482
 component in our model, we begin with the basic Instant- 483
 NGP [48] and sequentially reintroduce other components. 484
 The results in Tab. 3 show that our regularization and 485
 monocular cues improve the quality of both volume render- 486
 ing in NeRF and mesh rasterization. Additionally, we do 487
 observe a decline in rendering performance when convert- 488

489 ing NeRF into game engine-compatible meshes.

490 4.3. Video2Game

491 We have shown our approach’s effectiveness in rendering
492 quality and reconstruction accuracy across various setups.
493 Next, we demonstrate the construction of a web browser-
494 compatible game enabling player control and interaction
495 with the environment at over 100 FPS.

496 **Data preparation:** We build our environments based on
497 videos in Gardenvase [12], KITTI-360 [37] and VR-NeRF
498 [73] mentioned in Sec. 4.1, using our proposed approach.
499 The outcomes include executable environments with mesh
500 geometry, materials, and rigid-body physics, all encoded in
501 GLB and texture files.

502 **Game engine:** We build our game based on several key
503 components in our game engine mentioned in Sec. 3.4. By
504 leveraging them, our game generates a highly realistic vi-
505 sual rendering as well as physical interactions (see Fig. 1)
506 and runs smoothly at an interactive rate across various plat-
507 forms, browsers, and hardware setups (see Tab. 4). As for
508 other game engines (see Fig. 6), in Blender [17] we show-
509 case the compatibility of our exported assets with other
510 game engines. For Unreal [21], we further demonstrate a
511 real-time game demo where a humanoid robot can freely
512 interact within the Gardenvase scene, such as standing on
513 the table and kicking off the central vase. These prove the
514 compatibility of our proposed pipeline.

515 **Interactive game features:** *Movement in games:* Agents
516 can navigate the area freely within the virtual environment
517 where their actions follow real-world physics and are con-
518 strained by collision models. *Shooting game:* For realistic
519 shooting physics, we calculated the rigid-body collision dy-
520 namics for both the central vase and the surrounding scene
521 (see Fig. 3), separated using mesh semantic filtering. We
522 used a box collider for the vase and convex polygon collid-
523 ers for the background. The player shoots footballs with a
524 sphere collider at the vase on the table, causing it to fly off
525 and fall to the ground (see Fig. 1). *Temple-Run like game:*
526 The agent collects coins while running in the KITTI Loop
527 composed of four streets in KITTI-360. Obstructive chairs
528 on the road can be smashed thanks to pre-computed fracture
529 animations. The agent can also drive and push roadside ve-
530 hicles existing in the scene forward by crashing into them.
531 This interactivity is achieved through rigid-body dynamics
532 simulation and collision modeling.

533 **Robot simulation:** We demonstrate the potential of lever-
534 aging Video2Game for robot simulation using the VRNeRF
535 dataset. We reconstruct the scene and segment simulatable
536 rigid-body objects (e.g., the fruit bowl on the table). We
537 show two demos in Fig. 7: a Stretch Robot pushing the
538 bowl off the table and a Fetch Robot performing pick-and-
539 place actions. We employ PyBullet [18] to simulate the un-
540 derlying physics with the help of corresponding collision

	Platform	FPS (hz)	CPU-Usage (%)	GPU-Usage (%)
Mac M1 Pro	Mac OS, Chrome	102	34	70
Intel Core i9 + NV 4060	Windows, Edge	240	6	74
AMD 5950 + NV 3090	Linux, Chrome	144	20	40

Table 4. **Runtime Analysis.** Our interactive environment can run in real-time across various hardware setup and various platforms. User actions may slightly vary, which could lead to minor variations in runtime.

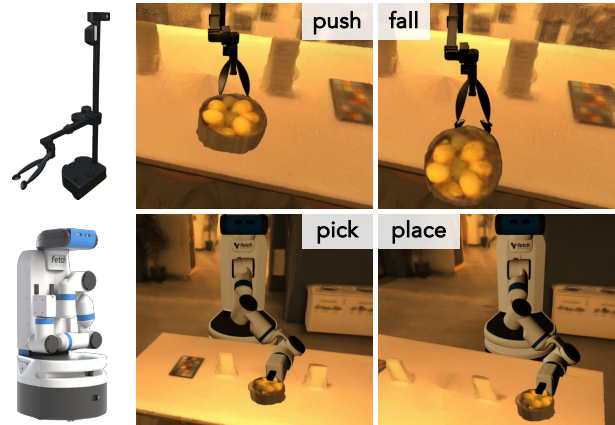


Figure 7. **Robot simulation in VRNeRF [73] dataset.** We demonstrate the possibility of conducting robot learning in our virtual environments using Stretch Robot [6] and Fetch Robot [5].

541 models. Since real-time grasping simulation is challenging,
542 following existing robot simulation frameworks [8, 9, 32],
543 objects near the Fetch gripper are automatically picked up.
544 This demonstrates our model’s ability to convert a real-time
545 video stream into a virtual environment, allowing robots to
546 rehearse before acting in the real environment.

5. Conclusion

547 We present a novel approach to converting real-world
548 video footage into playable, real-time, and interactive game
549 environments. Specifically, we combine the potential of
550 NeRF modeling with physics modeling and integrate them
551 into modern game engines. Our approach enables any in-
552 dividual to transform their surroundings into an interactive
553 digital environment, unlocking exciting possibilities for 3D
554 content creation, with promising implications for future
555 advancements in digital game design and robot simulation.
556

References

- 557
- 558
- 559
- 560
- 561
- 562
- 563
- 564
- 565
- [1] Cannon.js. <https://schteppe.github.io/cannon.js/>. 6
 - [2] GLSL. https://www.khronos.org/opengl/wiki/OpenGL_Shading_Language. 6
 - [3] Sketchbook. <https://github.com/swift502/Sketchbook>. 6
 - [4] WebGL. <https://www.khronos.org/webgl/>. 6

- 566 [5] Fetch Mobile Manipulator. <https://fetchrobotics.borealtech.com/robotics-platforms/fetch-mobile-manipulator/?lang=en>. 8
- 567
- 568
- 569 [6] Stretch® Research Edition. <https://hello-robot.com/product>. 8
- 570
- 571 [7] Xatlas. <https://github.com/mworchel/xatlas-python>. 5
- 572
- 573 [8] RoboTHOR: An Open Simulation-to-Real Embodied AI Platform. *CVPR*, 2020. 8
- 574
- 575 [9] ManipulaTHOR: A Framework for Visual Object Manipulation. In *CVPR*, 2021. 8
- 576
- 577 [10] Alexander Amini, Tsun-Hsuan Wang, Igor Gilitschenski, Wilko Schwarting, Zhijian Liu, Song Han, Sertac Karaman, and Daniela Rus. Vista 2.0: An open, data-driven simulator for multimodal sensing and policy learning for autonomous vehicles. In *ICRA*, 2022. 3
- 578
- 579
- 580
- 581
- 582 [11] Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kasten, and Tali Dekel. Text2live: Text-driven layered image and video editing. In *ECCV*, 2022. 2
- 583
- 584
- 585 [12] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, 2022. 4, 6, 7, 8
- 586
- 587
- 588 [13] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. *ECCV*, 2022. 2
- 589
- 590 [14] Shenchang Eric Chen and Lance Williams. View interpolation for image synthesis. In *Proceedings of the 20th annual conference on Computer graphics and interactive techniques*, 1993. 2
- 591
- 592
- 593
- 594 [15] Yun Chen, Frieda Rong, Shivam Duggal, Shenlong Wang, Xinchen Yan, Sivabalan Manivasagam, Shangjie Xue, Ersin Yumer, and Raquel Urtasun. Geosim: Realistic video simulation via geometry-aware composition for self-driving. In *CVPR*, 2021. 3
- 595
- 596
- 597
- 598
- 599 [16] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. *CVPR*, 2023. 2, 5, 7
- 600
- 601
- 602
- 603 [17] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 6, 8
- 604
- 605
- 606 [18] Erwin Coumans and Yunfei Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. <http://pybullet.org>, 2016–2021. 8
- 607
- 608
- 609 [19] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator, 2017. 1, 3
- 610
- 611
- 612 [20] Ainaz Eftekhari, Alexander Sax, Jitendra Malik, and Amir Zamir. Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10786–10796, 2021. 4
- 613
- 614
- 615
- 616
- 617 [21] Epic Games. Unreal engine. 1, 6, 8
- 618
- 619 [22] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinlong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 2
- 620
- 621 [23] Rockstar Games. Grand theft auto v. 2014. 1
- 622
- 623 [24] John K Haas. A history of the unity game engine. 2014. 1
- 624
- 625 [25] Benno Heigl, Reinhard Koch, Marc Pollefeys, Joachim Denzler, and Luc Van Gool. Plenoptic modeling and rendering from image sequences taken by a hand-held camera. In *DAGM-Symposium*, 1999. 2
- 626
- 627 [26] Ronghang Hu, Nikhila Ravi, Alexander C Berg, and Deepak Pathak. Worldsheet: Wrapping the world in a 3d sheet for view synthesis from a single image. In *ICCV*, 2021. 2
- 628
- 629
- 630 [27] Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18963–18974, 2022. 4
- 631
- 632
- 633
- 634 [28] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023. 5, 7
- 635
- 636
- 637
- 638 [29] Pradeep K Khosla and Takeo Kanade. Parameter identification of robot dynamics. In *IEEE conference on decision and control*, 1985. 3
- 639
- 640
- 641 [30] Seung Wook Kim, Yuhao Zhou, Jonah Philion, Antonio Torralba, and Sanja Fidler. Learning to simulate dynamic environments with gamegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1231–1240, 2020. 2, 3
- 642
- 643
- 644
- 645
- 646 [31] Seung Wook Kim, Jonah Philion, Antonio Torralba, and Sanja Fidler. Drivegan: Towards a controllable high-quality neural simulation. In *CVPR*, 2021. 2, 3
- 647
- 648
- 649 [32] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, Aniruddha Kembhavi, Abhinav Gupta, and Ali Farhadi. Ai2-thor: An interactive 3d environment for visual ai, 2022. 6, 8
- 650
- 651
- 652
- 653
- 654 [33] Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. Modular primitives for high-performance differentiable rendering. *ACM Transactions on Graphics*, 39(6), 2020. 5
- 655
- 656
- 657
- 658 [34] Yao-Chih Lee, Ji-Ze Genevieve Jang, Yi-Ting Chen, Elizabeth Qiu, and Jia-Bin Huang. Shape-aware text-driven layered video editing. *arXiv*, 2023. 2
- 659
- 660
- 661 [35] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996. 2
- 662
- 663
- 664 [36] Yuan Li, Zhi-Hao Lin, David Forsyth, Jia-Bin Huang, and Shenlong Wang. Climatenerf: Physically-based neural rendering for extreme climate synthesis. *arXiv*, 2022. 2, 4
- 665
- 666
- 667 [37] Yiyi Liao, Jun Xie, and Andreas Geiger. Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *IEEE TPAMI*, 2022. 4, 6, 7, 8
- 668
- 669
- 670 [38] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *ICCV*, 2021. 2
- 671
- 672
- 673 [39] Zhi-Hao Lin, Wei-Chiu Ma, Hao-Yu Hsu, Yu-Chiang Frank Wang, and Shenlong Wang. Neurmix: Neural mixture of planar experts for view synthesis. In *CVPR*, 2022. 2
- 674
- 675
- 676 [40] Zhi-Hao Lin, Bohan Liu, Yi-Ting Chen, David Forsyth, Jia-Bin Huang, Anand Bhattad, and Shenlong Wang. Urbanir: Large-scale urban scene inverse rendering from a single video, 2023. 2
- 677
- 678
- 679

- [41] Jeffrey Yunfan Liu, Yun Chen, Ze Yang, Jingkang Wang, Sivabalan Manivasagam, and Raquel Urtasun. Real-time neural rasterization for large scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8416–8427, 2023. 3
- [42] Fan Lu, Yan Xu, Guang Chen, Hongsheng Li, Kwan-Yee Lin, and Changjun Jiang. Urban radiance field representation with deformable neural mesh primitives, 2023. 3
- [43] Sivabalan Manivasagam, Shenlong Wang, Kelvin Wong, Wenyuan Zeng, Mikita Sazanovich, Shuhan Tan, Bin Yang, Wei-Chiu Ma, and Raquel Urtasun. Lidarsim: Realistic lidar simulation by leveraging the real world. In *CVPR*, 2020. 3
- [44] Willi Menapace, Stephane Lathuiliere, Sergey Tulyakov, Aliaksandr Siarohin, and Elisa Ricci. Playable video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10061–10070, 2021. 3
- [45] Willi Menapace, Stéphane Lathuiliere, Aliaksandr Siarohin, Christian Theobalt, Sergey Tulyakov, Vladislav Golyanik, and Elisa Ricci. Playable environments: Video manipulation in space and time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3584–3593, 2022. 3
- [46] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1, 2, 4
- [47] Thomas Müller. tiny-cuda-nn, 2021. 6
- [48] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *ACM TOG*, 2022. 2, 3, 4, 5, 7
- [49] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *CVPR*, 2021. 2
- [50] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, 2021.
- [51] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv*, 2021. 2
- [52] Cristina Portalés, José Luis Lerma, and Santiago Navarro. Augmented reality and photogrammetry: A synergy to visualize physical and virtual city environments. *ISPRS Journal of Photogrammetry and Remote Sensing*, 65(1):134–142, 2010. 1
- [53] Charalambos Poullis, Andrew Gardner, and Paul Debevec. Photogrammetric modeling and image-based rendering for rapid virtual environment creation. In *Proceedings of ASC*, 2004. 1
- [54] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE TPAMI*, 2020. 4
- [55] Christian Reiser, Richard Szeliski, Dor Verbin, Pratul P Srinivasan, Ben Mildenhall, Andreas Geiger, Jonathan T Barron, and Peter Hedman. Merf: Memory-efficient radiance fields for real-time view synthesis in unbounded scenes. *arXiv preprint arXiv:2302.12249*, 2023. 2, 5
- [56] Konstantinos Rematas, Andrew Liu, Pratul P Srinivasan, Jonathan T Barron, Andrea Tagliasacchi, Thomas Funkhouser, and Vittorio Ferrari. Urban radiance fields. In *CVPR*, 2022. 2
- [57] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv*, 2022. 2, 3
- [58] Uriel Singer, Shelly Sheynin, Adam Polyak, Oron Ashual, Iurii Makarov, Filippos Kokkinos, Naman Goyal, Andrea Vedaldi, Devi Parikh, Justin Johnson, et al. Text-to-4d dynamic scene generation. *arXiv*, 2023. 2, 3
- [59] Sanghyun Son, Yi-Ling Qiao, Jason Sewall, and Ming C Lin. Differentiable hybrid traffic simulation. *TOG*, 2022. 3
- [60] Pratul P Srinivasan, Richard Tucker, Jonathan T Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *CVPR*, 2019. 2
- [61] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022. 2
- [62] Richard Szeliski and Polina Golland. Stereo matching with transparency and matting. In *Sixth International Conference on Computer Vision*, 1998. 2
- [63] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *CVPR*, 2022. 4
- [64] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Justin Kerr, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, David McAllister, and Angjoo Kanazawa. Nerfstudio: A modular framework for neural radiance field development. *arXiv*, 2023. 5, 7
- [65] Jiaxiang Tang, Hang Zhou, Xiaokang Chen, Tianshu Hu, Er-rui Ding, Jingdong Wang, and Gang Zeng. Delicate textured mesh recovery from nerf via adaptive surface refinement. *arXiv preprint arXiv:2303.02091*, 2023. 2, 5
- [66] ADRIAN SAVARI Thomas, MOHD FAHRUL Hassan, MUSTAFFA Ibrahim, M Nasrull, A Rahman, SZ Sapuan, and F Ahmad. A study on close-range photogrammetry in image based modelling and rendering (imbr) approaches and post-processing analysis. *Journal of Engineering Science and Technology*, 14(4):1912–1923, 2019. 1
- [67] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *IROS*, 2012. 3
- [68] Shubham Tulsiani, Richard Tucker, and Noah Snavely. Layer-structured 3d scene inference via view synthesis. In *ECCV*, 2018. 2
- [69] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *CVPR*, 2022. 2

- 794 [70] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Vic-
795 tor Adrian Prisacariu. Nerf-: Neural radiance fields without
796 known camera parameters. *arXiv*, 2021. 2
- 797 [71] Bernhard Wymann, Eric Espié, Christophe Guionneau,
798 Christos Dimitrakakis, Rémi Coulom, and Andrew Sumner.
799 Torcs, the open racing car simulator. *Software available at*
800 *http://torcs.sourceforge.net*, 2000. 3
- 801 [72] Yuwen Xiong, Jingkang Ma, Wei-Chiu Wang, and Raquel
802 Urtaun. Ultralidar: Learning compact representations for
803 lidar completion and generation. *CVPR*, 2023. 3
- 804 [73] Linning Xu, Vasu Agrawal, William Laney, Tony Garcia,
805 Aayush Bansal, Changil Kim, Samuel Rota Bulò, Lorenzo
806 Porzi, Peter Kotschieder, Aljaž Božič, Dahua Lin, Michael
807 Zollhöfer, and Christian Richardt. VR-NeRF: High-fidelity
808 virtualized walkable spaces. In *SIGGRAPH Asia Conference*
809 *Proceedings*, 2023. 6, 8
- 810 [74] Zhenpei Yang, Yuning Chai, Dragomir Anguelov, Yin Zhou,
811 Pei Sun, Dumitru Erhan, Sean Rafferty, and Henrik Kret-
812 zschmar. Surfelgan: Synthesizing realistic sensor data for
813 autonomous driving. In *CVPR*, 2020. 3
- 814 [75] Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Mani-
815 vasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urta-
816 sun. Unisim: A neural closed-loop sensor simulator. *CVPR*,
817 2023. 2, 3
- 818 [76] Lior Yariv, Peter Hedman, Christian Reiser, Dor Verbin,
819 Pratul P Srinivasan, Richard Szeliski, Jonathan T Barron,
820 and Ben Mildenhall. Baked sdf: Meshing neural sdf for real-
821 time view synthesis. *Siggraph*, 2023. 2, 5, 7
- 822 [77] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and
823 Angjoo Kanazawa. Plenotrees for real-time rendering of
824 neural radiance fields, 2021. 4
- 825 [78] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sat-
826 tler, and Andreas Geiger. Monosdf: Exploring monocular
827 geometric cues for neural implicit surface reconstruction. *in*
828 *arXiv*, 2022. 4
- 829 [79] Haotian Zhang, Ye Yuan, Viktor Makovychuk, Yunrong
830 Guo, Sanja Fidler, Xue Bin Peng, and Kayvon Fatahalian.
831 Learning physically simulated tennis skills from broadcast
832 videos. *ACM Trans. Graph.* 3
- 833 [80] Haotian Zhang, Cristobal Sciuotto, Maneesh Agrawala, and
834 Kayvon Fatahalian. Vid2player: Controllable video sprites
835 that behave and appear like professional tennis players. *ACM*
836 *Transactions on Graphics (TOG)*, 40(3):1–16, 2021. 3
- 837 [81] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen
838 Koltun. Nerf++: Analyzing and improving neural radiance
839 fields. *arXiv*, 2020. 2
- 840 [82] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman,
841 and Oliver Wang. The unreasonable effectiveness of deep
842 features as a perceptual metric. In *CVPR*, 2018. 6
- 843 [83] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul De-
844 bevec, William T Freeman, and Jonathan T Barron. Nerfac-
845 tor: Neural factorization of shape and reflectance under an
846 unknown illumination. *ACM TOG*, 2021. 2
- 847 [84] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and An-
848 drew J. Davison. In-place scene labelling and understanding
849 with implicit scene representation, 2021. 4
- [85] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe,
and Noah Snavely. Stereo magnification: Learning view syn-
thesis using multiplane images. *arXiv*, 2018. 2
- [86] Yiming Zuo and Jia Deng. View synthesis with sculpted neu-
ral points. *in arXiv*, 2022. 2
- [87] Vlas Zyrianov, Xiyue Zhu, and Shenlong Wang. Learning to
generate realistic lidar point clouds. In *ECCV*, 2022. 3