

---

# DiVA-360: The Dynamic Visuo-Audio Dataset for Immersive Neural Fields

---

Anonymous Author(s)

Affiliation  
Address  
email

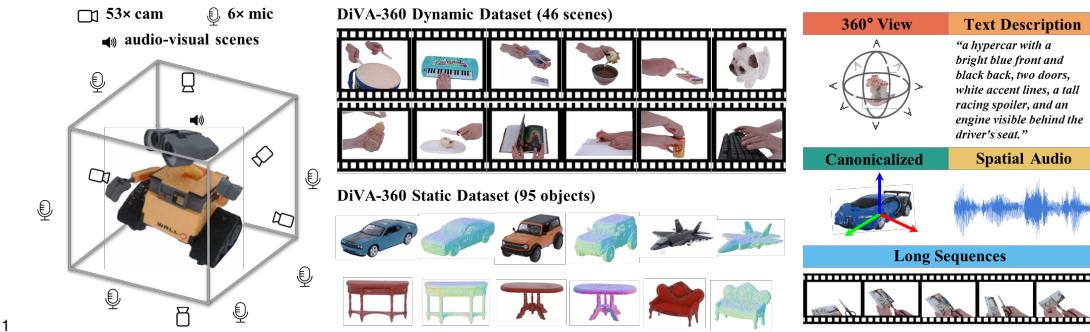


Figure 1: DiVA-360 is a **real-world** 360° multiview audio-visual dataset of dynamic and static scenes captured using a table-scale capture system consisting of 53 RGB cameras and 6 microphones. The DiVA-360 **dynamic** dataset consists of 46 prolonged (5s to 3 mins) dynamic scenes and interactions, while the DiVA-360 **static** dataset contains 95 objects across 11 categories and 30 static multi-object scenes. We also provide detailed text descriptions, 6 degrees of freedom (DoF) object pose canonicalization, spatial audio, and foreground-background segmentation. The goal of DiVA-360 is to enable continued progress in multimodal audio-visual neural field research.

## Abstract

2 Advances in neural fields are enabling high-fidelity capture of the shape and  
3 appearance of static and dynamic scenes. However, their capabilities lag behind  
4 those offered by representations such as pixels or meshes due to algorithmic  
5 challenges and the lack of large-scale real-world datasets. We address the dataset  
6 limitation with DiVA-360, a real-world 360° **dynamic visuo-audio** dataset with  
7 synchronized multimodal visual, audio, and textual information about table-scale  
8 scenes. It contains 46 dynamic scenes, 30 static scenes, and 95 static objects  
9 spanning 11 categories captured using a new hardware system using 53 RGB  
10 cameras at 120 FPS and 6 microphones for a total of 8.6M image frames and  
11 1360 s of dynamic data. We provide detailed text descriptions for all scenes,  
12 foreground-background segmentation masks, category-specific 3D pose alignment  
13 for static objects, as well as metrics for comparison. Our data, hardware and  
14 software, and code are available at <https://diva360.github.io/>.

## 15 1 Introduction

16 Neural fields [78], or neural implicit representations, have recently emerged as useful representations  
17 in computer vision, graphics, and robotics [78, 67] for capturing properties such as radiance [46,  
18 5, 4], shape [50, 81, 73, 49, 43, 38], dynamic motion [71, 35, 77, 20, 40, 7, 53], audio [41], and

19 language [32]. Their high fidelity, continuous representation, and implicit compression [19] properties  
20 make them attractive for immersive digital representation of our world.

21 However, despite their popularity, neural field capabilities remain far from that of conventional  
22 representations such as pixels (in 2D), and point clouds or meshes (in 3D). For instance, we can watch  
23 hours of videos with synchronized audio online, we can animate 3D meshes quickly on any device,  
24 and we have methods to quickly align 3D point clouds – tasks that currently cannot be achieved with  
25 neural fields. Recent work aims to enable these capabilities with most of it focusing on methods and  
26 algorithms [24, 79, 2] but **large-scale, real-world** datasets and benchmarks are equally important  
27 for continued progress [18, 13, 70, 8]. While some static [28, 55, 45, 75] and dynamic datasets [35]  
28 exist for neural fields, they have several limitations. First, existing dynamic datasets are limited to  
29 only a few scenes and only a few forward-facing cameras capturing for short durations. Second,  
30 static datasets may contain numerous objects and categories but lack within-category 3D alignment  
31 (aka *canonicalization*) – a common feature of synthetic 3D datasets like ShapenetCore [9, 56] that  
32 facilitates category-level learning [51, 66, 44, 25]. Third, many real-world datasets are captured with  
33 moving monocular cameras that cannot always provide sufficient multi-view cues for immersive  
34 reconstruction [21, 38, 53]. Finally, none of these datasets contain visually-adjacent modalities like  
35 audio and text similar to synthetic datasets [22, 23, 17].

36 We address these limitations by presenting **DiVA-360**, a real-world 360° dynamic visuo-audio dataset  
37 that contains synchronized multimodal visual, auditory, and textual information about table-scale  
38 objects and interactions. Rather than focus on the number of objects/scenes and categories, we instead  
39 focus on rich high-quality, synchronized, multimodal data about static and dynamic scenes. Our  
40 dynamic data includes high-resolution (1280×720), high-framerate (120 FPS), long (5s to 3 mins),  
41 and audio-synchronized videos captured simultaneously from 53 RGB cameras and 6 microphones  
42 spanning 360° volume within the capture space. Our static data includes high-resolution 53-view  
43 images and category-specific 6 degrees of freedom (DoF) pose alignment for object instances. Both  
44 static and dynamic scenes contain detailed text descriptions and foreground-background segmentation  
45 masks. In total, we provide **46 prolonged dynamic objects and interactions** spanning **1,360** seconds,  
46 **8.6M image frames**, **30 static multi-object scenes (5 clean and 25 messy)**, and **95 static objects**  
47 **from 11 categories**, all annotated with **8632 words of text descriptions** (see Table 2).

48 Capturing such large-scale multimodal data requires advances in capture systems, as well as bench-  
49 marking metrics. We have built a new capture system called **TRICS (Temporal Interaction Capture**  
50 **System**) which is designed to meet the multi-sensor synchronization, high-framerate, high-fidelity,  
51 and lighting requirements. For both the dynamic and static datasets we propose standardized metrics  
52 for reconstruction quality and runtime, and compare baseline methods on these metrics [47, 71].  
53 Our datasets, capture system, metric computation code, and annotations will be made publicly avail-  
54 able to the community at <https://diva360.github.io/>. To summarize, we make the following  
55 contributions:

- 56 • **TRICS**: A capture system specifically designed for 360° audio-visual capture of table-scale  
57 static and dynamic scenes with 53 RGB cameras and 6 microphones. We have developed our  
58 own hardware and software for sensor synchronization, capture, transfer, and calibration.
- 59 • **DiVA-360 Dynamic Dataset**: The largest audio-visual dataset for dynamic neural fields  
60 with 46 sequences (5s to 3 mins) captured at 120 FPS with synchronized spatial audio.
- 61 • **DiVA-360 Static Dataset**: We present a large static dataset of 30 scenes and 95 real-  
62 world objects spanning 11 categories captured in a category-aligned orientation and another  
63 random pose. This dataset includes information about the 6 DoF pose of objects.
- 64 • **Annotations**: For both dynamic and static scenes, we provide foreground-background  
65 segmentation masks, detailed text descriptions, trained models, and other metadata.

66 We believe our work can help the community take a leap from the current focus on static scenes and  
67 short dynamic videos toward a more holistic understanding of longer dynamic scenes, as well as  
68 text-to-4D scene generation [62], 3D object canonicalization [2], and audio-visual robotics [11].

## 69 2 Related Work

70 **Neural Fields:** Neural fields, coordinate-based neural networks, have generated considerable  
 71 interest in computer vision [78] because of their ability to represent geometry [44, 50, 12] and  
 72 appearance [46, 40, 63]. Neural radiance field (NeRF) [46] utilizes a Multilayer Perceptron (MLP)  
 73 to model density and color, leading to photorealistic novel view synthesis. Extensions of NeRF  
 74 have also been used to model shape with high-fidelity [81, 49, 37, 71]. Meanwhile, several methods  
 75 have made efforts to reduce the cost of constructing NeRF models [47, 59, 10]. Naturally, some  
 76 approaches have also turned their focus towards dynamic neural fields [71, 35, 53, 38, 52, 54, 20, 40].  
 77 However, these methods have thus far been limited to brief sequences and inadequate scene view  
 78 due to the limited camera capture range of the training data. A promising direction being explored is  
 79 the incorporation of au-  
 80 dio and language modal-  
 81 ities [41, 32, 26] with  
 82 neural fields. Our work  
 83 aims to facilitate the  
 84 broad spectrum of neu-  
 85 ral field research with  
 86 a more comprehensive  
 87 and richer dataset, from  
 88 higher-fidelity rendering  
 89 to faster training for 4D  
 90 dynamic field and from vi-  
 91 sual to multimodal visual-  
 92 audio-textual learning.

### 93 Multi-Camera Capture

94 **Systems:** Capturing  
 95 rich multimodal data re-  
 96 quires hardware and soft-  
 97 ware systems for captur-  
 98 ing rich data. The earliest

99 multi-camera capture systems were extensions of stereo cameras to 5–6 cameras [30] which were later  
 100 extended to capture a hemispherical volume [31] with up to 50 cameras for 3D and 4D reconstruction  
 101 using non-machine learning techniques [68]. Multi-camera systems have also been combined with  
 102 controllable lights to build *light stages* [16, 15]. Recent examples of multi-camera capture systems  
 103 include the panoptic camera systems [29, 29, 83]. Specialized systems have been built for table-scale  
 104 interactions, notably for hand interaction capture [85, 6]. However, these systems have a limited  
 105 number of cameras. Our TRICS system is specially designed for dense 53-view audio-visual capture  
 106 of table-scale scenes while also acting as a light stage.

107 **Large 3D Datasets:** Past advances in 3D learning have been largely driven by synthetic datasets  
 108 such as ShapeNet [9] and ModelNet [76], but their utility is somewhat curtailed by the absence of a  
 109 realistic appearance. Datasets for NeRF [46, 5], LLFF [45], and multiview Stereo (MVS) research  
 110 (DTU [1], Tanks & Temples [34], and BlendedMVS [80]) have provided sufficient multiview images,  
 111 yet their application is constrained to static scenes. Datasets like CO3D [55], OmniObject3D [75],  
 112 and ScanNeRF [14] and derivates like PeRFception [28] also focus on static objects and additionally  
 113 lack consistently orientation for objects.

114 Regarding dynamic datasets, BlockNeRF [65] contains hundreds of videos of street views, incor-  
 115 porating dynamic elements like cars and pedestrians but lacks a focus on objects, is short, and  
 116 lacks multiple views. DyNeRF [35], NDSD [82], ILFV [7], and Deep3DMV [39] use multiple  
 117 forward-facing cameras to take videos of dynamic activities so they lack views from behind. Besides,  
 118 DyNeRF and ILFV only contain short clips mostly around 10 s. Monocular dynamic view synthesis  
 119 datasets [53, 38, 52] capture monocular videos of human faces and human activities. However, the  
 120 use of a single camera restricts the visibility of all dynamic components from multiple viewpoints

Type	Dataset	Real	360° view	Dynamic	Caption	Canonical	Audio
Scene	DTU[27]	✓	✓	✗	✗	—	✗
	BlendedMVS[80]	✗	✓	✗	✗	—	✗
	ScanNet[13]	✓	✗	✗	✗	—	✗
	LLFF[45]	✓	✗	✗	✗	—	✗
	Mip-NeRF 360[5]	✓	✓	✗	✗	—	✗
	Block-NeRF[65]	✓	✗	✓	✗	—	✗
	DyNeRF[35]	✓	✗	✓	✗	—	✗
	HyperNeRF[53]	✓	✗	✓	✗	—	✗
	NDSD[82]	✓	✗	✓	✗	—	✗
	ILFV[7]	✓	✗	✓	✗	—	✗
	Deep3DMV[39]	✓	✗	✓	✗	—	✗
Object	ShapeNet[9]	✗	✓	✗	✗	✓	✗
	NeRF[46]	✗	✓	✗	✗	✓	✗
	CO3D[55]	✓	✓	✗	✗	✗	✗
	ScanNeRF[14]	✓	✓	✗	✗	✓	✗
	OmniObject3D[75]	✓	✓	✗	✗	✓	✗
Hybrid	ObjectFolder[23]	✗	✓	✗	✗	✓	✓
	PeRFception[28]	✓	✓	✗	✗	✗	✗
	Objaverse[17]	✗	✓	✓	✓	✓	✗
DiVA-360							

Table 1: We compare featured properties of our DiVA-360 with other multiview object-centric and scene-centric datasets. Our real-world 360° view dataset features the most comprehensive modalities and rich annotations of both static objects and dynamic scenes, aimed at promoting research in dynamic neural fields and 3D multimodal learning.

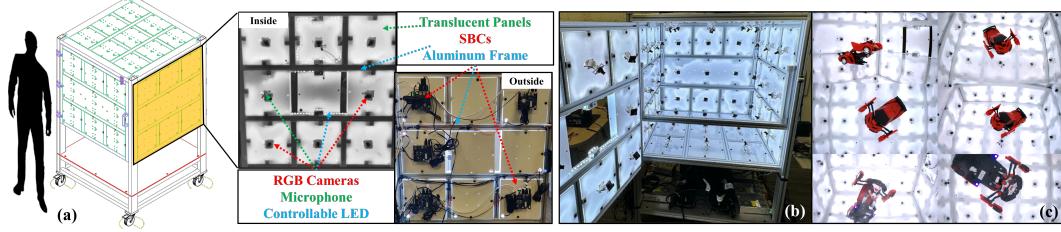


Figure 2: (a) TRICS is a refrigerator-sized aluminum frame that supports a  $1\text{m}^3$  capture volume mounted on wheels for mobility. Each side wall of the capture volume is divided into a  $3\times 3$  grid with each grid square containing sensors, LEDs, single-board computers (SBCs), and light diffusers. (b) Two walls of the capture volume act as doors for easy access to the capture volume. (c) We can acquire  $360^\circ$  RGB views of dynamic and static objects in this capture volume (6 views shown).

121 simultaneously, resulting in low effective multi-view factors (EMF)[21]. It is worth noting that  
 122 Objaverse [17] has created a substantial collection of top-tier 3D object models, characterized by their  
 123 wide-ranging categories and detailed annotations, inclusive of text descriptions and tags. Alongside  
 124 this, they introduced animated sequences to portray dynamic scenes. However, their data is not  
 125 sourced from the real world. Our dataset stands out by offering a  $360^\circ$  view of real dynamic scenes  
 126 with spatial audio and static objects captured in a consistent orientation within each category, all of  
 127 which are annotated with rich text descriptions (see Table 1).

### 128 3 Temporal Interaction Capture System (TRICS)

129 Our goal is to capture rich multimodal data of table-scale objects and interactions to enable further  
 130 research in high-fidelity audio-visual neural fields. To achieve this, we need a hardware and software  
 131 system that can capture high-framerate, high-resolution video and audio, and have the capability to  
 132 synchronize and calibrate these sensor streams. While commercial products exist for this purpose,  
 133 they do not meet all of our requirements. We therefore designed and built our own hardware and  
 134 software solution which we call the **Temporal Interaction Capture System (TRICS)**. Figure 2  
 135 shows our hardware system for capturing synchronized multimodal data. Please see the supplementary  
 136 document for more extensive details.

137 **TRICS Hardware:** Our system uses a mobile aluminum frame, housing a  $1\text{m}^3$  capture volume  
 138 outfitted with sensor panels across a  $3\times 3$  grid on each of its six sides (Figure 2 (a)). These panels  
 139 consist of RGB cameras, microphones, and programmable LED light strips, which together create  
 140 a versatile and uniformly light environment. The system is designed to handle large data output  
 141 through a custom communication setup that compresses and transmits data to a high-capacity control  
 142 workstation. This design, combining portability, comprehensive capture capabilities, and efficient  
 143 data management, allows for dynamic,  $360^\circ$  view capturing with low latency.

144 **TRICS Software:** While our hardware allows the capture of large-scale rich multimodal data,  
 145 controlling the sensors and LEDs, synchronizing and managing data, and camera calibration requires  
 146 specialized software which we have developed. For camera and microphone synchronization, we  
 147 adopt network-based synchronization [3] with an accuracy of 2–3ms. For camera calibration,  
 148 during each capture session, we affix transparent curtains with ArUco markers to the walls. Using  
 149 COLMAP [60, 61], we generate camera poses for the 53 cameras. The camera poses are further  
 150 refined using Instant-NGP’s [47] dense photometric loss for improved reconstruction quality. Finally,  
 151 we also built software for efficiently transferring terabytes of data from the control workstation to  
 152 cloud storage. We will release all of our software for community use.

### 153 4 DiVA-360 Dataset

154 We now describe our multimodal dynamic and static datasets that have been captured using TRICS.  
 155 While other datasets have focused on in-the-wild capture and a large number of categories [28, 55],

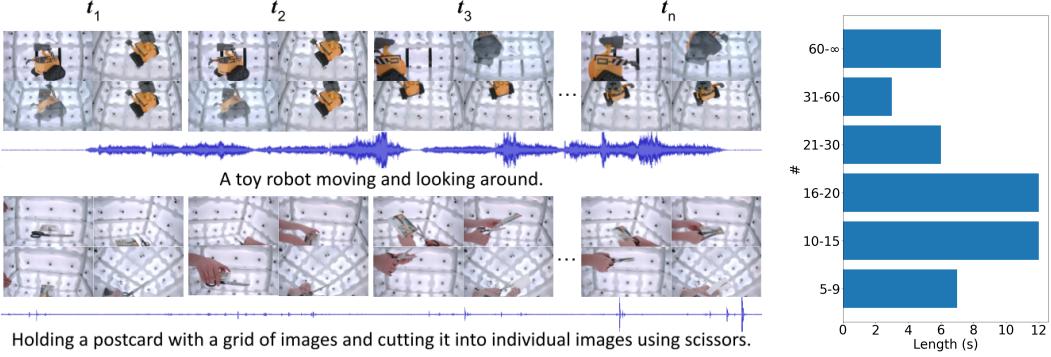


Figure 3: The DiVA-360 dynamic dataset contains multiview images and audio of dynamic objects (top row) and scenes (bottom), including routine human activities, musical instrument play, and objects moving or playing themselves. Our audio component (displaying 1 of 6 channels) contains both loud (e.g., top) and more subtle sounds (e.g., bottom) of objects in motion. On the right, we show the duration distribution of our sequences.

156 our goal is to focus on rich synchronized and multimodal data at high resolution, framerate, spanning  
 157 long duration, and captured from all 360° volume within the capture space. Our dataset contains  
 158 8.6 M image frames of 46 prolonged dynamic scenes over 1360 seconds, 95 static objects across 11  
 159 categories, and text descriptions totaling 8632 words.

#### 160 4.1 Dynamic Dataset

161 The DiVA-360 dynamic dataset contains synchronized long-duration audio-video of both moving  
 162 objects and hand interactions. Our goal is to make this dataset useful for learning long-duration  
 163 dynamic neural fields of appearance and audio – existing methods have been limited to only short  
 164 durations and lack audio [35, 71]. We captured 21 dynamic objects and 25 hand interactions with  
 165 objects for a total of 1360 seconds of audio-visual data from TRICS (see Figure 3). Our data  
 166 also contains masks for foreground-background segmentation and detailed text descriptions of each  
 167 sequence. To our knowledge, this is the largest-scale multimodal audio-visual dynamic dataset.

168 **Dynamic Objects:** We selected 21 dynamic objects that move and produce sounds. To be representative  
 169 of real-world motions, we chose scenes with different types of motion: (1) Slow motion: objects  
 170 that perform slow, continuous motions, e.g., music box and rotating world globe. (2) Fast motion:  
 171 objects that move or transform drastically, e.g., remote control cars and dancing toys. (3) Detailed  
 172 motion: objects that perform precise motions, e.g., a clock. (4) Repetitive motion: objects that repeat  
 173 the same motion pattern, e.g., Stirling engine and toys that sway left and right. (5) Random motion:  
 174 objects that perform indeterministic motions, e.g., plasma ball and sand in an hourglass. During  
 175 capture, all objects are placed on a transparent shelf for 360° view.

176 **Interactions:** In addition to dynamic objects, we also include 25 hand-object interaction scenes  
 177 representing real-world activities. The interactions included are hand activities commonly observed  
 178 in everyday life, such as flipping a book, replacing a toy’s batteries, and opening a lock. Most  
 179 interactions also generate subtle sounds, such as turning a page or opening a soda can. We hope these  
 180 hand-centric interactions encourage future modeling of complex hand dynamics. Similar to dynamic  
 181 objects, the objects to be manipulated are placed on a transparent shelf.

182 **Text Descriptions:** Each dynamic scene is accompanied by natural language descriptions at 3 levels  
 183 of detail. These descriptions are generated entirely by a human annotator without the assistance of  
 184 any automated tools. As such, we provide a human baseline for tasks that aim to align 3D visual  
 185 representations with natural language. The coarsest level aims to capture a broad summary of the  
 186 scene (“putting candy into a mug”), while finer levels increasingly describe appearance (“...the pieces  
 187 are in pink, green, orange, and black wrappers...”), relative position (“...candy scattered around a  
 188 black mug...”), number of hands, audio, and temporal progression. Across all 46 dynamic scenes, the

189 average length of the descriptions is 6.1, 18.4, and 38.7 words for the 3 levels of detail, amounting to  
190 a total of 2907 words.

191 **Foreground-Background Segmentation:** A major challenge in our dataset is the segmentation of  
192 foreground objects from background clutter. Manually segmenting every frame is infeasible due to  
193 the quantity and view inconsistency. Therefore, we developed a segmentation method using Instant  
194 NGP [47]. As preparation, we manually segment the foreground object in the first frame of one  
195 scene and train an I-NGP model on segmented images to refine coarse camera poses extracted from  
196 COLMAP[60, 61]. The refined pose is used for all downstream tasks. For each frame, we fit an  
197 I-NGP model that optimizes camera poses, lens distortion, and image latent vector. The model’s  
198 bounding box is then progressively reduced to remove background clutter. We then render trained  
199 I-NGP as binary masks. To further refine the masks, we removed connected components smaller than  
200 a threshold. Segmenting with this method is possible because all objects are placed around the center  
201 of TRICS. Since the segmentation is generated from I-NGP, the masks are multi-view consistent.

202 We analyze and compare the our DiVA-360 dynamic dataset with other dynamic datasets and two  
203 large object-centric  
204 datasets in Table 2.

205 DyNeRF [35] consists  
206 of only 6 short clips of  
207 forward-facing scenes.  
208 Block-NeRF [65] cap-  
209 tures a single long video  
210 of a street view from a  
211 moving vehicle. This

212 creates fleeting scenes  
213 that do not encompass  
214 full 360° camera cover-  
215 age. Though Objaverse

216 offers an expansive repository of animations, it lacks real-world scenes resulting in a domain gap.  
217 Our dataset incorporates extensive, prolonged dynamic sequences from 53 different viewpoints,  
218 effectively eliminating any blind spots within the dynamic components of the scene. In Section 5, we  
219 show how this data can be used to study dynamic neural field models and we provide metrics for  
220 comparison. We believe our synchronized audio and video data will spur new research.

Dataset	Object		Dynamic Scene		
	#Objects	#Categories	#Scenes	#Frames	Average length (s)
CO3D[55]	19k	50	—	—	—
OmniObject3D[75]	6k	190	—	—	—
Objaverse[17]	818k	21k	3k	$\infty$	—
DyNeRF[35]	—	—	6	37.8k	10
HyperNeRF[53]	—	—	17	13.8k	27
Deep3DMV[39]	—	—	96	3.8M	33
ILFV[7]	—	—	15	270.4k	13
Block-NeRF[65]	—	—	1	12k	100*
DiVA-360	95	11	46	8.6M	29.6

Table 2: Specifications of our DiVA-360 dataset and other dynamic datasets and large-scale object-centric datasets. \* indicates that despite Block-NeRF comprising a 100s-long video, it is made up of numerous transient street scenes, each with restricted view coverage.

## 221 4.2 Static Dataset

222 The DiVA-360 static dataset contains 360° 53-view images of 95 objects spanning 11 categories  
223 captured at two orientations, and 30 multi-object scenes captured (5 clean and 25 messy). Our goal is  
224 to make this dataset useful for categorical neural field learning for real-world objects. This dataset  
225 also contains language descriptions for each object instance at 3 levels of detail. Finally, to support  
226 category-level 3D learning, we provide the 6 degrees of freedom (DoF) pose of all object instances.  
227 In total, this data contains 11,660 images and 5725 words of text descriptions (see Figure 4).

228 **Objects:** The static dataset contains categories of objects similar to those found in ShapeNet [9]  
229 (e.g., cars, airplanes, cabinets), as well as common objects found in everyday life (e.g., utensils, mugs,  
230 keyboards). We also collected 30 multi-object scenes to resemble miniature rooms in both cluttered  
231 and cleanly-arranged settings [74]. For single objects, we first capture them in a pre-specified  
232 category-level canonical orientation followed by a random orientation. Similarly, for multi-objects,  
233 we first capture them in a “clean” state and rearrange the objects to capture increasingly messy states.

234 **Text Descriptions:** As in the dynamic dataset, we provide the natural language descriptions of the  
235 objects at 3 levels of detail generated entirely by a human annotator. For these static objects, the  
236 coarsest level provides a brief, generic description (“a black car”), while finer levels introduce details  
237 in appearance (“an all-black muscle car...”) and aim to differentiate the object within its class (“...with

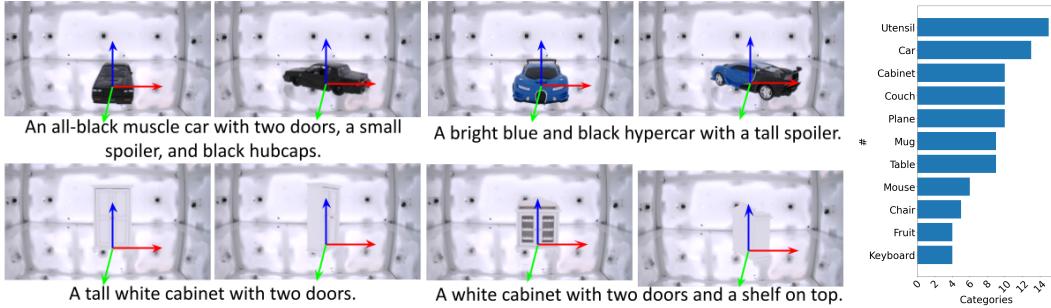


Figure 4: We showcase 4 pairs of captured static objects in a pre-defined canonical pose (left) and random poses (right) for each category. The histogram on the right shows our dataset includes a variety of objects from common shape categories.

238 two doors...”). Across all 125 static scenes, the average length of the descriptions is 4.8, 10.9, and  
239 30.2 words for the 3 levels of detail, amounting to a total of 5725 words.

240 **Canonicalization:** We provide category-level canonicalization for each of the static object categories,  
241 which provides an equivariant frame of reference that is consistent in position and orientation (3D  
242 pose) at the category level [56]. We automatically canonicalize objects using Canonical Fields  
243 Network (CaFi-Net) [2]. CaFi-Net uses a Siamese network architecture to extract equivariant field  
244 features from the neural field of object instances in arbitrary poses, and estimate a transformation that  
245 maps input to a canonical pose. We use CaFi-Net because of its self-supervised nature to estimate the  
246 canonical pose of each category.

247 While the scale of objects and categories in our dataset is smaller compared to CO3D and Om-  
248 niObject3D, it compensates with rich textual annotations, and each object category is consistently  
249 presented in a canonical pose, which facilitates the learning of category-level 3D representations. We  
250 hope that our DiVA-360 static dataset will accelerate research in object-centric learning in neural  
251 fields. In Section 5, we show how this data can be used to evaluate methods for static neural field  
252 reconstruction and provide metrics for 6 DoF pose canonicalization.

## 253 5 Benchmarks & Experiments

254 We show how DiVA-360 can be used to as a benchmark for neural field methods. We propose to  
255 standardize metrics for comparisons across methods and provide results of baselines on these metrics.  
256 All our experiments use Nvidia GPUs (2080Ti, 3090, 4090, A5000) for training and evaluation.

### 257 5.1 DiVA-360 Dynamic

258 Our goal is to evaluate existing methods for dynamic neural field reconstruction on image recon-  
259 struction quality. Specifically, we choose to compare two methods: (1) MixVoxels [71], a state-of-  
260 the-art method designed for dynamic radiance field reconstruction, and (2) Per-Frame I-NGP (PF  
261 I-NGP) [47], a static NeRF model which we fit to individual frames in all 46 sequences.

262 **Pre-processing:** We downsample the videos to 30 FPS and then segment all frames following  
263 Section 4. We select the top 35 out of 53 best cameras for training and hold out 6 cameras for testing.

264 **Metrics:** We use (a) Peak Signal-to-Noise Ratio (PSNR), (b) Structural Similarity Index Measure  
265 (SSIM) [58, 69], (c) Learned Perceptual Image Patch Similarity (LPIPS) [84] to measure the rendering  
266 quality, and Just Objectionable Difference (JOD) [42] to measure the visual difference between  
267 rendered video and ground truth, along with per-frame training/rendering time (s) for 6 testing views.

268 **Results and Analysis:** We quantitatively compare the two baseline methods MixVox-  
269 els and PF I-NGP in Table 3. Since PF I-NGP fits each frame individually and  
270 has more network capacity, its overall reconstruction performance is better than MixVox-

Baseline	PSNR↑	SSIM↑	LPIPS↓	JOD↑	Train (s/f)↓	Render (s/f)↓
MixVoxels[71]	$27.39 \pm 2.35$	$0.94 \pm 0.02$	$0.09 \pm 0.03$	$7.53 \pm 1.09$	$66.33 \pm 43.19$	$1.77 \pm 0.52$
PF I-NGP[47]	$28.13 \pm 3.50$	$0.95 \pm 0.03$	$0.08 \pm 0.04$	$7.61 \pm 0.93$	$48.85 \pm 4.73$	$0.67 \pm 0.18$

Table 3: We compare the rendering quality and train/render time of MixVoxels and PF I-NGP for dynamic scenes. Not surprisingly, PF I-NGP achieves higher rendering quality as it trains a model for each individual frame, though it lacks temporal coherence. In contrast, MixVoxels exhibits a much more compact form. The time is measured in seconds per frame (s/f).

271 els (see Figure 5 and Figure 6). However, MixVoxels only requires 2.7-4.7 MB  
 272 storage space per time step, which is six times smaller than PF I-NGP’s 29 MB.  
 273 Surprisingly, although  
 274 MixVoxels is designed  
 275 for dynamic scenes, its  
 276 training and inference times  
 277 are higher than PF I-NGP  
 278 (with a higher variance).  
 279 Besides, we also notice  
 280 that MixVoxel struggles  
 281 to capture the dynamic  
 282 components of the scenes,  
 283 leading to blurry and noisy  
 284 reconstruction (Figure 6  
 285 bottom). In the supplement,  
 286 we provide more details on  
 287 the differences between the  
 288 methods for various scenes.  
 289 Our primary objective is  
 290 to establish initial baseline  
 291 and serve as a resource for  
 292 future research aimed at enhancing these aspects.

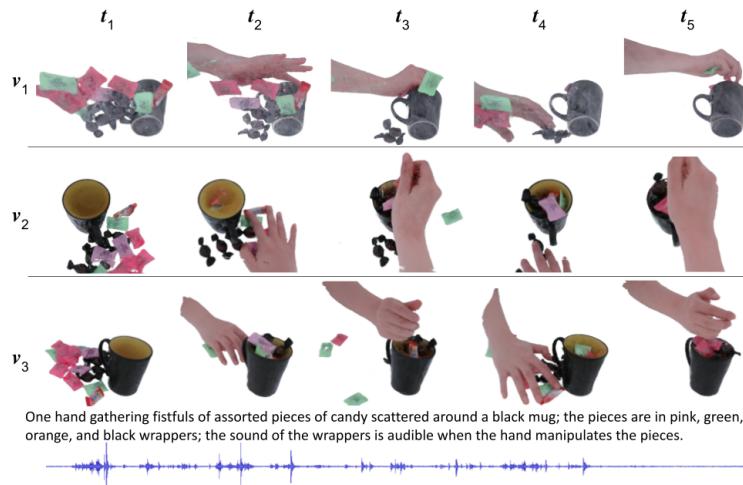


Figure 5: **Qualitative results of PF I-NGP on dynamic scene.** We show 3 views, 5 timestamps, text descriptions, and 1 audio channel. PF I-NGP has good per-frame quality but lacks temporal consistency.

## 293 5.2 DiVA-360 Static

294 To help build neural field representation for entire categories of objects, our dataset provides canonical-  
 295 oriented objects, which helps category-level understanding and generalization [64, 36, 48, 72, 57].  
 296 We use our object data captured in both canonical pose and random poses to benchmark neural object  
 297 canonicalization methods [2] to facilitate future research in category-level 3D perception. We also  
 298 provide a benchmark for rendering quality using I-NGP [47].

299 **Pre-processing:** We segment objects using the Segment Anything Model (SAM) [33] followed by  
 300 fitting an I-NGP model (Section 4).

301 **Metrics:** We train I-NGP on 34  
 302 best views and validate on 6 held-  
 303 out views. We use the same PSNR,  
 304 SSIM [58, 69], LPIPS [84] metrics for  
 305 evaluation. We provide a benchmark  
 306 on categorical neural object canonical-  
 307 ization using CaFi-Net [2]. For eval-  
 308 uation, we use the Instance-Level Con-  
 309 sistency (IC) and Category-Level Consistency metrics [56].

310 **Results and Analysis:** In Table 5, we notice that some category has much better results than others  
 311 leaving room for future improvements. Furthermore, segmentation artifacts appear in some scenes

Categories	CC ↓	IC ↓	Categories	CC ↓	IC ↓
chair	0.0411	0.0203	keyboard	0.1396	0.0719
table	0.0630	0.0292	car	0.0626	0.0224
cabinet	0.0736	0.0374	couch	0.0604	0.0343
mouse	0.0443	0.0494	plane	0.0538	0.0443
utensil	0.1536	0.1134			

Table 4: Object canonicalization benchmark using CaFiNet [2]. The performance is reasonably well for structured shapes such as chairs and tables but it struggles with smaller shape such as utensils.



Figure 6: **Qualitative comparison of PF I-NGP and MixVoxels and results on CaFiNet[2]**. We show rendering from PF I-NGP (top) and MixVoxels (bottom) on the left and show canonicalization results on the right (top: randomly oriented table NeRFs, bottom: canonicalized tables). Mix Voxels frequently generates artifacts like floaters or holes, particularly in the most dynamic components such as drumsticks and hands.

Category	PSNR $\uparrow$	SSIM* $\uparrow$	LPIPS* $\downarrow$	Category	PSNR $\uparrow$	SSIM* $\uparrow$	LPIPS* $\downarrow$
Chair	32.95 $\pm$ 2.11	98 $\pm$ 0.9	5 $\pm$ 2.3	Keyboard	31.57 $\pm$ 1.58	97 $\pm$ 0.8	7 $\pm$ 1.0
Table	38.11 $\pm$ 1.88	99 $\pm$ 0.3	2 $\pm$ 0.6	Car	34.06 $\pm$ 1.40	98 $\pm$ 0.4	3 $\pm$ 0.5
Cabinet	38.02 $\pm$ 1.65	99 $\pm$ 0.3	3 $\pm$ 0.8	Couch	32.87 $\pm$ 1.51	98 $\pm$ 0.6	4 $\pm$ 0.9
Mug	32.98 $\pm$ 2.12	98 $\pm$ 0.6	4 $\pm$ 0.7	Plane	34.04 $\pm$ 3.30	98 $\pm$ 0.8	4 $\pm$ 2.1
Fruit	34.91 $\pm$ 2.36	99 $\pm$ 0.3	3 $\pm$ 0.8	Utensil	36.45 $\pm$ 3.19	99 $\pm$ 0.4	6 $\pm$ 2.9
Mouse	35.03 $\pm$ 2.40	99 $\pm$ 0.2	3 $\pm$ 0.4	Scenes	27.41 $\pm$ 2.48	95 $\pm$ 2.0	8 $\pm$ 2.7

Table 5: Rendering quality for static object categories and scenes using Instant-NGP [47]. Metrics with (\*) are multiplied by  $10^2$  to accommodate space. Notice that some categories perform better than others, suggesting future methods that better generalize to a variety of shapes.

312 with more intricate and detailed geometries, encouraging future work in neural surface reconstruction.  
 313 The training time for each object instance was 60.51 s (35,000 iterations), and the mean rendering  
 314 time for 6 validation views was 0.29 s per frame. Table 4 shows the results of canonicalization  
 315 performance on various categories (we exclude two categories with few instances). CaFi-Net shows  
 316 good performance on categories with structured 3D shapes (e.g., tables) and sufficient instances.

## 317 6 Conclusion

318 We have introduced DiVA-360, a real-world 360° dynamic visuo-audio dataset that contains synchro-  
 319 nized multimodal visual, auditory, and textual information about table-scale objects and interactive  
 320 scenes. We propose a new TRICS capture system for rich multimodal data capture. DiVA-360  
 321 consists of a dynamic dataset of high-resolution, high-framerate, long (5s to 3 mins), and audio-  
 322 synchronized videos captured simultaneously from 53 RGB cameras and 6 microphones spanning  
 323 360° volume within the capture space. Our static data similarly includes high-resolution 53-view  
 324 images and category-specific 6 DoF pose alignment for object instances. In addition, both static and  
 325 dynamic scenes contain detailed text descriptions of the scene or interaction and mask annotations to  
 326 separate the foreground from the background.

327 **Limitations and Future Work:** Our focus is intentionally on high-quality multimodal data rather  
 328 than the number of scenes or objects, but large-scale learning is also essential and provided by datasets  
 329 like [55, 75]. Our metrics and evaluation are limited to images – in future work we will consider  
 330 metrics for audio and text. TRICS cannot capture scenes larger than table-scale – we plan to expand  
 331 this to larger volumes and in-the-wild capture in the future.

332 **Societal/Ethical Impact:** Our dataset does not reveal any private information and presents limited  
 333 means for misuse. However, future extensions of our work could contain private information that  
 334 can be misused. Furthermore, the multimodal nature of the dataset presents challenges in reducing  
 335 misuse and leaking of personal data that future research should explore.

336 **Checklist**

- 337 1. For all authors...
  - 338 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
339 contributions and scope? [Yes]
  - 340 (b) Did you describe the limitations of your work? [Yes]
  - 341 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
  - 342 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
343 them? [Yes]
- 344 2. If you are including theoretical results...
  - 345 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
  - 346 (b) Did you include complete proofs of all theoretical results? [N/A]
- 347 3. If you ran experiments (e.g. for benchmarks)...
  - 348 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
349 mental results (either in the supplemental material or as a URL)? [Yes]
  - 350 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
351 were chosen)? [Yes]
  - 352 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
353 ments multiple times)? [Yes]
  - 354 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
355 of GPUs, internal cluster, or cloud provider)? [Yes]
- 356 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - 357 (a) If your work uses existing assets, did you cite the creators? [Yes]
  - 358 (b) Did you mention the license of the assets? [Yes] MIT License
  - 359 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
  - 360 (d) Did you discuss whether and how consent was obtained from people whose data you're  
361 using/curating? [N/A]
  - 362 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
363 information or offensive content? [N/A]
- 364 5. If you used crowdsourcing or conducted research with human subjects...
  - 365 (a) Did you include the full text of instructions given to participants and screenshots, if  
366 applicable? [N/A]
  - 367 (b) Did you describe any potential participant risks, with links to Institutional Review  
368 Board (IRB) approvals, if applicable? [N/A]
  - 369 (c) Did you include the estimated hourly wage paid to participants and the total amount  
370 spent on participant compensation? [N/A]

371 **References**

- 372 [1] Henrik Aanæs, Rasmus Jensen, George Vogiatzis, Engin Tola, and Anders Dahl. Large-scale data for  
373 multiple-view stereopsis. *International Journal of Computer Vision*, 120, 11 2016.
- 374 [2] Rohith Agaram, Shaurya Dewan, Rahul Sajnani, Adrien Poulenard, Madhava Krishna, and Srinath Sridhar.  
375 Canonical fields: Self-supervised learning of pose-canonicalized neural fields. In *The IEEE Conference on  
376 Computer Vision and Pattern Recognition (CVPR)*, June 2023.
- 377 [3] Sameer Ansari, Neal Wadhwa, Rahul Garg, and Jiawen Chen. Wireless software synchronization of  
378 multiple distributed cameras. In *2019 IEEE International Conference on Computational Photography  
379 (ICCP)*, pages 1–9. IEEE, 2019.
- 380 [4] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P.  
381 Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. *ICCV*, 2021.
- 382 [5] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360:  
383 Unbounded anti-aliased neural radiance fields. *CVPR*, 2022.

- 384 [6] Samarth Brahmbhatt, Cusuh Ham, Charles C Kemp, and James Hays. Contactdb: Analyzing and predicting  
 385 grasp contact via thermal imaging. In *Proceedings of the IEEE/CVF conference on computer vision and*  
 386 *pattern recognition*, pages 8709–8719, 2019.
- 387 [7] Michael Broxton, John Flynn, Ryan Overbeck, Daniel Erickson, Peter Hedman, Matthew DuVall, Jason  
 388 Dourgarian, Jay Busch, Matt Whalen, and Paul Debevec. Immersive light field video with a layered mesh  
 389 representation. 39(4):86:1–86:15, 2020.
- 390 [8] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran  
 391 Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments.  
 392 *International Conference on 3D Vision (3DV)*, 2017.
- 393 [9] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio  
 394 Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An  
 395 Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University  
 396 — Princeton University — Toyota Technological Institute at Chicago, 2015.
- 397 [10] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. Mar  
 398 2022.
- 399 [11] Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vicenc Amengual Gari, Ziad Al-Halah, Vamsi Krishna  
 400 Ithapu, Philip Robinson, and Kristen Grauman. Soundspaces: Audio-visual navigation in 3d environments.  
 401 In *ECCV*, 2020.
- 402 [12] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. *Proceedings of IEEE*  
 403 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- 404 [13] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner.  
 405 Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern*  
 406 *Recognition (CVPR)*, IEEE, 2017.
- 407 [14] Luca De Luigi, Damiano Bolognini, Federico Domeniconi, Daniele De Gregorio, Matteo Poggi, and  
 408 Luigi Di Stefano. Scannerf: a scalable benchmark for neural radiance fields. In *Winter Conference on*  
 409 *Applications of Computer Vision*, 2023. WACV.
- 410 [15] Paul Debevec. The light stages and their applications to photoreal digital actors. *SIGGRAPH Asia*, 2(4):1–6,  
 411 2012.
- 412 [16] Paul Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, Westley Sarokin, and Mark Sagar.  
 413 Acquiring the reflectance field of a human face. In *Proceedings of the 27th annual conference on Computer*  
 414 *graphics and interactive techniques*, pages 145–156, 2000.
- 415 [17] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt,  
 416 Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects.  
 417 *arXiv preprint arXiv:2212.08051*, 2022.
- 418 [18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical  
 419 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255.  
 420 Ieee, 2009.
- 421 [19] Emilien Dupont, Adam Goliński, Milad Alizadeh, Yee Whye Teh, and Arnaud Doucet. Coin: Compression  
 422 with implicit neural representations. *arXiv preprint arXiv:2103.03123*, 2021.
- 423 [20] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic  
 424 monocular video. In *Proceedings of the IEEE International Conference on Computer Vision*, 2021.
- 425 [21] Hang Gao, Ruilong Li, Shubham Tulsiani, Bryan Russell, and Angjoo Kanazawa. Monocular dynamic  
 426 view synthesis: A reality check. In *NeurIPS*, 2022.
- 427 [22] Ruohan Gao, Yen-Yu Chang, Shivani Mall, Li Fei-Fei, and Jiajun Wu. Objectfolder: A dataset of objects  
 428 with implicit visual, auditory, and tactile representations. *arXiv preprint arXiv:2109.07991*, 2021.
- 429 [23] Ruohan Gao, Zilin Si, Yen-Yu Chang, Samuel Clarke, Jeannette Bohg, Li Fei-Fei, Wenzhen Yuan, and  
 430 Jiajun Wu. Objectfolder 2.0: A multisensory object dataset for sim2real transfer. In *Proceedings of the*  
 431 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10598–10608, 2022.
- 432 [24] Lily Goli, Daniel Rebain, Sara Sabour, Animesh Garg, and Andrea Tagliasacchi. nerf2nerf: Pairwise  
 433 registration of neural radiance fields. *arXiv preprint arXiv:2211.01600*, 2022.
- 434 [25] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan Russell, and Mathieu Aubry. AtlasNet: A  
 435 Papier-Mâché Approach to Learning 3D Surface Generation. In *Proceedings IEEE Conf. on Computer*  
 436 *Vision and Pattern Recognition (CVPR)*, 2018.
- 437 [26] Yudong Guo, Keyu Chen, Sen Liang, Yong-Jin Liu, Hujun Bao, and Juyong Zhang. Ad-nerf: Audio driven  
 438 neural radiance fields for talking head synthesis. In *Proceedings of the IEEE International Conference on*  
 439 *Computer Vision (ICCV)*, 2021.
- 440 [27] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engil Tola, and Henrik Aanæs. Large scale multi-view  
 441 stereopsis evaluation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages  
 442 406–413. IEEE, 2014.
- 443 [28] Yoonwoo Jeong, Seungjoo Shin, Junha Lee, Chris Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park.  
 444 Perfception: Perception using radiance fields. 2022.

- 445 [29] Hanbyul Joo, Hao Liu, Lei Tan, Lin Gui, Bart Nabbe, Iain Matthews, Takeo Kanade, Shohei Nobuhara, and  
 446 Yaser Sheikh. Panoptic studio: A massively multiview system for social motion capture. In *Proceedings of*  
 447 *the IEEE International Conference on Computer Vision*, pages 3334–3342, 2015.
- 448 [30] Takeo Kanade, Hiroshi Kano, Shigeru Kimura, Atsushi Yoshida, and Kazuo Oda. Development of a  
 449 video-rate stereo machine. In *Proceedings 1995 IEEE/RSJ International Conference on Intelligent Robots*  
 450 and *Systems. Human Robot Interaction and Cooperative Robots*, volume 3, pages 95–100. IEEE, 1995.
- 451 [31] Takeo Kanade and PJ Narayanan. Virtualized reality: perspectives on 4d digitization of dynamic events.  
 452 *IEEE Computer Graphics and Applications*, 27(3):32–40, 2007.
- 453 [32] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language  
 454 embedded radiance fields. *arXiv preprint arXiv:2303.09553*, 2023.
- 455 [33] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao,  
 456 Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything.  
 457 *arXiv:2304.02643*, 2023.
- 458 [34] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking  
 459 large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017.
- 460 [35] Tianye Li, Mira Slavcheva, Michael Zollhöfer, Simon Green, Christoph Lassner, Changil Kim, Tanner  
 461 Schmidt, Steven Lovegrove, Michael Goesele, and Zhaoyang Lv. Neural 3d video synthesis. *CoRR*,  
 462 abs/2103.02597, 2021.
- 463 [36] Xiaolong Li, He Wang, Li Yi, Leonidas Guibas, A Lynn Abbott, and Shuran Song. Category-level  
 464 articulated object pose estimation. *Proceedings of the IEEE Conference on Computer Vision and Pattern*  
 465 *Recognition*, 2020.
- 466 [37] Zhaoshuo Li, Thomas Müller, Alex Evans, Russell H Taylor, Mathias Unberath, Ming-Yu Liu, and Chen-  
 467 Hsuan Lin. Neuralangelo: High-fidelity neural surface reconstruction. In *IEEE Conference on Computer*  
 468 *Vision and Pattern Recognition (CVPR)*, 2023.
- 469 [38] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time  
 470 view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
 471 *Pattern Recognition (CVPR)*, 2021.
- 472 [39] Kai-En Lin, Lei Xiao, Feng Liu, Guowei Yang, and Ravi Ramamoorthi. Deep 3d mask volume for view  
 473 synthesis of dynamic scenes. In *ICCV*, 2021.
- 474 [40] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh.  
 475 Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–  
 476 65:14, July 2019.
- 477 [41] Andrew Luo, Yilun Du, Michael Tarr, Josh Tenenbaum, Antonio Torralba, and Chuang Gan. Learning  
 478 neural acoustic fields. *Advances in Neural Information Processing Systems*, 35:3165–3177, 2022.
- 479 [42] Rafał K Mantiuk, Gyorgy Denes, Alexandre Chapiro, Anton Kaplanyan, Gizem Rufo, Romain Bachy,  
 480 Trisha Lian, and Anjul Patney. Fovvideovdp: A visible difference predictor for wide field-of-view video.  
 481 *ACM Transactions on Graphics (TOG)*, 40(4):1–19, 2021.
- 482 [43] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy  
 483 networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on*  
 484 *Computer Vision and Pattern Recognition (CVPR)*, 2019.
- 485 [44] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy  
 486 networks: Learning 3d reconstruction in function space. In *Proceedings IEEE Conf. on Computer Vision*  
 487 and *Pattern Recognition (CVPR)*, 2019.
- 488 [45] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortíz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi,  
 489 Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling  
 490 guidelines. *ACM Transactions on Graphics (TOG)*, 2019.
- 491 [46] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren  
 492 Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- 493 [47] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives  
 494 with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022.
- 495 [48] David Novotny, Nikhila Ravi, Benjamin Graham, Natalia Neverova, and Andrea Vedaldi. C3dpo: Canonical  
 496 3d pose networks for non-rigid structure from motion. In *Proceedings of the IEEE International Conference*  
 497 *on Computer Vision*, 2019.
- 498 [49] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and  
 499 radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference*  
 500 *on Computer Vision*, pages 5589–5599, 2021.
- 501 [50] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf:  
 502 Learning continuous signed distance functions for shape representation. In *The IEEE Conference on*  
 503 *Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- 504 [51] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf:  
 505 Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF*  
 506 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

- 507 [52] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz,  
 508 and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021.
- 509 [53] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman,  
 510 Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for  
 511 topologically varying neural radiance fields. *ACM Trans. Graph.*, 40(6), dec 2021.
- 512 [54] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance  
 513 fields for dynamic scenes. *arXiv preprint arXiv:2011.13961*, 2020.
- 514 [55] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David  
 515 Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction.  
 516 In *International Conference on Computer Vision*, 2021.
- 517 [56] Rahul Sajnani, Adrien Poulenard, Jivitesh Jain, Radhika Dua, Leonidas J. Guibas, and Srinath Sridhar.  
 518 Condor: Self-supervised canonicalization of 3d pose for partial shapes. In *The IEEE Conference on*  
 519 *Computer Vision and Pattern Recognition (CVPR)*, June 2022.
- 520 [57] Rahul Sajnani, AadilMehti Sanchawala, Krishna Murthy Jatavallabhula, Srinath Sridhar, and K. Madhava  
 521 Krishna. Draco: Weakly supervised dense reconstruction and canonicalization of objects, 2020.
- 522 [58] Umme Sara, Morium Akter, and Mohammad Sharif Uddin. Image quality assessment through fsim, ssim,  
 523 mse and psnr—a comparative study. *Journal of Computer and Communications*, 7(3):8–18, 2019.
- 524 [59] Sara Fridovich-Keil and Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa.  
 525 Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022.
- 526 [60] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on*  
 527 *Computer Vision and Pattern Recognition (CVPR)*, 2016.
- 528 [61] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view  
 529 selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016.
- 530 [62] Uriel Singer, Shelly Sheynin, Adam Polyak, Oron Ashual, Iurii Makarov, Filippos Kokkinos, Naman Goyal,  
 531 Andrea Vedaldi, Devi Parikh, Justin Johnson, and Yaniv Taigman. Text-to-4d dynamic scene generation.  
 532 *arXiv:2301.11280*, 2023.
- 533 [63] Vincent Sitzmann, Michael Zollhoefer, and Gordon Wetzstein. Scene representation networks: Continuous  
 534 3d-structure-aware neural scene representations. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-  
 535 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32.  
 536 Curran Associates, Inc., 2019.
- 537 [64] Srinath Sridhar, Davis Rempe, Julien Valentin, Sofien Bouaziz, and Leonidas J. Guibas. Multiview  
 538 aggregation for learning category-specific shape reconstruction. In *Advances in Neural Information*  
 539 *Processing Systems (NeurIPS)*, 2019.
- 540 [65] Matthew Tancik, Vincent Casser, Xincheng Yan, Sabeek Pradhan, Ben Mildenhall, Pratul Srinivasan,  
 541 Jonathan T. Barron, and Henrik Kretzschmar. Block-NeRF: Scalable large scene neural view synthesis.  
 542 *arXiv*, 2022.
- 543 [66] Maxim Tatarchenko\*, Stephan R. Richter\*, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox.  
 544 What do single-view 3d reconstruction networks learn? 2019.
- 545 [67] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Treitschke, Wang Yifan, Christoph  
 546 Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, et al. Advances in neural rendering.  
 547 In *Computer Graphics Forum*, volume 41, pages 703–735. Wiley Online Library, 2022.
- 548 [68] Christian Theobalt, Marcus A Magnor, Pascal Schüller, and Hans-Peter Seidel. Combining 2d feature  
 549 tracking and volume reconstruction for online video-based human motion capture. *International Journal*  
 550 *of Image and Graphics*, 4(04):563–583, 2004.
- 551 [69] Paul Upchurch, Noah Snavely, and Kavita Bala. From a to z: Supervised transfer of style and content  
 552 using deep neural network generators. *CoRR*, abs/1603.02003, 2016.
- 553 [70] Kashi Venkatesh Vishwanath, Amin Vahdat, Ken Yocom, and Diwaker Gupta. Modelnet: Towards a  
 554 datacenter emulation environment. In Henning Schulzrinne, Karl Aberer, and Anwitaman Datta, editors,  
 555 *Peer-to-Peer Computing*, pages 81–82. IEEE, 2009.
- 556 [71] Feng Wang, Sinan Tan, Xinghang Li, Zeyue Tian, and Huaping Liu. Mixed neural voxels for fast multi-view  
 557 video synthesis. *arXiv preprint arXiv:2212.00190*, 2022.
- 558 [72] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J. Guibas. Normal-  
 559 ized object coordinate space for category-level 6d object pose and size estimation. In *The IEEE Conference*  
 560 *on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- 561 [73] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. Neus2:  
 562 Fast learning of neural implicit surfaces for multi-view reconstruction, 2022.
- 563 [74] QiuHong Anna Wei, Sijie Ding, Jeong Joon Park, Rahul Sajnani, Adrien Poulenard, Srinath Sridhar, and  
 564 Leonidas Guibas. Lego-net: Learning regular rearrangements of objects in rooms. In *Proceedings of the*  
 565 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19037–19047, 2023.
- 566 [75] Tong Wu, Jiarui Zhang, Xiao Fu, Yuxin Wang, Jiawei Ren, Liang Pan, Wayne Wu, Lei Yang, Jiaqi  
 567 Wang, Chen Qian, Dahua Lin, and Ziwei Liu. Omniobject3d: Large-vocabulary 3d object dataset for  
 568 realistic perception, reconstruction and generation. *IEEE/CVF Conference on Computer Vision and Pattern*  
 569 *Recognition (CVPR)*, 2023.

- 570 [76] Zhirong Wu, Shuran Song, Aditya Khosla, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d  
571 shapenets: A deep representation for volumetric shape modeling. In *IEEE Conference on Computer Vision*  
572 and *Pattern Recognition (CVPR)*, Boston, USA, June 2015.
- 573 [77] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for  
574 free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
575 *Recognition (CVPR)*, pages 9421–9431, 2021.
- 576 [78] Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico Tombari,  
577 James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural fields in visual computing and beyond.  
578 *Computer Graphics Forum*, 2022.
- 579 [79] Guandao Yang, Serge Belongie, Bharath Hariharan, and Vladlen Koltun. Geometry processing with neural  
580 fields. *Advances in Neural Information Processing Systems*, 34:22483–22497, 2021.
- 581 [80] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan.  
582 Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. *Computer Vision and*  
583 *Pattern Recognition (CVPR)*, 2020.
- 584 [81] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In  
585 *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- 586 [82] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of  
587 dynamic scenes with globally coherent depths from a monocular camera. In *Proceedings of the IEEE/CVF*  
588 *Conference on Computer Vision and Pattern Recognition*, pages 5336–5345, 2020.
- 589 [83] Zhixuan Yu, Jae Shin Yoon, In Kyu Lee, Prashanth Venkatesh, Jaesik Park, Jihun Yu, and Hyun Soo  
590 Park. Humbi: A large multiview dataset of human body expressions. In *Proceedings of the IEEE/CVF*  
591 *Conference on Computer Vision and Pattern Recognition*, pages 2990–3000, 2020.
- 592 [84] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
593 effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- 594 [85] Christian Zimmermann, Duygu Ceylan, Jimei Yang, Bryan Russell, Max Argus, and Thomas Brox.  
595 Freihand: A dataset for markerless capture of hand pose and shape from single rgb images. In *Proceedings*  
596 *of the IEEE/CVF International Conference on Computer Vision*, pages 813–822, 2019.