

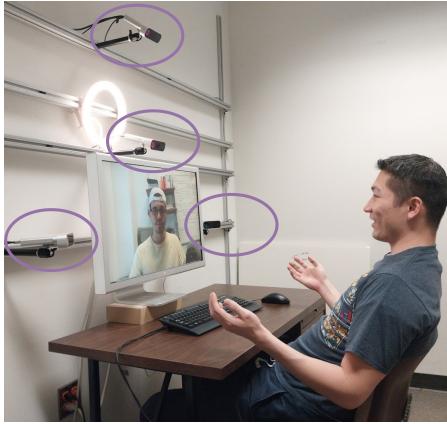
# Bringing Telepresence to Every Desk

Shengze Wang Ziheng Wang Ryan Schmelzle  
YoungJoong Kwon Liujie Zheng Soumyadip Sengupta Henry Fuchs

UNC Chapel Hill

shengzew@cs.unc.edu wzh@unc.edu rysch01@live.unc.edu youngjoong@cs.unc.edu  
liujiez@email.unc.edu {ronisen,fuchs}@cs.unc.edu

Affordable Desktop System with  
4 RGBD cameras



Free Viewpoint Video for Personal 3D Video Conferencing



Figure 1: We present a capturing and rendering system designed for personal telepresence. (left) Our system utilizes four RGBD cameras to render high-resolution free-viewpoint videos, which are crucial to immersive 3D video conferencing. (right) We show synthesized images from different novel viewpoints across different frames.

## Abstract

*In this paper, we work to bring telepresence to every desktop. Unlike commercial systems, personal 3D video conferencing systems must render high-quality videos while remaining financially and computationally viable for the average consumer. To this end, we introduce a capturing and rendering system that only requires 4 consumer-grade RGBD cameras and synthesizes high-quality free-viewpoint videos of users as well as their environments.*

*Experimental results show that our system renders high-quality free-viewpoint videos without using object templates or heavy pre-processing. While not real-time, our system is fast and does not require per-video optimizations. Moreover, our system is robust to complex hand gestures and clothing, and it can generalize to new users. This work provides a strong basis for further optimization, and it will help bring telepresence to every desk in the near future. The code and dataset will be made available.*

## 1. Introduction

In recent years, video conferencing has become ubiquitous in our daily lives, facilitated by software such as Zoom,

Gather, and Meet. However, 2D video conferencing fails to provide immersive experiences and realistic conversations. Thus, researchers have worked for decades to develop 3D telepresence systems that enable remote users to virtually share the same space [47, 12, 17, 25, 34].

Unlike systems using cumbersome headsets [41], encumbrance-free systems [27, 67, 34, 47] leverage regular/autostereo displays to provide immersive experiences without requiring the users to wear headsets. Recent examples such as Project Starline [27] and VirtualCube [67] have achieved unprecedented levels of realism. However, these commercial systems require intricate hardware setups and dedicated physical spaces (*i.e.* booths or rooms) making them inaccessible to the general public.

Unlike commercial systems, personal telepresence must be affordable and accessible to a wider range of users. Our system achieves this through the use of sparse and consumer-quality sensors. Specifically, our system employs only 4 Microsoft Azure Kinect RGBD cameras, which cost approximately \$1600 in total - less than half the price of AR glasses such as Magic Leap and Hololens, which can cost over \$3500. To use our system, a user only needs to install the four cameras around their monitor, without requiring

any additional hardware or dedicated capture spaces/booths. This setup provides two benefits: first, it enhances the sense of realism by placing the users in the context of their actual surroundings, creating a more immersive telepresence experience and contributing to a sense of personal connections between users. Second, this lightweight system reduces the financial and physical burden on consumers. However, this lightweight setup also presents several challenges:

1. *Sparse viewpoints*: The sparsity of viewpoints results in wider baselines and larger perspective changes, making accurate reconstructions more challenging.
2. *Inaccurate and noisy depth sensors*: Depth measurements from RGBD cameras are biased, inducing alignment errors (Fig.2). Moreover, depth values vary across frames, resulting in flickering.
3. *Background*: Without dedicated booths, the system needs to synthesize high-quality renderings of the background in addition to the foreground.

To this end, we propose a relatively capturing and rendering system. Our system synthesizes high-quality novel view images of human subjects and backgrounds given 4 input RGBD streams. To improve the reconstruction quality, we introduce the Multi-layer Point Cloud (MPC), a novel volumetric representation designed for biased depth inputs. MPC is constructed by sweeping point clouds from the input viewpoints. Compared to the conventional novel-view-depth-sweeping, MPC enables more accurate reconstruction of slanted surfaces and contour regions, thus reducing flickering. To further improve temporal smoothness, our temporal neural renderer aggregates information across frames. To achieve high-resolution video synthesis under limited GPU memory, we introduce the Spatial Skip Connection inspired by UNet[51] skip connections. To study the efficacy of our system, we created a dataset tailored for personal 3D video conferencing, *i.e.* the Personal Telepresence Dataset. Experiments show that our system outperforms baseline methods. Ablation studies further show that each proposed module improves the stability and accuracy of the results. Our work can be summarized as follows:

- We present a relatively affordable capturing and rendering system for desktop telepresence using only 4 RGBD cameras. Our setup can be easily replicated on any desk without a dedicated space or booth.
- We designed the Multi-Layer Point Cloud (MPC), a new volumetric representation that improves reconstruction from RGBD inputs. We further improve the rendering stability and memory efficiency via our temporal renderer and Spatial Skip Connections.
- Experiments show that our system outperforms recent competitive methods in the synthesis of free-viewpoint videos. Our method is fast and does not require object

templates or heavy pre-processing. Moreover, it is robust to complex hand gestures and clothing, and it can generalize to new users.

## 2. Related Work

### 2.1. 3D Video Conferencing Systems

Since the pioneering works of Cruz-Neira *et al.* [9] and Raskar *et al.* [47]), there has been a plethora of work [22, 35, 36, 24, 65, 27, 67] on headset-free 3D video conferencing.

**Personal Vs. Commercial Systems** The key difference between commercial and personal systems is that personal systems must achieve high-quality results while remaining financially and computationally viable for consumers. As a result, the quality of sensors and complexity of the systems tend to be lower than commercial alternatives. Moreover, recent commercial systems such as Starline [27] and VirtualCube [67] resemble photo-booths blocking the background. However, in personal usage, it is impractical to dedicate a space specifically for video conferencing. Moreover, the background is often crucial to creating realistic, informal, and intimate personal interactions. Therefore, we envision personal systems that can be installed in a typical room setting and capture background environments.

### 2.2. Novel View Synthesis

**Dynamic View synthesis** There are many different approaches to synthesize free-viewpoint videos; Yoon *et al.* [63] utilizes multi-view stereo and monocular depth estimators to generate 3D videos without per-video optimization or prior knowledge of the scene. Additionally, many NeRF-based approaches [61, 29, 16, 57] encode dynamic scenes as spatio-temporal radiance fields. [61, 29, 16, 57] learn a radiance field and a motion field for each frame. TöRF [3] uses Time-of-Flight sensors to achieve better modeling of both static and dynamic scenes. Approaches like Nerfies [43], HyperNeRF [44], and Neural 3D Videos (N3V) [28] use latent codes to help model dynamic contents. However, each of these approaches require extensive training time. Moreover, all these approaches require buffering the video in advance, making them unsuitable for instant live-streaming applications. ENERF [30] demonstrated notable improvements on generalization over prior works [64, 4, 59, 7] while achieving real-time rendering without optimization. LookinGood [38] uses multiple RGBD cameras to reconstruct a human mesh in real-time and uses a neural network to render the colored mesh from high-quality novel views. While this approach heavily depends on the quality of the preprocessed reconstruction, our approach uses raw RGBD frames to eliminate the need for preprocessing and naturally renders the background and foreground at the same time.

**Human-Specific Approaches.** Some recent works [46, 31, 54, 45, 26] exclusively focus on animating clothed hu-

Table 1: **Prior Telepresence Systems:** A high-level comparison between our system and prior approaches. Features can theoretically be implemented but have not yet been shown in prior works; these features are indicated with \*

Systems	Desktop	Affordable	Generalizable	Background	Quality	Real-Time	Complexity
Maimone <i>et al.</i> [36]	✓	✓	Yes	Yes	Low	✓	Low
Holoportation [41]	✗	✗	Yes	*	Medium	✓	High
LookinGood[38]	✗	✗	*	*	High	✓	High
Project Starline[27]	✗	✗	Yes	No	High	✓	High
VirtualCube[67]	✓	✓	Yes	*	Medium	✓	Low
Ours	✓	✓	Yes	Yes	High	0.19-1.28fps	Low

mans. They often use colored videos as inputs and leverage human body templates (*e.g.* SMPL[32] and STAR[42]) and deep textures. On the other hand, HVSNet [40] uses monocular RGBD videos to render human subjects from feature point clouds. Our work requires reconstruction of general scenes including both humans and objects, making these human-only approaches inapplicable. Moreover, our application requires a detailed depiction of the upper body and hands, presenting a challenge for human models designed for portrait or full-body views. Therefore, these approaches would be ill-suited for our application.

**Preprocessed 3D Representations.** Some works [49, 50, 23, 52, 2] rely on raw 3D geometry (*e.g.* point clouds and meshes) generated from multi-view stereo software such as COLMAP[53]. Such approaches are less prone to generating fog-like artifacts common in NeRF-based methods. Methods like [49, 50] also demonstrate good generalization to new scenes. Point cloud-based neural rendering approaches [23, 52, 60] show impressive sharpness and details in large scenes with thin structures. However, these approaches either require lengthy training for each frame or do not show competitive results.

### 3. Motivation for Multi-Layer Point Cloud

In this section, we discuss the motivation for our novel volumetric representation, *i.e.*, Multi-layer Point Cloud.

**Conventional Novel View Sweep Volume.** Many RGB-only approaches [4, 62, 18, 6, 64, 30, 13, 68, 20, 8, 15] estimate scene geometry by depth-sweeping from the novel view camera. This approach is proven effective for RGB-only scenarios and is thus naturally extensible to RGBD-based novel view synthesis. For example, VirtualCube [67] first generates an averaged novel view depth map from input depth maps. It then generates depth candidates via a novel view depth-sweep around the averaged depth map. However, we show that this simple extension can cause inaccurate reconstructions due to biases in depth sensors.

**Challenges Induced by Depth Bias.** Modern commodity depth sensors (*e.g.* Microsoft Azure Kinect) can now achieve good quality but still suffer from noise and bias. Depth bias is especially important because it can cause notable misalignment between views. Fig. 2(a) shows an example of such misalignment between point clouds from two cameras, implying that the depth values are offset from

the true geometry. While conventional volumes can locally search for the true surfaces, their sweep/search volumes might not cover the true surface. As illustrated by Fig. 2(b)-(d), the direction of depth bias lies along the viewing direction of the input camera, but the depth sweep is in the novel view direction. As a result, the sweep volume (green) inevitably misses the true surface, leading to incomplete reconstructions. This effect worsens as the surface becomes more slanted in the novel view (Fig. 2(c)-(d)). Moreover, the sweep volume (green) might completely miss the true surface when the surface is too steep in the novel view (Fig. 2(d)). This problem is amplified by the unstable depth values that vary between frames, inducing flickering around object contours and steep surfaces. As a result, previous methods require careful post-processing to stabilize generated videos (*e.g.* averaging adjacent frames [67, 38]).

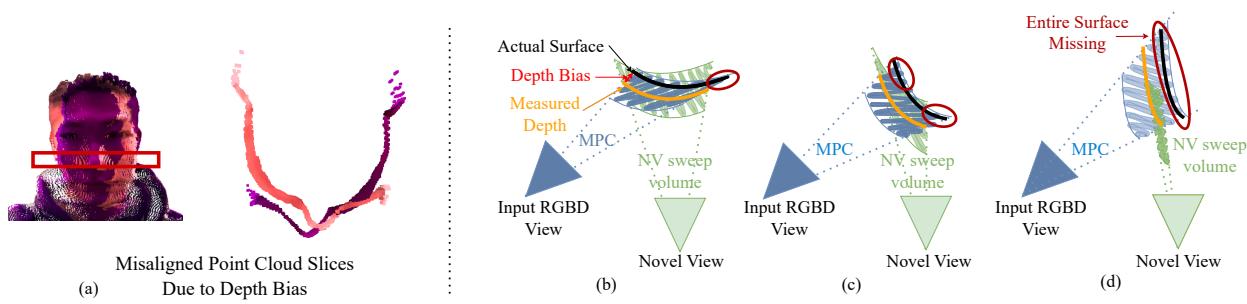
**Multi-layer Point Cloud (MPC).** Depth biases are difficult to remove because they depend on a wide range of factors [14, 21, 58, 19, 56]: the viewing angles, materials, and distances to the captured surfaces, the strength and color of lighting, etc. Therefore, to address the aforementioned issues, we designed the Multi-layer Point Cloud, a new volumetric representation more suitable for RGBD inputs. Contrary to conventional novel view sweep volumes, MPC volumes are generated from the input views instead of the novel view. Therefore, MPC volumes are always aligned with the direction of depth biases. As a result, MPC volumes (blue, Fig. 2(b)-(d)) ensure better coverage of true surfaces despite the steep viewing angles and can improve the reconstruction of slanted surfaces (*e.g.* object contours). More algorithmic details in Sec. 4.2.1.

## 4. Method

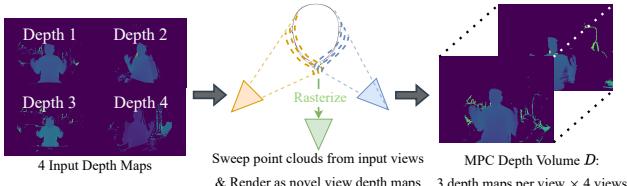
Our goal is to render stable and photorealistic free-viewpoint videos for personal telepresence using a few (*e.g.* 4) RGBD streams. In Sec. 4.1, we describe our capture system. In Sec. 4.2, we describe our rendering system.

### 4.1. Personal Desktop Capture System.

Our capture system is shown in Figure. 1. We enhance the typical video conferencing setup (*i.e.* one camera mounted near the screen) by placing 4 Microsoft Azure Kinect RGBD cameras near a 23-inch display. We place 3 of the cameras to the left, right, and above the display to



**Figure 2: MPC Volumes Vs. Conventional Novel View Sweep Volumes.** (a) Depth measurements from RGBD cameras are biased (*i.e.* offset from the true geometry), leading to misalignment between point clouds from different views. (b) Sweeping from the novel view (green volumes) may miss the true surfaces (identified by red circles). The error worsens as the surface gets more slanted in the novel view as shown in (c) and (d). In contrast, our Multi-layer Point Cloud (blue volume) better covers the true geometry, leading to better reconstruction.



**Figure 3: Construction of Multi-layer Point Cloud (MPC) Volume.** Each of the  $K = 4$  input depth maps is perturbed/swept with  $N = 3$  offset values. The  $N$  perturbed depth maps are then lifted into point clouds and rasterized into depth maps in the novel view, producing the MPC depth volume  $D \in \mathbb{R}^{K \times N \times H \times W}$  that stores  $K \times N$  depth candidates for each pixel in the novel view.

provide complete coverage of the user, and the 4th camera right above the monitor to provide details. The cameras are hardware-synchronized via audio cables. Each camera captures color images at  $2048 \times 1536$  and depth images, which are later merged into a single RGBD image via reprojection and rasterization. Additionally, a ring light is placed behind the middle camera in order to reduce shadows and improve the visibility of facial details. The Azure Kinect cameras are priced at \$399 each (totaling \$1596 for 4 cameras), similar to high-end VR headsets like Meta Quest Pro (\$1500) and cheaper than typical AR glasses (*e.g.* Magic Leap and HoloLens at over \$3500).

## 4.2. Rendering System.

Fig. 4 shows the overview of our rendering system. We first describe the general rendering pipeline that utilizes Multi-Layer Point Cloud (Sec. 4.2.1). Then, we describe our Temporal Renderer (Sec. 4.2.4) and Spatial Skip Connection (Sec. 4.2.5) that further improve the video quality.

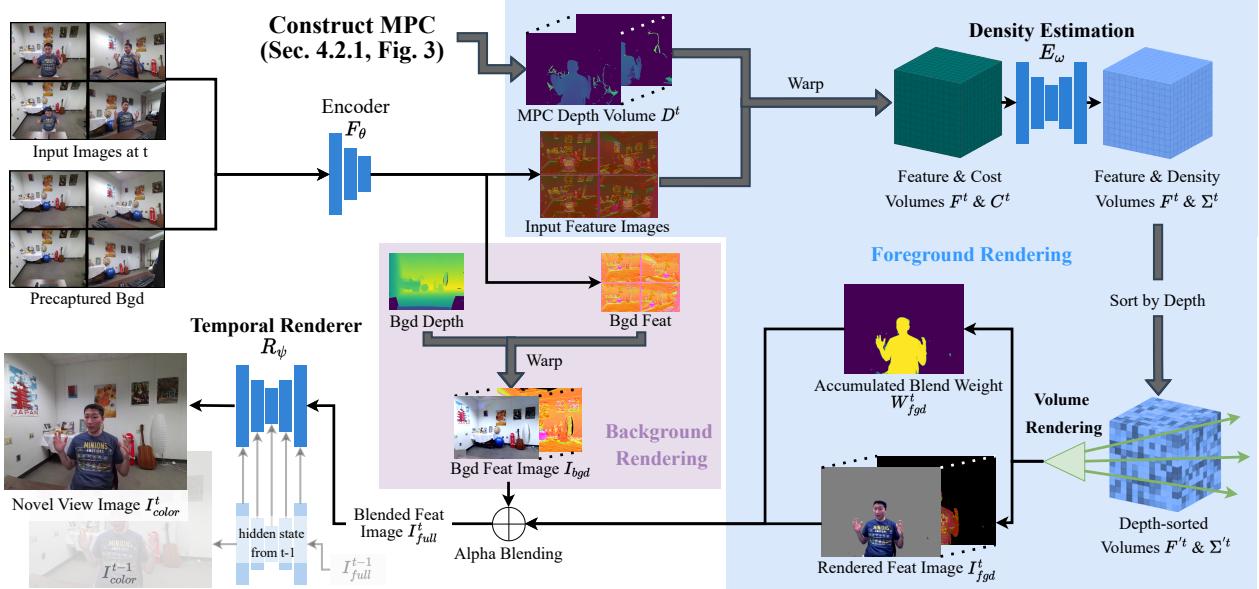
### 4.2.1 Constructing the MPC Volume

As discussed in Sec. 3, MPC improves reconstruction because it is robust to depth biases in RGBD cameras. In Fig. 3, we illustrate the construction of MPC volumes. For each

of the  $K = 4$  input views, we first generate  $N = 3$  copies of the input depth maps, each perturbed by a value of  $\Delta d$ ; we use  $\Delta d \in [-1cm, 0cm, +1cm]$  in our experiments. These perturbed input depth maps are then lifted into 3D space, resulting in  $N$  layers of point clouds for each of the  $K$  views. These point clouds are then rasterized into the novel view as depth maps, *i.e.* the MPC depth volume. Thus the MPC depth volume  $D \in \mathbb{R}^{K \times N \times H \times W}$  stores  $K \times N$  depth candidates for each pixel in the novel view. Notice that the MPC depth volume is ordered by cameras (*i.e.* the first  $N$  depth maps are from view 1, and the next  $N$  from view 2, etc.), and thus the  $K \times N$  depth candidates for a pixel are not sorted by depth. During our experiments, these depth volumes are preprocessed to avoid redundant computation and accelerate training and testing.

### 4.2.2 Density Estimation

The MPC depth volume  $D$  provides a set of depth candidates for each pixel in the novel view image. We then estimate density values for each candidate similar to prior works [4, 62, 30, 18, 6]. More specifically, we lift and project each candidate into all input views. The averaged feature of a candidate across views is stored in the feature volume  $F$ , and its feature variance across views is stored in the cost volume  $C$ . Good candidates (*i.e.* close to the true surface) are more likely to have low variances (*i.e.* high multi-view consistency) than bad candidates. To account for occlusions, only the unoccluded views (*i.e.* the candidate's projected depth is consistent with the input view's depth map) are used during averaging and variance calculation. Finally, a 3D ConvNet  $E_\psi$  estimates density values for each candidate based on the feature volume  $F$  and cost volume  $C$ , producing the density volume  $\Sigma$ . Notice that  $F$ ,  $C$  and  $\Sigma$  are not sorted by depth because the MPC depth volume is ordered by cameras (Sec. 4.2.1).



**Figure 4: System Overview:** (1) Foreground (Blue): Given the 4 input RGBD images at frame  $t$ , we first construct an MPC depth volume  $D^t$ .  $D^t$  is used to warp the input feature images to the novel view and construct the feature volume  $F^t$  and the cost volume  $C^t$ . Given  $F^t$  and  $C^t$ , A 3D ConvNet  $E_\psi$  then estimates a density volume  $\Sigma^t$ . To enable volume rendering, volumes  $\Sigma^t$  and  $F^t$  are then sorted by depth based on  $D^t$ , producing  $F'^t$  and  $\Sigma'^t$ . Volume rendering on  $F'^t$  thus produces the feature image  $I_{fgd}^t$  and the accumulated blend weight map  $W_{fgd}^t$ . (2) Foreground+Background(Purple): using the novel view background depth, we warp and average input feature images of the background into a background feature image  $I_{bgd}^t$ . We then alpha blend  $I_{bgd}^t$  and  $I_{fgd}^t$  into a full feature image  $I_{full}^t$  using  $W_{fgd}^t$ . Our temporal renderer then leverages the hidden state from  $t - 1$  to render the novel view image  $I_{color}^t$ .

#### 4.2.3 Rendering

To achieve high-quality results, we leverage volumetric rendering, a key component to recent advances in novel view synthesis [1]. Given the depth-sorted density volume  $\Sigma'^t$ , volume rendering blends all candidates along the ray  $r$  of a pixel  $p$  into one feature pixel in the image  $I_{fgd}^t$ :

$$T_i = \exp\left(-\sum_{j=1}^{K \times N-1} \sigma_j\right) \quad (1)$$

$$I_{fgd}^t(r) = \sum_{i=1}^{K \times N} T_i (1 - \exp(-\sigma_i)) f_i \quad (2)$$

$f_i$  is the feature stored in  $F'^t$  for the  $i$ th candidate of pixel  $p$ . The result is a blended feature image  $I_{feat}^t \in \mathbb{R}^{H \times W \times F}$  and an accumulated blend weight map  $W_{feat}^t \in \mathbb{R}^{H \times W}$  (please refer to [39] for more details). Given the warped and averaged background feature image (Sec. 4.2.6), we alpha blend the background and foreground feature images into a full feature image  $I_{full}^t$  using  $W_{fgd}^t$ .  $I_{full}^t$  and  $W_{fgd}^t$  are then fed to a UNet-based renderer to generate a color image  $I_{color}^t$  and a feature image  $I_{out}^t$  at time  $t$ .

#### 4.2.4 Stabilization via Temporal Renderer

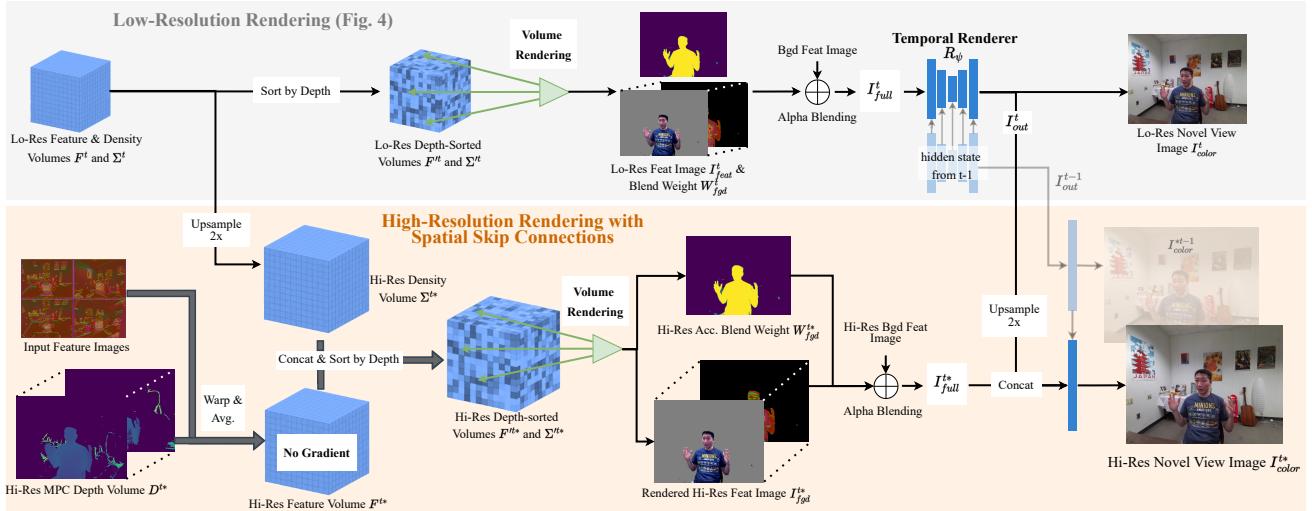
Our rendering system is capable of generating high-quality images without knowledge of prior frames. However, flickering artifacts between adjacent frames may occur due to

the noise in depth measurements from RGBD cameras. Some prior approaches perform stabilization by averaging adjacent frames, e.g., LookinGood [38] and Virtual-Cube [67]. However, this approach produces ghosting artifacts under fast motions, and flickering persists on steep surfaces. To mitigate these artifacts, we use a recurrent renderer  $R_\psi$  that exploits the temporal information between frames to suppress temporal inconsistencies.

First, our renderer  $R_\psi$  (a GRU-UNet) synthesizes the previous frame  $I_{color}^{t-1}$  at the viewpoint of interest. During the synthesis, the GRU captures useful temporal information and stores it as its hidden state. Then, the renderer uses the hidden state to condition the synthesis of  $I_{color}^t$  at the same viewpoint. This formulation enables our renderer to selectively utilize or discard information from the prior frame and synthesize more stable and accurate videos, even for fast motions.

#### 4.2.5 High-Resolution via Spatial Skip Connections

Our rendering system requires 3D convolution for volume rendering, and high-resolution rendering can easily exceed the GPU memory capacity. A straightforward solution to achieve high resolution is to perform convolution on an upsampled low resolution input. However, this technique relies on convolution layers to hallucinate high-resolution details, which can be unreliable. To overcome this challenge,



**Figure 5: Spatial Skip Connections(SSL):** Similar to skip connections in UNet[51], SSL directly propagates low-level details to a convolution layer by concatenating a high-resolution feature map  $I_{full}^{t*}$  with an upsampled feature map  $I_{out}^t$ . To calculate  $I_{full}^{t*}$ , we perform volume rendering using high-resolution feature and density volumes  $F^{t*}$  and  $\Sigma^{t*}$ . To avoid the memory heavy 3D convolution,  $\Sigma^{t*}$  is upsampled from  $\Sigma^t$  in order to approximate the actual density volume. Additionally,  $F^{t*}$  is constructed without gradient calculation.

we propose the Spatial Skip Connection (SSC) to enhance the details. Similar to skip connections in UNet[51], SSC directly propagates low-level details to a convolution layer by concatenating a high-resolution feature map  $I_{full}^{t*}$  with an upsampled feature map  $I_{out}^t$ .

To calculate the skip connection  $I_{full}^{t*}$ , we perform high-resolution rendering using the method in previous sections, but with two changes that save the memory: (1) Instead of performing costly 3D convolution to estimate a high-resolution density volume, we approximate it by upsampling  $\Sigma^t$  to  $\Sigma^{t*}$ . (2) we do not calculate gradient for the high-resolution feature volume  $F^{t*}$ . Volume rendering then generates the high-resolution blend weight map  $W_{fgd}^{t*}$  and foreground feature image  $I_{fgd}^{t*}$ .  $I_{fgd}^{t*}$  could thus be calculated by alpha blending  $I_{fgd}^{t*}$  and the high-resolution background feature image (Sec. 4.2.6).  $I_{fgd}^{t*}$  is then concatenated with  $I_{out}^t$  (Sec. 4.2.3) before being decoded into the final high-resolution novel view image  $I_{color}^{t*}$ .

#### 4.2.6 Background Pre-Capture

To reduce occlusion and improve stability, we pre-capture the background at the start of each capturing session without the users. The background depth maps are lifted into point clouds and rasterized in the novel view as depth maps. These rasterized background depth maps are then used to warp features from the input views into the novel view in order to generate a background feature image (Sec. 4.2.3 and Sec. 4.2.5).

## 5. Training

The input to our system are synchronized RGBD frames from the input cameras, and the output is a free-viewpoint video. Given that there are no existing dataset that uses multiple RGBD cameras to capture video conferencing sessions, we produced our own dataset: the Personal Telepresence Dataset.

### 5.1. Personal Telepresence Dataset

The dataset is captured using the hardware configuration described in Sec. 4.1. The dataset contains 35 training sequences of a single target user performing various actions in their office. There are also 10 test sequences of various users performing actions in the same office. Each sequence contains 4 input RGBD videos and 2 ground truth RGB-only videos, each lasting 20 to 30 seconds. The ground truth cameras are placed near the middle input camera. There are a total of 6 different camera setups (*i.e.* different relative positions and orientations between cameras), creating more diversity in the training data. We also introduce changes to the background decoration to ensure that models can account for variation in the user's environment across sessions. Given that the system is intended for personal use in a specific room, we include only moderate variation to the background. Before each capture session, we capture the background without the user in order to reduce occlusion. Before capturing with a new camera setup, we calibrate the camera poses with a checkerboard and the OpenCV calibration tool. We use the camera intrinsics provided by the manufacturer.

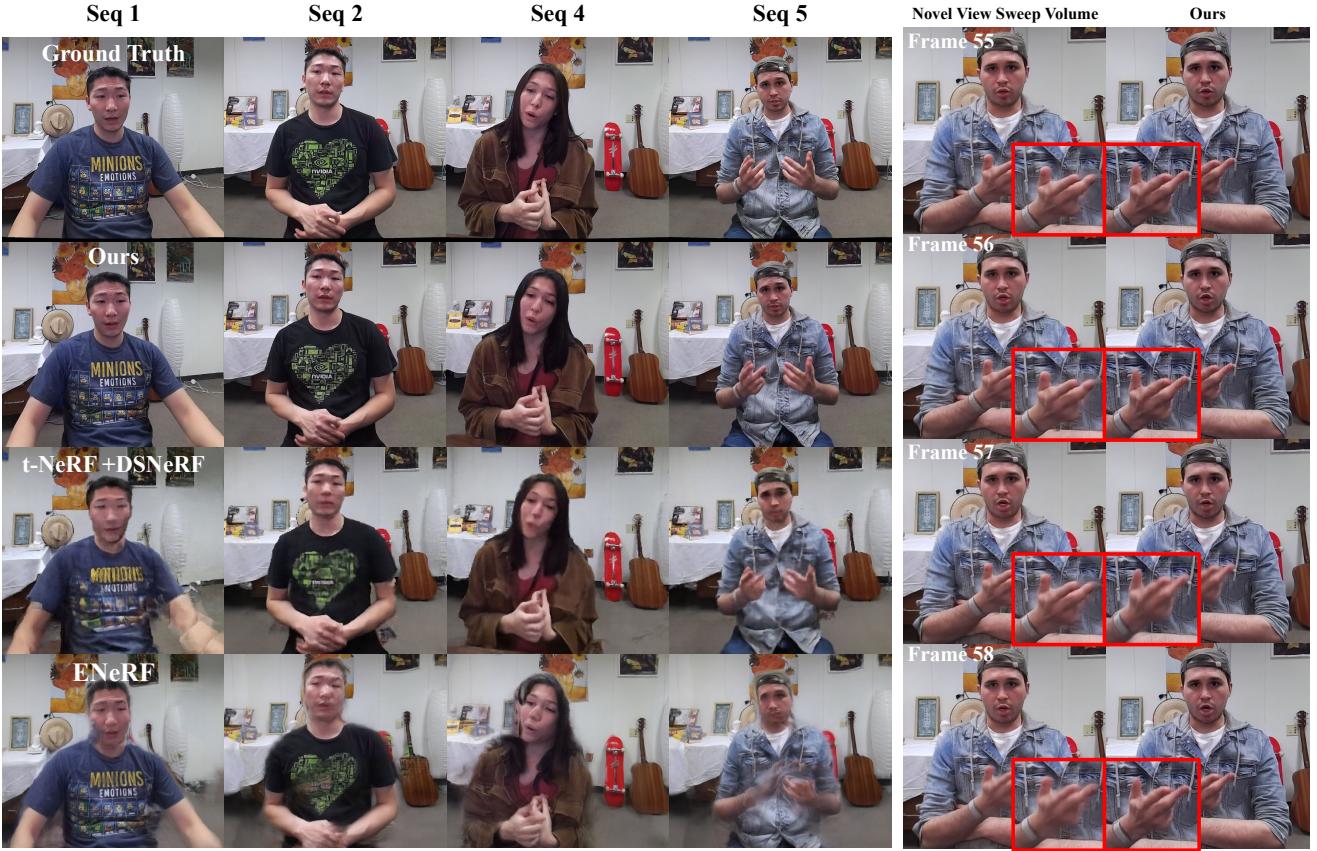


Figure 6: (a) Comparison with baseline methods. Our neural renderer recover sharp and fine-grained details in both the foreground and background. (b) Ablation. Compared to conventional depth sweeping from the novel view, our approach shows more accurate reconstruction of challenging structures like fingers.

## 5.2. Training Loss

There are three trainable modules in our rendering network: (1) the image feature encoder  $F_\theta(\cdot)$ . (2) the density estimation network  $E_\omega(\cdot)$ , and (3) the temporal renderer  $R_\psi(\cdot)$ . To save memory, we use a frozen EfficientNetV2s [55] pretrained on "ImageNet1K V1" [10] as our encoder  $F_\theta(\cdot)$ . There is no noticeable impact to performance. We train the rest of the network ( $E_\omega(\cdot)$  and  $R_\psi(\cdot)$ ) on our Personal Telepresence Dataset using the following losses:

$$L(\tilde{I}, I) = \|\tilde{I} - I\|_1 + \sum_l \lambda_l \|\phi_l(\tilde{I}) - \phi_l(I)\|_1 \quad (3)$$

$$L_{perc} = L(\tilde{I}^t, I_{color}^t) + L(\tilde{I}^{t*}, I_{color}^{t*}) \quad (4)$$

Eq. 3 is the perceptual loss function proposed by [5]. Eq. 4 calculates perceptual loss on both low and high-resolution images.  $\tilde{I}^t$  and  $I_{color}^t$  are the ground truth and rendered image at time  $t$ .  $*$  denotes a high-resolution image.

$$L_{fgd\ perc} = L(W \circ \tilde{I}^t, W \circ I_{color}^t) + L(W^* \circ \tilde{I}^{t*}, W^* \circ I_{color}^{t*}) \quad (5)$$

Eq. 5 calculates the perceptual loss on the foreground by masking out the background.  $W$  is the foreground mask

where a pixel is 1 if its accumulated weight in  $W_{fgd}^t$  is greater than  $> 0.5$ , and 0 otherwise. The final loss term is the summation of the general perceptual loss and the foreground perceptual loss:

$$L_{final} = L_{perc} + L_{fgd\ perc} \quad (6)$$

## 5.3. Variants

**Our Full Model** uses MPC, the temporal renderer and Spatial Skip Connection (SSC). It is trained in two stages in order to save memory: Stage (1): the model (*i.e.*  $E_\omega(\cdot)$  and  $R_\psi(\cdot)$ ) is first trained for 12 epochs to generate a novel view image for a single frame. Stage (2): With the density network  $E_\omega(\cdot)$  trained and frozen, the renderer  $R_\psi(\cdot)$  is re-initialized and re-trained for 6 epochs to produce 2 novel view images at the same viewpoint from 2 frames of input. In this way, the training memory consumption at Stage 2 is greatly reduced.

We also evaluate multiple variants: **Ours(MPC only)**, **Ours(MPC+Temporal)**, and **Novel View Sweep Volume**(no temp, SSC, MPC). These variants omit various components for ablation purposes, and they are trained with

Table 2: **Quantitative Comparisons.** Values ordered by JOD  $\uparrow$  / LPIPS  $\downarrow$  / PSNR (db)  $\uparrow$  / SSIM  $\uparrow$ .  $\uparrow$  = higher is better,  $\downarrow$  = lower is better.

Methods	General	Sequence 1	Sequence 2	Sequence 3	Sequence 4	Sequence 5
	-izable	JOD $\uparrow$ / LPIPS $\downarrow$ / PSNR (db) $\uparrow$ / SSIM $\uparrow$	JOD $\uparrow$ / LPIPS $\downarrow$ / PSNR (db) $\uparrow$ / SSIM $\uparrow$	JOD $\uparrow$ / LPIPS $\downarrow$ / PSNR (db) $\uparrow$ / SSIM $\uparrow$	JOD $\uparrow$ / LPIPS $\downarrow$ / PSNR (db) $\uparrow$ / SSIM $\uparrow$	JOD $\uparrow$ / LPIPS $\downarrow$ / PSNR (db) $\uparrow$ / SSIM $\uparrow$
t-NeRF	$\times$	3.82 / 0.252 / 19.15 / 0.65	4.70 / 0.112 / 23.11 / 0.80	4.47 / 0.159 / 21.49 / 0.75	3.72 / 0.153 / 18.98 / 0.76	4.67 / 0.128 / 22.60 / 0.79
t-NeRF+DSNeRF[11]	$\times$	4.60 / 0.173 / 22.24 / 0.75	4.64 / 0.122 / 23.31 / 0.81	4.50 / 0.146 / 22.35 / 0.77	4.79 / 0.117 / 23.85 / 0.80	4.75 / 0.124 / 23.60 / 0.80
Gao <i>et al.</i> [16]+DSNeRF	$\times$	3.87 / 0.264 / 20.47 / 0.70	4.27 / 0.194 / 21.22 / 0.75	4.10 / 0.242 / 20.09 / 0.69	4.19 / 0.200 / 21.56 / 0.74	4.02 / 0.283 / 22.24 / 0.74
ENeRF[30]	$\checkmark$	4.50 / 0.274 / 17.85 / 0.65	4.24 / 0.284 / 18.01 / 0.66	4.44 / 0.284 / 18.01 / 0.66	4.24 / 0.266 / 18.35 / 0.68	4.56 / 0.266 / 18.35 / 0.68
Ours Full Model	$\checkmark$	<b>7.53 / 0.041 / 27.92 / 0.92</b>	<b>7.63 / 0.037 / 27.04 / 0.94</b>	<b>7.22 / 0.061 / 26.54 / 0.92</b>	<b>7.24 / 0.050 / 26.55 / 0.92</b>	<b>7.59 / 0.042 / 26.47 / 0.92</b>

Table 3: **Ablation Studies** Values ordered by JOD  $\uparrow$  / LPIPS  $\downarrow$  / PSNR (db)  $\uparrow$  / SSIM  $\uparrow$ .  $\uparrow$  = higher is better,  $\downarrow$  = lower is better.

Methods	Averrage on Sequences 1-5	
	JOD $\uparrow$ / LPIPS $\downarrow$ / PSNR (db) $\uparrow$ / SSIM $\uparrow$	JOD $\uparrow$ / LPIPS $\downarrow$ / PSNR (db) $\uparrow$ / SSIM $\uparrow$
Ours Full (MPC+Temporal+SSC)	<b>7.473 / 0.045 / 26.917 / 0.925</b>	
Ours (MPC+Temporal)	7.422 / 0.047 / 26.924 / 0.923	
Ours (MPC only)	7.401 / 0.050 / <b>27.100</b> / 0.923	
Novel View Sweep Volume	7.349 / 0.053 / 26.946 / 0.919	

Table 4: **Description of Test Sequences** The term **"same"** refers to something that is part of the training data. The term **"new"** refers to something that did not appear in the training data and thus new to the model. All sequences except for **"Sequence 1"** is evaluated from new novel viewpoints. This enables a rough understanding of how the performance changes when the novel viewpoint is not covered by training data.

	User	Clothing	Input Views	Evaluation Views
Sequence 1	Same	New	Same	Same
Sequence 2	Same	New	New	New
Sequence 3	New	New	New	New
Sequence 4	New	New	New	New
Sequence 5	New	New	New	New

Stage 1 only. Note that the **Novel View Sweep** variant uses conventional cost volumes, *i.e.* depth sweep from the novel view camera. Similarly to VirtualCube, we first lift input depth maps from all views into point clouds, which are then rasterized into novel view depth maps. The average of these depth maps are then used as the center of the local depth sweep. For a fair comparison, we use the same number (*i.e.* 12) of sweeps with a step size of 1cm (*i.e.* between -6cm and +6cm of the average depth map).

Notice that all variants use the **same pre-captured background** as our full model and thus the **same background feature image** during alpha blending, resulting in very similar performances on the background.

**Implementation Details** We train the network on 4 NVidia RTX 3090 GPUs using the AdamW [33] optimizer with a learning rate of  $10^{-4}$ , and a batch size of 4. To reduce memory consumption during training, we render random crops of size  $384 \times 512$  for the low-resolution synthesis and  $768 \times 1024$  for high-resolution. During testing and evaluation, the low-resolution rendering is of size  $480 \times 640$  and the high-resolution rendering is of size  $960 \times 1280$ . We preprocess MPC depth volumes and store them locally. While there are many real-time depth-map/point-cloud renderers, we use Pytorch3D[48] to render the depth maps due to its simplicity.

## 6. Evaluation

Our model is trained on a single male target user in his office using 4 RGBD cameras as inputs and 2 RGB-only cameras as ground truth supervision. During testing and evaluation, performances on the 2 ground truth cameras are averaged to provide a more robust assessment.

**Evaluation Metrics** We use LPIPS [66], SSIM, and PSNR to measure per-frame image quality. We also use JOD [37] as an important metric to evaluate video stability.

**Test Sequences** We also consider the different factors that may affect the evaluation results: (1) whether or not the novel viewpoints are included in the training data, (2) robustness to new clothing, (3) robustness to new users. Therefore, we employ different settings for the 5 test sequences as shown in Table 4. We also include results on 5 additional test sequences in our supplementary materials. Since this is a personal system, we do not introduce significant differences to the environment during testing.

### 6.1. Baselines

Our goal is to generate free-viewpoint videos from a few RGBD cameras in a video conferencing setup. Given that few competitive baselines are specifically designed for this application, we adapt recent neural rendering methods to this application in order to study a wider spectrum of possible solutions to the problem. These baselines include: (1) **t-NeRF**: NeRF [39] with time as an additional input, (2) **Gao *et al.* [16]**: models the scene with neural radiance and scene-flow fields. We also enhance baselines (1) and (2) with DSNeRF [11]’s depth loss term, which is a KL divergence term used to encourage uni-modal density distribution near the input depth value. (3) **ENeRF** [30]: a real-time generalizable radiance field optimized with a plane-sweep. (4) **Novel View Sweep Volume** (Tab. 3): Similar to Microsoft VirtualCube [67] (code not available), this baseline generates cost volumes via depth sweeps from the novel view camera. This is a variant of our system without MPC, SSC, or temporal rendering.

### 6.2. Results and Comparison

**Quantitative Results.** As shown in Table 2, our system outperforms all NeRF-based approaches in test sequences 1-5 by a large margin; this is likely due to the sparsity of viewpoints. Moreover, our system maintains a high quality even when applied to “new” users not in the training data

(Sequences 3-4). Additionally, Sequences 1 and 2 achieve similar results despite different camera setups (*i.e.* input and ground truth camera poses), implying that our system performs well with new viewpoints.

**Qualitative Results.** Fig. 6 shows novel view renderings for various sequences. Our system is robust to sparse viewpoints, reconstructs complex hand gestures with high fidelity, and generalizes well to new users, whereas our baselines struggle with sparse viewpoints. While it is not possible to directly compare with VirtualCube [67], we show in supplementary videos that our system renders more stable videos and is more robust to large fast movements.

**Ablation Studies.** Table. 3 shows the performance of different variants of averaged over the 5 test sequences. Note that the *pre-captured background is shared across all models*. This implies very similar results on a large portion of the images, and that the quantitative improvements mostly arise from the foreground. Despite sharing the backgrounds, our proposed modules still show consistent improvements in accuracy (LPIPS, SSIM) and stability (JOD). In the supplementary videos, we provide clear visualizations of improvements in the foreground area. The most significant improvement comes from our MPC volume, and our temporal renderer and Spatial Skip Connection further improves the stability and accuracies.

## 7. Limitation and Future Work

(1) Although our rendering system is fast, it is not real-time yet. Currently, the construction of cost volumes consumes the majority of the rendering time. To achieve real-time rendering, we will experiment with smaller volumes, sparse volumes, sparse convolution, and lighter networks. (2) Our current system does not render free-viewpoint videos via immersive display technologies (*e.g.* autostereo displays). In the future, we will track the user’s pupil positions and leverage autostereo displays to render separate videos for each eye to create depth cues essential for immersive telepresence.

## 8. Conclusion

We presented a capture and rendering system designed for personal telepresence systems. Our capturing system only requires a few RGBD cameras and thus can be easily installed on typical desks unlike recent commercial systems. Our rendering system renders high-quality free-viewpoint videos, and it outperforms recent view synthesis methods in terms of both accuracy and stability. It can accurately reconstruct complex hand gestures, fast body movements, and rich environmental details without object templates, complex pre-processing, or costly optimization.

By proposing a personal system, we work to democratize immersive telepresence experiences. The lightweight nature of this setup facilitates more intimate and informal connections between remote individuals without requiring

the resources needed by commercial counterparts. As such, our system helps increase access to telepresence for the general public and encourages future work in this area.

## References

- [1] Alex Yu and Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks, 2021. 5
- [2] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. 2020. 3
- [3] Benjamin Attal, Eliot Laidlaw, Aaron Gokaslan, Changil Kim, Christian Richardt, James Tompkin, and Matthew O’Toole. Törf: Time-of-flight radiance fields for dynamic scene view synthesis. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [4] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. *arXiv preprint arXiv:2103.15595*, 2021. 2, 3, 4
- [5] Qifeng Chen and Vladlen Koltun. Photographic image synthesis with cascaded refinement networks. *CoRR*, abs/1707.09405, 2017. 7
- [6] Shuo Cheng, Zexiang Xu, Shilin Zhu, Zhiwen Li, Li Erran Li, Ravi Ramamoorthi, and Hao Su. Deep stereo using adaptive thin volume representation with uncertainty awareness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2524–2534, 2020. 3, 4
- [7] Julian Chibane, Aayush Bansal, Verica Lazova, and Gerard Pons-Moll. Stereo radiance fields (srf): Learning view synthesis from sparse views of novel scenes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2021. 2
- [8] R.T. Collins. A space-sweep approach to true multi-image matching. In *Proceedings CVPR IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 358–363, 1996. 3
- [9] Carolina Cruz-Neira, Daniel J. Sandin, and Thomas A. DeFanti. Surround-screen projection-based virtual reality: The design and implementation of the cave. In *Proceedings of the 20th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH ’93, page 135–142, New York, NY, USA, 1993. Association for Computing Machinery. 2
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 7
- [11] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. *arXiv preprint arXiv:2107.02791*, 2021. 8
- [12] John V Draper, David B Kaber, and John M Usher. Telepresence. *Human factors*, 40(3):354–375, 1998. 1
- [13] John Flynn, Michael Broxton, Paul E. Debevec, Matthew DuVall, Graham Fyffe, Ryan S. Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. *CoRR*, abs/1906.07316, 2019. 3

- [14] Luigi Gallo, Kyis Essmael, Ernesto Damiani, Giuseppe De Pietro, and Albert Dipanda. Comparative evaluation of methods for filtering kinect depth data. *Multimedia Tools and Applications*, 74, 05 2014. 3
- [15] David Gallup, Jan-Michael Frahm, Philippos Mordohai, Qingxiong Yang, and Marc Pollefeys. Real-time plane-sweeping stereo with multiple sweeping directions. 06 2007. 3
- [16] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. *CoRR*, abs/2105.06468, 2021. 2, 8
- [17] Simon J Gibbs, Constantin Arapis, and Christian J Breiteneder. Teleport-towards immersive copresence. *Multimedia Systems*, 7(3):214–221, 1999. 1
- [18] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. 2019. 3, 4
- [19] Ying He, Bin Liang, Yu Zou, Jin He, and Jun Yang. Depth errors analysis and correction for time-of-flight (tof) cameras. *Sensors*, 17:92, 01 2017. 3
- [20] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3
- [21] Jiyoung Jung, Joon-Young Lee, Yekeun Jeong, and In So Kweon. Time-of-flight sensor calibration for a color and depth camera pair. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(7):1501–1513, 2015. 3
- [22] Peter Kauff and Oliver Schreer. An immersive 3d video-conferencing system using shared virtual team user environments. pages 105–112, 09 2002. 2
- [23] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. *CoRR*, abs/2109.02369, 2021. 3
- [24] Claudia Kuster, Nicola Ranieri, Agustina Agustina, H. Zimmer, J.C. Bazin, C. Sun, Tiberiu Popa, and M. Gross. Towards next generation 3d teleconferencing systems. pages 1–4, 10 2012. 2
- [25] Claudia Kuster, Nicola Ranieri, Henning Zimmer, Jean-Charles Bazin, Chengzheng Sun, Tiberiu Popa, Markus Gross, et al. Towards next generation 3d teleconferencing systems. In *2012 3DTV-Conference: The True Vision-Capture, Transmission and Display of 3D Video (3DTV-CON)*, pages 1–4. IEEE, 2012. 1
- [26] Youngjoong Kwon, Dahun Kim, Duygu Ceylan, and Henry Fuchs. Neural human performer: Learning generalizable radiance fields for human performance rendering. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [27] Jason Lawrence, Dan B Goldman, Sreeth Achar, Gregory Major Blashevich, Joseph G. Deslodge, Tommy Fortes, Eric M. Gomez, Sascha Häberling, Hugues Hoppe, Andy Huibers, Claude Knaus, Brian Kuschak, Ricardo Martin-Brualla, Harris Nover, Andrew Ian Russell, Steven M. Seitz, and Kevin Tong. Project starline: A high-fidelity telepresence system. *ACM Transactions on Graphics (Proc. SIGGRAPH Asia)*, 40(6), 2021. 1, 2, 3
- [28] Tianye Li, Mira Slavcheva, Michael Zollhöfer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, and Zhaoyang Lv. Neural 3d video synthesis. *CoRR*, abs/2103.02597, 2021. 2
- [29] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [30] Haotong Lin, Sida Peng, Zhen Xu, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Efficient neural radiance fields for interactive free-viewpoint video. In *SIGGRAPH Asia Conference Proceedings*, 2022. 2, 3, 4, 8
- [31] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Trans. Graph.(ACM SIGGRAPH Asia)*, 2021. 2
- [32] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015. 3
- [33] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2017. 8
- [34] Andrew Maimone, Jonathan Bidwell, Kun Peng, and Henry Fuchs. Enhanced personal autostereoscopic telepresence system using commodity depth cameras. *Computers & Graphics*, 36(7):791–807, 2012. 1
- [35] Andrew Maimone and Henry Fuchs. Encumbrance-free telepresence system with real-time 3d capture and display using commodity depth cameras. In *2011 10th IEEE International Symposium on Mixed and Augmented Reality*, pages 137–146, 2011. 2
- [36] Andrew Maimone, Xubo Yang, Nate Dierk, Andrei State, Mingsong Dou, and Henry Fuchs. General-purpose telepresence with head-worn optical see-through displays and projector-based lighting. In *2013 IEEE Virtual Reality (VR)*, pages 23–26, 2013. 2, 3
- [37] Rafal K. Mantiuk, Gyorgy Denes, Alexandre Chapiro, Anton Kaplanyan, Gizem Rufo, Romain Bachy, Trisha Lian, and Anjul Patney. Fovvideovdp: A visible difference predictor for wide field-of-view video. *ACM Trans. Graph.*, 40(4), jul 2021. 8
- [38] Ricardo Martin-Brualla, Rohit Pandey, Shuoran Yang, Pavel Pidlypenskyi, Jonathan Taylor, Julien P. C. Valentin, Sameh Khamis, Philip L. Davidson, Anastasia Tkach, Peter Lincoln, Adarsh Kowdle, Christoph Rhemann, Dan B. Goldman, Cem Keskin, Steven M. Seitz, Shahram Izadi, and Sean Ryan Fanello. Lookingood: Enhancing performance capture with real-time neural re-rendering. *CoRR*, abs/1811.05029, 2018. 2, 3, 5
- [39] 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 *European conference on computer vision*, pages 405–421. Springer, 2020. 5, 8
- [40] Phong Nguyen, Nikolaos Sarafianos, Christoph Lassner, Janne Heikkila, and Tony Tung. Human view synthesis using a single sparse rgb-d input. 2021. 3

- [41] Sergio Orts-Escalano, Christoph Rhemann, Sean Fanello, Wayne Chang, Adarsh Kowdle, Yury Degtyarev, David Kim, Philip L. Davidson, Sameh Khamis, Mingsong Dou, Vladimir Tankovich, Charles Loop, Qin Cai, Philip A. Chou, Sarah Mennicken, Julien Valentin, Vivek Pradeep, Shenlong Wang, Sing Bing Kang, Pushmeet Kohli, Yuliya Lutchny, Cem Keskin, and Shahram Izadi. Holoporation: Virtual 3d teleportation in real-time. In *Proceedings of the 29th Annual Symposium on User Interface Software and Technology*, UIST ’16, page 741–754, New York, NY, USA, 2016. Association for Computing Machinery. 1, 3
- [42] Ahmed A A Osman, Timo Bolkart, and Michael J. Black. STAR: A sparse trained articulated human body regressor. In *European Conference on Computer Vision (ECCV)*, pages 598–613, 2020. 3
- [43] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021. 2
- [44] 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. *ACM Trans. Graph.*, 40(6), dec 2021. 2
- [45] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, 2021. 2
- [46] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *CVPR*, 2021. 2
- [47] Ramesh Raskar, Greg Welch, Matt Cutts, Adam Lake, Lev Stesin, and Henry Fuchs. The office of the future: A unified approach to image-based modeling and spatially immersive displays. In *Proceedings of the 25th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH ’98, page 179–188, New York, NY, USA, 1998. Association for Computing Machinery. 1, 2
- [48] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 8
- [49] Gernot Riegler and Vladlen Koltun. Free view synthesis. *CoRR*, abs/2008.05511, 2020. 3
- [50] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12216–12225, 2021. 3
- [51] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 2, 6
- [52] Darius Rückert, Linus Franke, and Marc Stamminger. Adop: Approximate differentiable one-pixel point rendering. *arXiv preprint arXiv:2110.06635*, 2021. 3
- [53] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 3
- [54] Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin. A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose. In *Advances in Neural Information Processing Systems*, 2021. 2
- [55] Mingxing Tan and Quoc Le. Efficientnetv2: Smaller models and faster training. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 10096–10106. PMLR, 18–24 Jul 2021. 7
- [56] Kenichiro Tanaka, Yasuhiro Mukaigawa, Takuya Funatomi, Hiroyuki Kubo, Yasuyuki Matsushita, and Yasushi Yagi. Material classification from time-of-flight distortions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PP:1–1, 09 2018. 3
- [57] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2021. 2
- [58] Michal Tölgessy, Martin Dekan, Lubos Chovanec, and Peter Hubinský. Evaluation of the azure kinect and its comparison to kinect v1 and kinect v2. *Sensors*, 21:413, 01 2021. 3
- [59] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, 2021. 2
- [60] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. SynSin: End-to-end view synthesis from a single image. In *CVPR*, 2020. 3
- [61] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9421–9431, 2021. 2
- [62] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo, 2018. 3, 4
- [63] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5336–5345, 2020. 2
- [64] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images, 2020. 2, 3
- [65] Cha Zhang, Qin Cai, Philip A. Chou, Zhengyou Zhang, and Ricardo Martin-Brualla. Viewport: A distributed, immersive teleconferencing system with infrared dot pattern. *IEEE MultiMedia*, 20(1):17–27, 2013. 2

- [66] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 8
- [67] Yizhong Zhang, Jiaolong Yang, Zhen Liu, Ruicheng Wang, Guojun Chen, Xin Tong, and Baining Guo. Virtualcube: An immersive 3d video communication system, 2021. 1, 2, 3, 5, 8, 9
- [68] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018. 3