# INV: Towards Streaming Incremental Neural Videos

Shengze Wang[1,2]    Alexey Supikov[2]    Joshua Ratcliff[2]    Henry Fuchs[1]    Ronald Azuma[2]

[1]UNC Chapel Hill                                    [2]Intel Labs

{shengzew,fuchs}@cs.unc.edu    {Alexei.Soupikov,joshua.j.ratcliff,ronald.t.azuma}@intel.com

## Abstract

*Recent works in spatiotemporal radiance fields can produce photorealistic free-viewpoint videos. However, they are inherently unsuitable for interactive streaming scenarios (e.g. video conferencing, telepresence) because have an inevitable lag even if the training is instantaneous. This is because these approaches consume videos and thus have to buffer chunks of frames (often seconds) before processing.*

*In this work, we take a step towards interactive streaming via a frame-by-frame approach naturally free of lag. Conventional wisdom believes that per-frame NeRFs are impractical due to prohibitive training costs and storage. We break this belief by introducing Incremental Neural Videos (INV), a per-frame NeRF that is efficiently trained and streamable. We designed INV based on two insights:*

*(1) Our main finding is that MLPs naturally partition themselves into Structure and Color Layers, which store structural and color/texture information respectively.*

*(2) We leverage this property to retain and improve upon knowledge from previous frames, thus amortizing training across frames and reducing redundant learning.*

*As a result, with negligible changes to NeRF, INV can achieve good qualities ($> 28.6db$) in 8min/frame. It can also outperform prior SOTA in $19\%$ less training time. Additionally, our Temporal Weight Compression reduces the per-frame size to 0.3MB/frame (6.6% of NeRF). More importantly, INV is free from buffer lag and is naturally fit for streaming. While this work does not achieve real-time training, it shows that incremental approaches like INV present new possibilities in interactive 3D streaming. Moreover, our discovery of natural information partition leads to a better understanding and manipulation of MLPs. Code and dataset will be released soon.*

## 1. Introduction

Recent advances in photorealistic 3D video generation hint at an exciting future where moments can be captured and experienced in 3D. One application of peculiar potential is interactive 3D streaming. It is the basis of telepresence
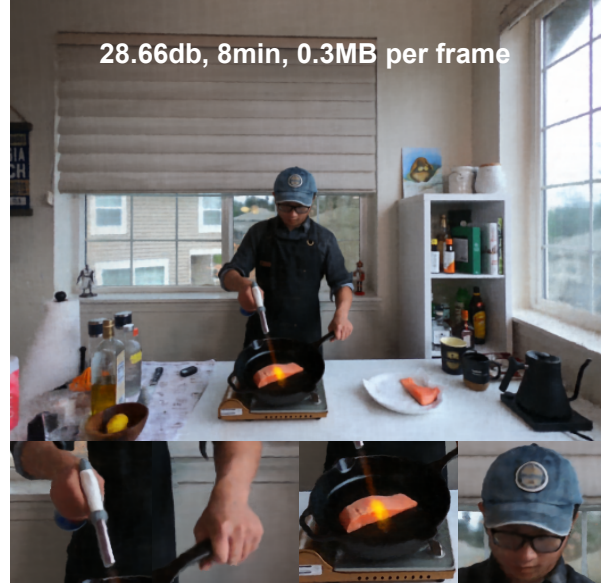


Figure 1. An example rendered novel view rendered by INV on ***Plenoptic Video Datasets [16]***. INV uses the original NeRF as the rendering backbone and averages over 28.7db PSNR with only 0.3MB of storage and 8 minutes of training per frame.

systems, which can democratize high-quality education via immersive remote classrooms, replace business traveling, and allow for face-to-face reunions between distant friends and families. In this work, we take a step towards this exciting future by introducing an efficient representation suitable for interactive 3D video streaming.

Spatiotemporal NeRFs [3, 11, 16, 17]) are a popular approach to synthesizing 3D videos due to their impressive photorealism. However, they suffer from an inherent lag because they consume videos and thus have to wait for chunks of frames (often seconds) before processing. Due to this buffer lag, these approaches are not suitable for interactive streaming scenarios like telepresence and immersive cloud gaming. In this work, we introduce Incremental Neural Videos (INV), a frame-by-frame approach that is naturally free of lag. While it is a common belief that per-frame NeRFs are impractical due to their prohibitive training costs

and storage size, INV breaks this belief by showing the possibility of training each frame in minutes and maintaining a streamable size. INV achieves this with negligible changes to NeRF, and it is a general framework that can be combined with more optimized approaches.

Moreover, prior approaches mostly show results on short videos lasting several seconds. This is because they model the entire 3D video as a single spatiotemporal object, and thus are limited by small model capacities. While one can divide longer videos into chunks (as in Neural 3D Videos [16]), it could months of GPU hours to generate a minute-long 3D video. The high training costs thus prevent the general public from using them.

We argue that the cause of such high training costs, in addition to the commonly discussed sampling bottleneck, can also be attributed to the learning of duplicated information. Prior works [11,16,17,26] learn the mapping from input 4D points $[x, y, z, t]$ (or equivalent) to their radiance. This input formulation essentially views 3D points from different frames as different 4D points. As a result, even if the scene is static, the same 3D points would need to be re-learned when the time variable changes. In our work, we retain and improve upon previously learned information to avoid training each frame from scratch and to continually improve over time. Moreover, we discovered that MLPs naturally store structural information in earlier Structure Layers and color information in later Color Layers. Since the color/texture space is often consistent in many scenarios (*e.g.* teleconferencing and live streaming), we significantly reduce storage by maintaining a constant color/texture space and learning only structural information.

Based on these observations, we propose Incremental Neural Videos (INV), a neural representation that can be efficiently learned and stored frame-by-frame. Our contributions can be summarized as follows:

- Natural Partition of Structure and Color Layers: we found that MLPs naturally store structural in earlier layers and color/texture information in later layers. We perform numerous experiments on both 2D and 3D videos to showcase this phenomenon. This discovery leads to a clearer understanding and more effective manipulation of MLPs.

- We leverage this phenomenon to design Incremental Neural Videos (INV), which is naturally free of lag and suitable for interactive streaming. INV consists of two types of sub-modules: (1) a shared color module that is shared across frames and encodes the color/texture of the scene, and (2) per-frame structure modules that encode the changing structures of the dynamic scene and are stored frame-by-frame.

- We propose Structure Transfer, an incremental training scheme that significantly accelerates training. With Structure Transfer, INV outperforms the state-of-the-art on per-frame quality metrics with less training. Moreover, INV can achieve good qualities in minutes.

- We propose Temporal Weight Compression to further compress the already concise representation. A compressed INV frame is only 0.3MB, 6.6% of NeRF.

## 2. Related Works

### 2.1. Static Novel View Synthesis

**Layered Representations.** Existing methods differ in how they represent 3D scenes. Multiplane Images (MPI) [4,6,39,41,44,54] are variable-viewpoint images that represent a scene as a set of fronto-parallel images, layered-meshes, or multilayer-spheres. However, these methods usually exhibit stack-of-card artifacts when viewed at a close distance or from steep angles.

**Neural Radiance Fields.** NeRF [26] introduced a simple but powerful representation based on Multi-Layer Perceptrons (MLPs) and differentiable volumetric rendering. It achieved unprecedented photorealism while allowing free viewpoints rendering. The high-level idea is to use a neural network to model the radiance of 3D points in the target scene. Notably, positional encodings [38,42] play a vital role in improving the quality of the reconstructed radiance field. Many works (*e.g.* [5,8,45,52]) have been inspired by NeRF. However, static NeRF faced two main challenges:

(1) Slow training and rendering. Multiple approaches (*e.g.* [1,7,10,12,13,18,27,49]) have been proposed to improve training and rendering speed, but long training time is still the main bottleneck of NeRF-based models.

(2) Generalization to new scenes. Many works like [8,9,45,50] have tried to achieve generalization by adding ConvNets or Transformers to NeRF, but it remains a challenge to learn generalizable neural radiance fields.

**Explicit 3D Representations.** One of the fundamental difficulties in view synthesis is to model the scene in 3D based on 2D images. To alleviate the burden from view synthesis algorithms, some works [2,14,34–36] rely on preprocessed raw 3D geometry (*e.g.* point clouds and meshes), obtained from multi-view stereo software such as COLMAP [37]. Such approaches show fewer fog-like artifacts which are common in NeRF-based methods. Methods like [34,35] also demonstrate good generalization to new scenes. Point cloud-based neural rendering approaches [14,36,46] show impressive sharpness and details in large scenes with thin structures. However, the aforementioned approaches rely on good input 3D models to render high-quality images.

### 2.2. Dynamic Novel View Synthesis

**Human-Specific Approaches.** Some recent works [15, 21, 32, 33, 40] focus on animating clothed humans only.

They often use colored videos as inputs and leverage human body templates (*e.g.* SMPL [23] and STAR [29]) and deep textures. On the other hand, LookinGood [25] uses multiple RGBD cameras to reconstruct human mesh in real-time. They then use a neural network to render colored mesh into high-quality novel views. HVSNet [28] uses monocular RGBD videos and renders from feature point clouds. Our work focuses on general scenes with complex human-object interactions. Therefore, approaches limited to human subjects would not apply to our target scenario.

**Dynamic Novel View Synthesis of General Scenes** Yoon et al. [48] utilizes multi-view stereo and monocular depth estimators to estimate spatially and temporally coherent depth maps and motion fields. Their approach generates 3D videos without per-video optimization or prior knowledge of the scene. Many NeRF-based approaches [11, 17, 43, 47] encode dynamic scenes as spatio-temporal radiance fields. Many learn a radiance field and a motion field for each frame. The motion field is then used to establish 3D correspondences between frames, thus allowing for temporal consistency regularizations. Optical flow and monocular depth estimators are often used to guide the optimization to a good local optimum. TöRF [3] uses Time-of-Flight sensors to achieve better modeling of both static and dynamic scenes. Approaches like Nerfies [30], Hyper-NeRF [31], and Neural 3D Videos (N3V) [16] use latent codes to help model dynamic contents. N3V also shows that smart sampling strategies can accelerate training and improve rendering qualities. However, all of these approaches require hours of training per frame, leading to weeks of training for a short video lasting several seconds. Our work shows how to significantly accelerate training and reduce storage via a concise incremental representation that naturally extends to streaming applications.

## 3. Method

We discuss our work in 3 parts: (1) the Incremental Transfer (I.T.) training scheme that drastically accelerates the training. I.T. also leads to (2) our discovery of Structure and Color Layers in MLPs, and (3) the Incremental Neural Videos (INV) representation that leverages this discovery to drastically reduce storage size and increase stability.

### 3.1. Transferring Information Across Frames

Prior works [11, 16, 17, 43, 47] can render impressively realistic 3D videos after hours of training per frame. The cause of such high training costs, aside from the well-explored sampling bottleneck, can also be attributed to the redundancy in training. Many works map 4D points $[x, y, z, t]$ (or equivalent) to $RGB\sigma$ values. This input formulation essentially views 3D points from different frames as different 4D points. Therefore, even if the scene is static, the same 3D points would need to be re-learned when the

time variable changes. While loss terms enforce consistencies between mappings for different frames, duplicated mappings still consume model capacity and training resources. A similar argument also applies to dynamic points.

We notice that, when the frame rate is infinitely high, adjacent frames should be almost identical. Therefore, given the MLP representation $\Phi_t$ for a frame at time $t$, it should take almost no effort to learn $\Phi_{t'}$ for the next frame at $t'$. Thus, the difference $\Delta\Phi = \Phi_{t'} - \Phi_t$ should decrease as time difference $\Delta t = t' - t$ decreases.

$$\lim_{\Delta t \to 0} \Delta\Phi = \mathbf{0}$$

Therefore, a model should reuse most of the previous information given high frame rates. Following this intuition, we propose Incremental Transfer (I.T.) to explicitly reuse information from frame to frame. We train one model per frame, with the earlier frame as the initialization for the later frame. As shown in Fig.3 and 5, I.T. leads to significantly less training by directly improving upon previous frames. The model thus gradually converges to a good quality similar to NeRF with only a fraction of the training. However, such a naive training scheme leads to significant temporal artifacts in rendered 3D videos, like flickering. In later sections, we show how INV reduces such artifacts.

Moreover, in Table. 1, we show that incrementally trained InstantNGP ("Ours NGP") acheives good results with merely 6 sec/frame of training.

### 3.2. Discovery of Structure and Color Layers

Another benefit of Incremental Transfer is that we can now analyze the effect of different $\Delta\Phi$'s. This leads to our core discovery: When mapping from image/volume coordinates to colors (and/or densities), an MLP naturally partitions itself into Structure Layers and Color Layers. Structure Layers are early layers of the model that store structural information. Color Layers are later layers that store color/style information (*i.e.* object color, shadows, etc.). Moreover, the mapping spaces learned by each layer are similar across nearby frames. We perform several experiments to show this interesting behavior.

#### 3.2.1 2D Structure Swap

We first examine how changes in Structure Layers can affect the structural contents of 2D images. First, we perform Incremental Transfer to learn MLPs **(with positional encodings)** $\Phi_A$ and $\Phi_B$ for frames A and B. Second, we perform Structure Swap, where we replace the 1st layer of $\Phi_A$ with that of $\Phi_B$. As shown in Fig. 2 (Left), **without any further training**, the whole scene moves to the right (*e.g.* more letters on the car plate are covered) and the kid's legs move closer. However, replacing later layers does not induce such structural changes. This means that 2D MLPs

## Structure Layers & Color Layers in 2D

### Structural Info In Early Structure Layers



**Total Layers:** 5
**P.E.:** Regular (NeRF)
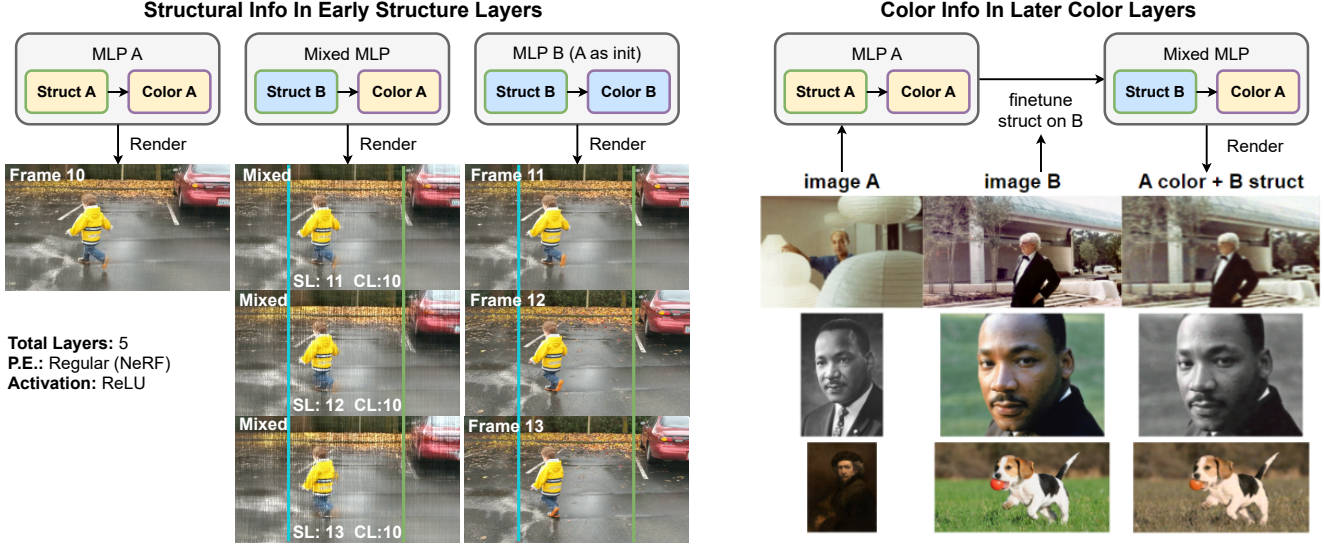**Activation:** ReLU

### Color Info In Later Color Layers



Figure 2. We found that MLPs naturally store structural and color information separately in **earlier** Structure Layers and **later** Color Layers. **Structure Swap (Left):** When an MLP is trained **incrementally** from frame A to B, color information remains similar among nearby frames but structural information varies. As a result, replacing A's Structure Layers (Struct A) with B's Structure Layers (Struct B) would induce meaningful structural changes/movements **without training**. Empirically, we observe 1 Structure Layer for 2D MLPs. **Color Scheme transfer (Right):** One can often transfer color schemes between images by (1) training MLP A on image A, and then (2) finetuning A's SLs (with CLs frozen) on image B. During the early stages of finetuning, the color scheme of A is often preserved.
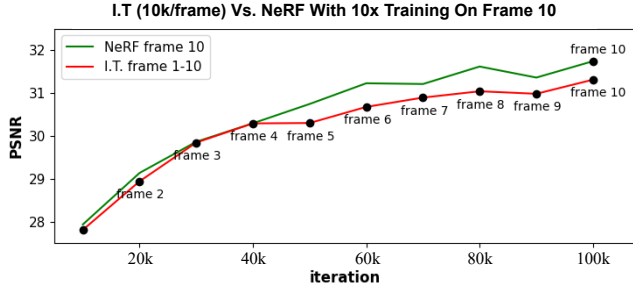


Figure 3. **NeRF 100k vs Incremental Transfer (I.T.) 10k:** On the "cut roasted steak" sequence in *Plenoptic Video Datasets* [16], I.T. is trained for 10k iterations (8 min) for each of the 10 frames, thus accumulating the same 100k iterations as NeRF (from scratch). On frame 10, I.T. achieves 31.31db PSNR, and NeRF achieves 31.75db. I.T. continually improves upon previously learned information, and thus quickly generates good results with a fraction of the training.

store structure information in earlier Structure Layers (the $1st$ layer in this case), and changes in Structure Layers induces structural changes in the rendered image. Moreover, since nearby Structure Layers can be swapped without further fine-tuning, the mappings learned by each layer are similar for nearby frames. This observation is important to support the idea of sharing the later layers.

### 3.2.2 2D Color Scheme Transfer

We have shown that 2D MLPs **(with positional encodings)** store structural information in early Structure Layers. But what do later layers store? We explore the answer via the Color Scheme Transfer experiment. As shown in Fig. 2 (Right), we first train MLP $\Phi_A$ to reconstruct image A. Then, we finetune the Structure Layer (1st layer in this case) on image B while freezing later Layers. During the early stages of finetuning (*e.g.* first 400 iters), the color scheme of A is often preserved while the structure quickly changes to B. Although the color scheme would eventually also change to B after longer training (*e.g.* 1k iters), this phenomenon shows that color information is stored in the later Color Layers of a 2D MLP.

### 3.2.3 3D Structure Swap and Color Scheme Transfer

We perform 3D analysis using our *Immersive Telepresence Dataset* (Fig. 4 (Left)), captured with 3 camera bars, each containing 5 compact cameras. Our images are much closer to the subject, enabling much easier visualization and analysis. We observed the following behaviors:

(1) NeRFs often contain more than one Structure Layer (ranging from 3 to 7 layers depending on the content).

(2) Replacing incrementally more Structural Layers induces continual and meaningful changes/motions in 3D. For

## Structure Layers & Color Layers in 3D

### Structural Info In Early Structure Layers
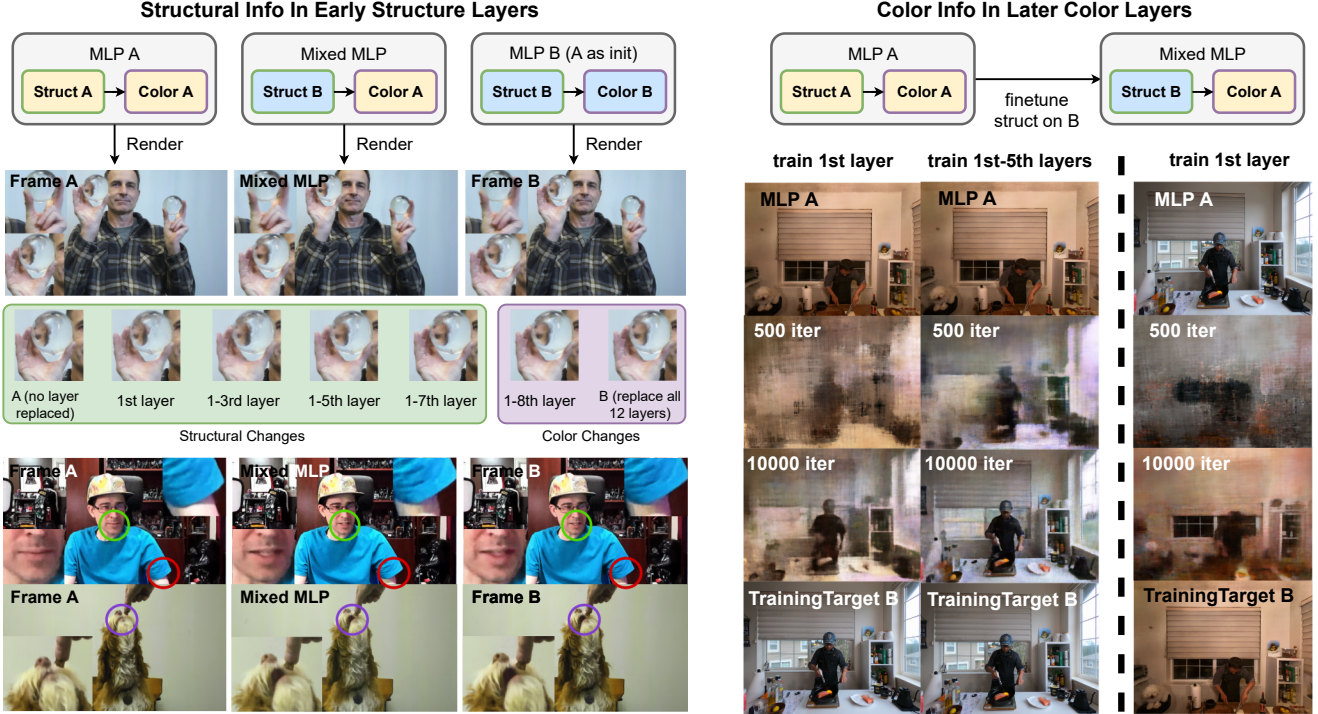


### Color Info In Later Color Layers



Figure 4. Natural information partitions also exist in NeRF. Different from 2D cases, there are more Structure Layers in the 3D case.

example in Fig. 4 (Left), when more Structure Layers from B are swapped into A, the person's head rotates more and a smile gradually appears. This means that the mappings learned by each layer stay consistent across time, allowing for layer replacements without finetuning.

(3) Replacing Colors Layers results in color changes (*e.g.* tone, shadows, *etc*.) but no obvious structural changes.

(4) 3D Color Scheme Transfer experiments show the same findings as 2D experiments.

This phenomenon is intuitively true since an MLP is a chain of non-linear mappings. The first layer maps coordinate/structural inputs to a latent space, which is then gradually transformed into the final color space through the chain of mappings. Earlier layers would thus retain more properties related to the input structural space.

3D Color Scheme Transfer experiments follow the same procedure described in Sec. 3.2.2. In Fig. 4 (Right), we show that the novel view renderings demonstrate expected behaviors. In other words, during the early stages of training (*e.g.* first 500 iters), the color scheme of the "original" is preserved while the structure quickly changes to the "target". Then, much longer training is needed to correct the color scheme since later layers are frozen (notice that with more trainable layers, this process becomes faster and easier). This consistency between 2D and 3D MLPs further supports our finding that MLPs **(with pos. enc.)** naturally

store structure and color information separately.

These findings also hint at a different view of Instant-NGP [27]. **NGP's hash grid can be seen as the output of Structure Layers stored in 3D, and the MLP can be seen as the Color Layers**. In this way, NGP achieves the same "position → latent space" mapping as the Structure Layers without going through any linear layers, thus directly accessing the Color Layers and thus accelerating the training.

### 3.3. Efficient Incremental Neural Videos

In Sec. 3.1, we showed that Incremental Transfer (I.T.) significantly reduces the training cost for the original NeRF while achieving good per-frame image quality. However, I.T. has two main disadvantages: (1) the storage cost is extremely high. A 5-minute 30FPS 3D video would require 40.1GB of storage. (2) There is significant flickering despite good per-frame image quality. To address these two issues, we propose the Incremental Neural Videos (INV) Representation that drastically reduces storage cost while providing high-quality visual results with notably improved temporal stability over I.T.

#### 3.3.1 The INV Representation

Our representation consists of 2 types of modules: (1) Per-Frame Structure Layers (SLs) which are stored frame-by-
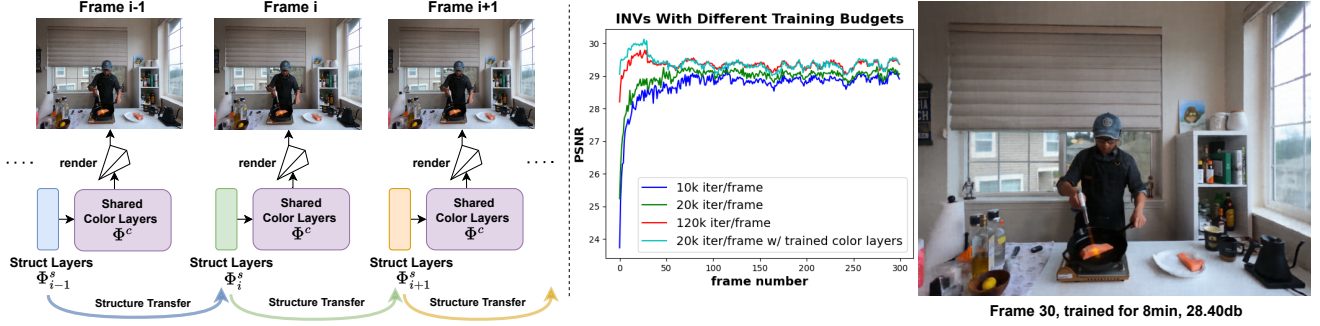
Figure 5. **Overview. Left:** Our representation consists of 2 types of modules: (1) Per-Frame Structure Layers (SL) $\Phi_i^s$ ($i \in 1...N$), which encode the structure for a specific frame $i$, and (2) Shared Color Layers (SCL) $\Phi^c$, which are shared by all frames to describe the color style of the scene. A video of N frames results in N Per-Frame SLs $\Phi_i^s$ ($i \in 1...N$) and 1 SCL $\Phi^c$. In this way, we maintain a constant color/style space while allowing the scene structure to change from frame to frame. **Right:** With only minutes of training per frame (8min in this example), INV quickly and continually improves in the first few seconds and then stabilizes to good quality.

frame, and (2) Shared Color Layers (SCL) which are shared by all frames. Per-Frame SLs encode the structure for a specific frame. They map the coordinate inputs (after positional encoding) to latent vectors $v$. SCL then maps these latent vectors $v$ into $RGB\sigma$. As a result, a video of N frames would be converted into N Per-Frame SLs $\Phi_i^s$ ($i \in 1...N$) and 1 SCL $\Phi^c$. In this way, we have a constant color/style space for the scene while allowing the scene structure to change from frame to frame. Moreover, storing only Structure Layers brings storage savings. For the original NeRF of 12 layers (including $RGB$ and $\sigma$ heads, no skip connection), storing 3 Structure Layers would reduce the weight size to 1.12MB (24.6% of the complete model of 4.43MB).

### 3.3.2 Training INVs via Structure Transfer

There are two stages to training INVs:

(1) **Warm-Up**: Incremental Transfer (all layers trained) from frame to frame until frame $i$. Color Layers at frame $i$ are then stored and shared as the Shared Color Layers $\Phi^c$.

(2) **Structure Transfer**: Incremental Transfer only on Structure Layers, with Shared Color Layers frozen. Structural information from previous frames is thus explicitly reused in later frames. Structure Layers' outputs are then converted into colors and densities via Shared Color Layers, which provide a constant color/style space. Only Structure Layers are stored for each frame (about 1.12MB/frame). Due to the freezing, videos are notably more stable as seen in the supplementary videos.

(Optional) **SplitINV: Video Stabilization with Background NeRF**: When INVs receive short training (*e.g.* 8min/frame), videos can be flickery. To alleviate this, one can use a separate frozen NeRF to model the static backgrounds and focus computation on dynamic contents. Implementation details are included in the supplementary.

To render an INV frame $i$, one could first recover the full model $\Phi$ by concatenating the Per-Frame Structure Layers $\Phi_i^s$ with the Shared Color Layers $\Phi^c$. The recovered model $\Phi_i$ is then used to render the 3D frame $i$.

### 3.3.3 Temporal Weight Compression (TWC)

To further reduce the size, we propose Temporal Weight Compression. Floating-point data compression algorithms (*e.g.* [19, 20]) can provide significant compression rate while maintaining high fidelity for structured data (*e.g.* 2D images, 3D geometry, *etc.*), but they struggle to compress near-random data that lack clear structure, such as the MLPs weights. However, there could be temporal structures in how the weights change through time. Therefore, we construct temporal weight matrices by concatenating weights from consecutive frames. We then perform FPZIP [20] at 16bits precision on the temporal weight matrix. As shown in Table. 1 "INV+TWC", this simple approach maintains video quality while compressing the Structure Layers down to 0.3MB/frame. Therefore, once the Shared Color Layers (3.29MB) are transmitted, it only costs 0.3MB to transmit any new INV frame, *i.e.* 72 MegaBits/second at 30FPS.

## 4. Implementation Details

Our baseline NeRF is a faithful PyTorch re-implementation of NeRF that reproduces the original results. We use the same NeRF for INV, except that we disable skip connections to the middle of the MLP. Although skip connections provide some quality improvement, they activate more Structure Layers and thus increase the storage size. All models are trained on a single NVIDIA RTX3090 GPU with the same settings as the original NeRF (learning rate $5 \times 10^{-4}$, 1024 ray samples, 64 depth samples for coarse NeRF and 128 samples for fine NeRF). In our tables, 8 minutes of INV training amounts to 10k iterations, 15

Table 1. **Comparisons and Ablations On Plenoptic Video Dataset [16]** With 50 minutes/frame less ($-19\%$) training than the SOTA Neural 3D Videos, INV achieves better per-frame qualities (PSNR & LPIPS). INV also achieves good quality in only 8min/frame. **"INV+TWC":** Temporal Weight Compression maintains video quality while compressing INV to 0.3MB/frame (6.6% of NeRF). **"INV⋆":** with well-trained Shared Color Layers, it takes just 15min/frame to achieve similar quality to 90min/frame. **Size↓** and **Training↓** are per-frame. **Warm-Up** is the number of warm-up frames, which are of lower quality but still included in the evaluation.

| Method | Variation | Warm-Up | # of SLs | PSNR↑ | LPIPS↓ | JOD↑ | Size↓ | Training↓ |
|---|---|---|---|---|---|---|---|---|
| MVS | - | - | - | 19.1213 | 0.2599 | - | - | - |
| NV [22] | - | - | - | 22.7975 | 0.2951 | 6.50 | 2.58MB | - |
| LLFF | - | - | - | 23.2388 | 0.2346 | 6.48 | - | - |
| NeRF-T | - | - | - | 28.4487 | 0.1000 | 7.73 | - | - |
| DyNeRF [16] | - | - | - | **29.5808** | **0.0832** | **8.07** | **0.09MB** | 260min |
| NeRF [26] | - | - | - | 24.6232 | 0.3588 | 4.73 | 4.56MB | **8min** |
| InstantNGP [27] | - | - | - | 24.4650 | 0.2977 | 4.87 | - | 6sec |
| Ours NGP | Incremental | 30 | - | 26.5933 | 0.0944 | 6.29 | - | 6sec |
| INV | - | 120 | first 3 layers | 28.7220 | 0.1215 | 6.61 | 1.12MB | **8min** |
| INV | - | 60 | first 3 layers | 28.9544 | 0.1111 | 6.79 | 1.12MB | 15min |
| INV | - | 30 | first 3 layers | 29.3491 | 0.0964 | 7.08 | 1.12MB | 90min |
| INV | - | 30 | first 3 layers | **29.8471** | **0.0818** | **7.42** | 1.12MB | 210min |
| INV | - | 60 | first 6 layers | 28.9236 | 0.1194 | 6.11 | 2.63MB | **8min** |
| INV+TWC | TWC | 120 | first 3 layers | 28.8698 | 0.1226 | 6.61 | **0.29MB** | **8min** |
| SplitINV | Bgd NeRF | 30 | first 3 layers | 28.7816 | 0.1157 | 6.92 | 1.12MB | 10min |
| INV+TWC | TWC | 30 | first 3 layers | 29.3364 | 0.1062 | **7.11** | 0.31MB | 90min |
| INV ⋆ | Trained SCL | 30 | first 3 layers | **29.3791** | **0.0898** | 6.85 | 1.12MB | 15min |

Table 2. **INV Evaluation on other sequences**: We evaluate INV on 3 other sequences released by Neural 3D Videos(DyNeRF) [16]. We are not able to compare to DyNeRF on these sequences because DyNeRF did not report the performances on these sequences and the code is not released.

| Seq | PSNR↑ | LPIPS↓ | JOD↑ | Training↓ |
|---|---|---|---|---|
| cut roasted beef | 32.0214 | 0.0644 | 7.17 | 20k (15min) |
| sear steak | 31.7662 | 0.0518 | 7.51 | 120k (90min) |
| coffee martini | 28.1384 | 0.0963 | 6.90 | 120k (90min) |

min is 20k, 90 min is 120k, and 210 min is 280k.

Temporal Weight Compression batches all INV frames (after warmup) into a single temporal weight matrix, then we use the *fpzip* python library to perform floating-point compression at 16 bits resolution.

# 5. Experiments

We evaluate the quantitative and qualitative performance of INV against State-Of-The-Art approaches.

## 5.1. Datasets

The ***Plenoptic Video Datasets (PVD) [16]*** is released by Neural 3D Videos (DyNeRF) as their benchmark dataset. It is captured with 21 GoPro cameras at $2028 \times 2704$ and 30 FPS. It contains challenging volumetric effects (*e.g.* flames), view-dependent effects (e.g. specularity from silverware, transparency from bottles), complex actions (*e.g.*

cooking), changing topologies (*e.g.* cutting beef, pouring liquid), etc. Same as DyNeRF, INV is evaluated at $1352 \times 1014$ on the same and only sequence that DyNeRF reported quantitative results on, *i.e.* the first 10 seconds of the "flame salmon 1" sequence. We report quantitative results on other sequences in the supplementary.

***Immersive Telepresence Dataset***. As mentioned in Sec. 3.2.3, we also perform analysis of Structure and Color Layers on our custom dataset. The dataset is captured with 3 camera bars each containing 5 compact cameras, and we follow *PVD [16]* to use 14 views for training and 1 middle view for evaluation. Different from *PVD [16]*, our dataset is captured closer to the subject and thus much more suitable for analyzing and visualizing Structure Swap. However, since compact cameras have lower image qualities, this dataset is only suitable for visualization and analysis but quantitative evaluation. Please refer to the supplementary for more details.

## 5.2. Metics and Competing Methods

We follow DyNeRF to use these metrics:(1) Peak Signal-To-Noise Ratio (PSNR), (2) Perceptual Score (LPIPS [53]), and (3) Just-Objectionable-Difference(JOD) [24]. We omit FLIP and DSSIM since we deem the three widely used metrics enough for reliable measurement. The following methods are evaluated:

**NeRF [26]**: Our baseline NeRF is trained from scratch and stored for each frame separately. This baseline helps to
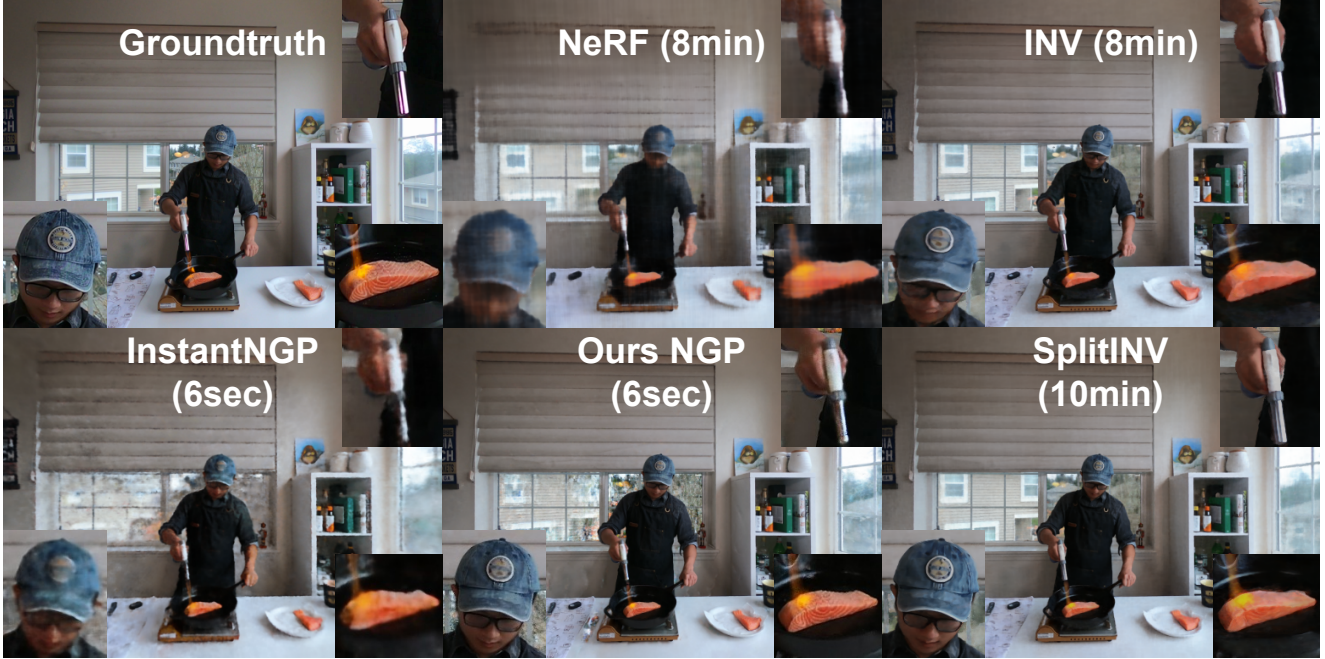
Figure 6. **Visual Comparisons**: INV achieves significantly better results than NeRF training just 8min/frame, and SplitINV renders even better foreground with just 10min/frame. "Ours NGP" is trained incrementally and achieves significantly better results than the original InstantNGP. Notice that "Ours NGP" achieves good LPIPS score but low PSNR score because it achieves high qualities on large portions of the scene, but often suffers from large amounts of floating artifacts (*e.g.* windows, around the person)

quantify the training acceleration and storage savings provided by INV (which uses the same NeRF internally).

**Neural 3D Videos (DyNeRF) [16]**: DyNeRF proposed to use learnable latent codes as inputs for each frame (in addition to positional encoding) to improve the rendering quality. DyNeRF also accelerates training via hierarchical training and ray importance sampling. Longer videos are divided into chunks of 10-second clips for better results. The authors also released the *Plenoptic Video Datasets [16]*. They show showed superior performance to their baselines. Since the **DyNeRF code is not available** and only results on the "flame salmon" sequence are reported, we cannot provide side-by-side visual comparisons and can only compare with them on that sequence.

**InstantNGP [27]**: We train InstantNGP on the scene with 800 iterations (6̃sec) from scratch for each frame.

**INV**: INV uses the original NeRF with the same settings as our NeRF baseline, but we disable skip connections. This is because they activate more Structure Layers, leading to higher storage costs while providing insignificant improvements. In most experiments, we train 3 Structure Layers. Although the frozen and shared layers also include some Structure Layers, we count them as part of the Shared Color Layers for simplicity. "INV+TWC" uses TWC to batches and compress all frames after Warm-Up. "SplitINV" uses a separate frozen NeRF to model the backgrounds, and INV

is used to capture the foreground. All metrics are measured on all 300 frames, including the Warm-Up frames which show worse results.

**SplitINV**: We use a separate NeRF to model the static background, while the foreground is encoded by INV. This allows INV to focus computation on the foreground and also helps stabilize the background. Please refer to the supplementary for more details on training.

**Ours NGP**: We train InstantNGP on the scene with 800 iterations/frame with the previous frame as the initialization for the MLP and hash grid. Importantly, **the hash grid of NGP can be seen as the output of Structure Layers stored in 3D, and the MLP can be seen as the Color layers**. In this way, NGP skips the Structure Layers and directly accesses the Color layers. Therefore, after 30 frames of warmup, we fix the MLP and train the hash grid only on the foreground pixels (as estimated by an optical flow estimator [51]). This allows NGP to focus computation on the foreground and also helps stabilize the background.

**Multi View Stereo (MVS), Local Light Field Fusion (LLFF), NeuralVolumes (NV), and NeRF-T**: We adopt the numbers as reported by DyNeRF [16]. Please refer to [16] for more details.

Figure 7. **Free-Viewpoint Renderings**: We show free viewpoint renderings on the 200th frames of sequences "flame salmon 1", "cut roasted beef", "sear steak", and "coffee martini".

## 5.3. Results

As shown in Fig. 3 and 5, INV's quality quickly and continually improves during the Warm-Up Stage. During the Structure Transfer Stage, INV performance stabilizes at a good quality. Moreover, INV can model challenging visual effects (*e.g.* flames, reflections, *etc.*) even though the Color Layers are frozen and shared between frames.

As shown in Table. 1, INV achieves state-of-the-art results on per-frame quality metrics (PSNR & LPIPS) with 50 minutes/frame less ($-19\%$) training than DyNeRF. Moreover, with merely 8min/frame of training, we achieve an average PSNR of 28.77db ("INV+TWC"). We also notice that more trainable Structure Layers result in better per-frame image metrics, but it almost completely removes the stabilization effect of freezing later layers as seen in the table and supplementary videos. This is because more structural contents are allowed to change with no guarantee that the change is consistent over time.

In "INV⋆", we show that well-trained Shared Color Layers can notably improve quality. During Warm-Up, the models are trained for 210min/frame. During Structural Transfer, the Structure Layers are only trained for 15min/frame. However, compared to the basic INV trained for 15min/frame throughout the video, PSNR increases by 0.42db (from 28.95db to 29.38db). This quality is similar

to the basic INV trained for 90min/frame. Pretraining INV could potentially be useful for many live-streaming settings (*e.g.* game streaming, lectures, and remote meetings). Since many live-streaming sessions capture the same background and people, it could be helpful to maintain or continually improve a long-term memory via a color/style space that is shared across sessions. Moreover, "SplitINV" increases the stability via its frozen background NeRF.

Due to its incremental nature, INV can synthesize 3D videos on-the-fly and does not have to wait for video chunks before processing. This property makes it naturally suitable for interactive 3D streaming. Additionally, INV doesn't make strong assumptions about the modeling backbone, thus it can likely be enhanced by more sophisticated techniques.

## 6. Limitations and Future Work

The main limitation of our work is visual stability. INV models that are trained longer (*e.g.* 210min/frame) generally show good stability. However, short training (*e.g.* 8min/frame) results in poor stability despite good per-frame qualities. While background models could improve stability on static contents, it is still challenging to stabilize dynamic contents effectively. It is likely helpful to enforce consistency for temporal correspondences (*e.g.* as done in

[11, 16, 17]). Moreover, it could improve stability and quality by training an intelligently selected subset of neurons, instead of always freezing certain layers. We leave this for future research.

While INV shows promising training acceleration, more needs to be done to achieve real-time training in the future. One advantage of INV is that it retains and continually improves upon previous information. This property could be combined with other approaches (*e.g.* smarter sampling, human-specific modeling, voxel-based approaches, etc.) to further reduce training time.

Additionally, we did not address the need for real-time rendering in this work. Since there are many prior works [12, 13, 27, 49] on real-time rendering for NeRFs, one might be able to achieve real-time streaming and rendering by combining INV with these approaches, assuming that the INV frames are already trained.

Another future work is to handle large and/or abrupt changes, as well as moving cameras. A simple solution is to detect these sudden changes and update the entire model. It could also be beneficial to dynamically determine which layers or neurons to train. However, this is likely less of a problem for many live-streaming scenarios with still cameras and smoothly varying content (*e.g.* game streaming, lectures, *etc.*). For general videos, it is also possible to store dedicated Shared Color Layers for each scene.

## 7. Conclusions

In this work, we introduced INV, an efficient frame-by-frame representation naturally suitable for interactive 3D streaming. It significantly reduces training while compressing storage sizes to the streamable realm. Additionally, we find that MLPs naturally partition themselves into Structure and Color Layers. Our findings can thus help the community better understand and manipulate MLP models. We believe that this is a solid step toward the future of interactive 3D streaming.

## References

[1] Alex Yu and Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks, 2021. 2

[2] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. 2020. 2

[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. 1, 3

[4] Benjamin Attal, Selena Ling, Aaron Gokaslan, Christian Richardt, and James Tompkin. Matryodshka: Real-time 6dof video view synthesis using multi-sphere images. *CoRR*, abs/2008.06534, 2020. 2

[5] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields, 2021. 2

[6] Michael Broxton, John Flynn, Ryan Overbeck, Daniel Erickson, Peter Hedman, Matthew DuVall, Jason Dourgarian, Jay Busch, Matt Whalen, and Paul Debevec. Immersive light field video with a layered mesh representation. 39(4):86:1–86:15, 2020. 2

[7] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision (ECCV)*, 2022. 2

[8] 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

[9] 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

[10] 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. 2

[11] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. *CoRR*, abs/2105.06468, 2021. 1, 2, 3, 10

[12] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. *arXiv preprint arXiv:2103.10380*, 2021. 2, 10

[13] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. *ICCV*, 2021. 2, 10

[14] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. *CoRR*, abs/2109.02369, 2021. 2

[15] 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

[16] 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. 1, 2, 3, 4, 7, 8, 10

[17] 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. 1, 2, 3, 10

[18] D. B.* Lindell, J. N. P.* Martel, and G. Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *Proc. CVPR*, 2021. 2

[19] Peter Lindstrom. Fixed-rate compressed floating-point arrays. *IEEE Transactions on Visualization and Computer Graphics*, 20, 08 2014. 6

[20] Peter Lindstrom and Martin Isenburg. Fast and efficient compression of floating-point data. *IEEE transactions on visualization and computer graphics*, 12:1245–50, 09 2006. 6

[21] 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

[22] Stephen Lombardi, Tomas Simon, Jason M. Saragih, Gabriel Schwartz, Andreas M. Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *CoRR*, abs/1906.07751, 2019. 7

[23] 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

[24] Rafał 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. 7

[25] 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. 3

[26] 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. 2, 7

[27] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv:2201.05989*, Jan. 2022. 2, 5, 7, 8, 10

[28] Phong Nguyen, Nikolaos Sarafianos, Christoph Lassner, Janne Heikkila, and Tony Tung. Human view synthesis using a single sparse rgb-d input. 2021. 3

[29] 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

[30] 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. 3

[31] 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. 3

[32] 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

[33] 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

[34] Gernot Riegler and Vladlen Koltun. Free view synthesis. *CoRR*, abs/2008.05511, 2020. 2

[35] 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. 2

[36] Darius Rückert, Linus Franke, and Marc Stamminger. Adop: Approximate differentiable one-pixel point rendering. *arXiv preprint arXiv:2110.06635*, 2021. 2

[37] 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. 2

[38] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *arXiv*, 2020. 2

[39] Pratul P Srinivasan, Richard Tucker, Jonathan T Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 175–184, 2019. 2

[40] 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

[41] Richard Szeliski and Polina Golland. Stereo matching with transparency and matting. In *Sixth International Conference on Computer Vision (IEEE Cat. No. 98CH36271)*, pages 517–524. IEEE, 1998. 2

[42] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *NeurIPS*, 2020. 2

[43] 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. 3

[44] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 2

[45] 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

[46] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. SynSin: End-to-end view synthesis from a single image. In *CVPR*, 2020. 2

[47] 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. 3

[48] 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. 3

[49] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 2, 10

[50] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images, 2020. 2

[51] Feihu Zhang, Oliver J. Woodford, Victor Adrian Prisacariu, and Philip H.S. Torr. Separable flow: Learning motion cost volumes for optical flow estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10807–10817, 2021. 8

[52] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *CoRR*, abs/2010.07492, 2020. 2

[53] 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. 7

[54] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018. 2