

# LIM: Large Interpolator Model for Dynamic Reconstruction

Anonymous CVPR submission

Paper ID 1064

## Abstract

Reconstructing dynamic assets from video data is central to many in computer vision and graphics tasks. Existing 4D reconstruction approaches are limited by category-specific models or slow optimization-based methods. Inspired by the recent Large Reconstruction Model (LRM) [15], we present the Large Interpolation Model (LIM), a transformer-based feed-forward solution, guided by a novel causal consistency loss, for interpolating implicit 3D representations across time. Given implicit 3D representations at times  $t_0$  and  $t_1$ , LIM produces a deformed shape at any continuous time  $t \in [t_0, t_1]$  delivering high-quality interpolations in seconds (per frame). Furthermore, LIM allows explicit mesh tracking across time, producing a consistently uv-textured mesh sequence ready for integration into existing production pipelines. We also use LIM, in conjunction with a diffusion-based multiview generator, to produce dynamic 4D reconstructions from monocular videos. We evaluate LIM on various dynamic datasets, benchmarking against image-space interpolation methods (e.g., FiLM [41]) and direct triplane linear interpolation, and demonstrate clear advantages. In summary, LIM is the first feed-forward model capable of high-speed tracked 4D asset reconstruction across diverse categories.

## 1. Introduction

Reconstructing dynamic 4D assets from video data is a fundamental problem in computer vision and graphics, with many virtual and augmented reality applications. Existing 4D reconstructors follow two main paradigms: category-specific articulated reconstruction and image-or-text conditioned 4D distillation. Hence, they are either restricted to a specific class of objects such as humans [32] and pets [3, 44], or are optimization-based [42, 63] making them slow requiring minutes to hours per reconstruction.

Recently, in the context of static reconstruction, the large reconstruction model (LRM) [16] has been proposed as an elegant feed-forward network that, starting from a fixed rig of multiview images, directly produces 3D implicit repre-

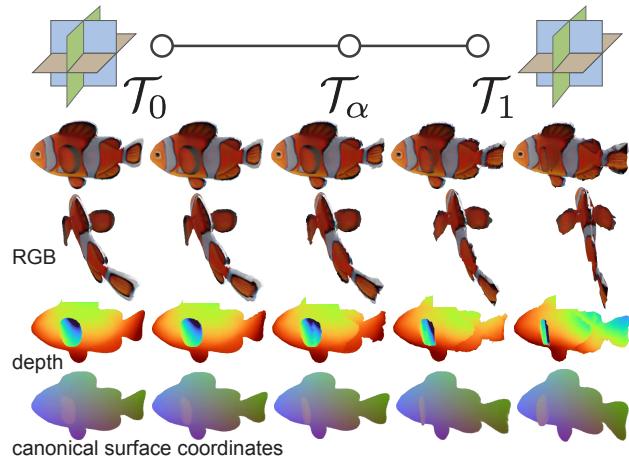


Figure 1. **Large Interpolator Model (LIM)** outputs a 4D video reconstruction by interpolating 3D implicit representations between two consecutive keyframes at times  $t = 0$  and  $t = 1$ .

sentation, which can then be rendered for novel view generation. In this work, in the context of dynamic reconstruction, we ask if a similar feed-forward approach can be developed to reconstruct a tracked explicit representation across time.

Here, L4GM [43] proposed a feedforward 4D video reconstructor which, for each video keyframe, accepts few views of the reconstructed object and outputs a mixture of 3D Gaussian Splats [23]. However, this approach has limitations as it can only reconstruct the keyframes at their exact timesteps without the ability to interpolate the shape through time. Additionally, establishing correspondences between Gaussian mixtures from different timesteps is challenging, which complicates tracing the deformation of the underlying object geometry through time. This limitation hinders many important downstream applications, such as gaming, where we require the 3D shape and texture of a single mesh to be defined in a static canonical pose while only its geometry (i.e., vertices) is allowed to be deformed arbitrarily.

We thus present Large Interpolation Model (LIM) as a transformer-based feed-forward solution that accepts an im-

038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059

060 plicit representation of an object at two different keyframe  
061 times  $t_0$  and  $t_1$  of a video, and interpolates between the two  
062 at any continuous intermediate timestep  $t \in [t_0, t_1]$ . We  
063 enable this with a novel self-supervised *causal consistency*  
064 *loss* that allows us to meaningfully interpolate continuously  
065 in time, even when supervised with keyframes from distinct  
066 time stamps. LIM is not only an efficient interpolator, but  
067 can also track a source mesh across time producing a func-  
068 tional deformable 3D asset with a shared uv texture map.  
069 Here, LIM tracks the mesh by means of an additional vol-  
070 umetric function that maps each time-specific 3D implicit-  
071 surface point to a unique coordinate on the intrinsic (time-  
072 invariant) surface of the object. This is unique – unlike any  
073 other competing dynamic reconstructor [43], LIM outputs  
074 a mesh with time-invariant texture and topology, and time-  
075 dependent vertex deformation. This renders LIM directly  
076 applicable in existing production setups.

077 Our LIM module also enables dynamic reconstruction  
078 from monocular video. Specifically, given keyframes of a  
079 monocular video, a pretrained image diffusion model gen-  
080 erates additional object views which, using a multiview  
081 LRM, we convert to keyframe-specific implicit 3D repre-  
082 sentations. Then, LIM directly interpolates the 3D repre-  
083 sentations yielding a dynamic 4D asset.

084 Our experiments demonstrate that LIM outperforms ex-  
085 isting alternatives in terms of the overall quality of the  
086 implicit-shape interpolations while being several times  
087 faster. Furthermore, we also evaluate the quality of the  
088 mesh tracing, where LIM records significant performance  
089 improvements.

## 090 2. Related Work

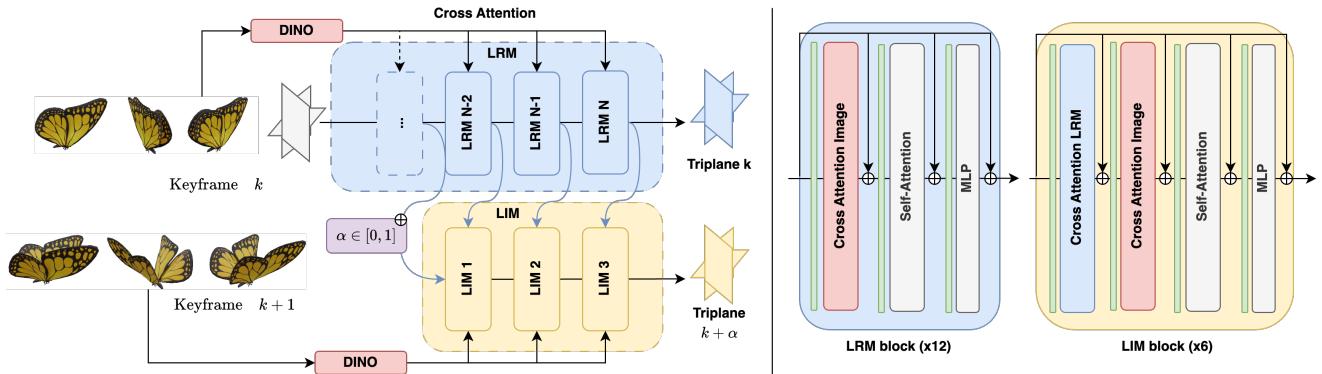
091 **3D Reconstruction.** Early work, introduced by Dream-  
092 Fusion [39] optimizes a 3D scene via *score distillation*  
093 *sampling* from a pretrained text-to-image diffusion model  
094 [35, 40, 48]. However, these methods are slow to opti-  
095 mize and suffer from inconsistencies (like the Janus prob-  
096 lem). Zero123 [45] learns to condition diffusion models  
097 on a single-view image and camera transformation, which  
098 allows novel view generation. Multiple novel views of  
099 a single object can then be used to optimize a NeRF re-  
100 construction which, however, is often impaired by view-  
101 inconsistencies of the novel-view generator. SyncDreamer  
102 [31] proposes an extension that improves the consistency of  
103 novel views and transitively of the generated 3D shapes.  
104 Due to the highly-challenging task of reconstructing any  
105 3D asset from a single image, several works [13, 14, 19–  
106 21, 25, 26, 28, 28, 29, 52, 55, 56, 60] learn 3D recon-  
107 structors of a specific category which simplifies learning of  
108 shape, deformation, and appearance priors. We also note  
109 some works on learning a generalizable dynamic radiance  
110 field from monocular videos [47, 49] which, however, are  
111 not designed for outputting a time-deforming 3D mesh. Re-

cent methods [15, 53], trained on large 3D datasets such  
112 as *Objaverse* [9, 10], propose feed-forward reconstructors  
113 which directly predict 3D representation of an object, con-  
114 ditioning on a single or multiple views. These methods dra-  
115 matically reduce reconstruction speed as they don’t rely on  
116 any optimization loop.

117  
**4D Representations.** Extending the popular research on  
118 representing static 3D scenes with implicit shapes, recent  
119 works proposed new time-deforming alternatives. *Dynerv*  
120 [12] extends static neural radiance field [36] with an addi-  
121 tional compact latent code to represent time deformation.  
122 However, similar to the static implicit shape reconstruc-  
123 tors [36, 61], its optimization process is relatively slow.  
124 [4, 8, 11] factorize a dynamic representation into multiple  
125 low-rank components, which dramatically speeds up the op-  
126 timization. Notably, *Hexplane* [4] proposes a 6-plane rep-  
127 resentation which extends the spacial triplane representation  
128 [7] to a spatio-temporal one. With the emergence of 3D  
129 *Gaussian Splatting* [22] (3DGS), [34, 54] propose its ex-  
130 tension to dynamic scenes, relying either on a per-frame op-  
131 timization with dynamic constraints, or on a temporal net-  
132 work to deform the gaussians in time.

133  
**4D Generation and Reconstruction.** Several works focus  
134 on text-to-4D generation: MAV3D [46] optimizes a Hex-  
135 plane [4] representation via score distillation sampling from  
136 a text-to-image and a text-to-video diffusion model. 4Dify  
137 [2] introduces a 3D-aware text-to-image diffusion model,  
138 and parameterizes the representation with a multi-resolution  
139 hash encoding [37]. However, these methods tend to pro-  
140 duce very limited and simple motions. TC4D [1] proposes  
141 an extension to decompose movement into local deforma-  
142 tion and global rigid motion. [30] applies same SDS super-  
143 vision with Gaussian splatting.

144  
Similar to us, recent work focused on video-to-4D re-  
145 construction. Consistent4D [18] generates 4D content from  
146 monocular video via SDS supervision, optimizing a Cas-  
147 cade DyNerf [12]. Simiarly, 4DGen [62] and DreamGaus-  
148 sian4D [42] encode the 4D asset as a set of static 3D gaus-  
149 sians and a regularized deformation field. [38, 64] leverages  
150 diffusion models to generate frames across views and times-  
151 tamps, and optimizes a dynamic gaussian splats based on  
152 these frames. [17, 57] decompose motion and appearance  
153 in gaussian splatting: instead of predicting a deformation  
154 for each gaussian in a canonical frame, they deform Gaus-  
155 sians by means of sparse control points. All above methods  
156 are relatively slow due to the 2nd reconstruction stage that  
157 optimizes each 4D asset from scratch. Furthermore, these  
158 methods cannot easily trace the resulting 4D asset through  
159 time, which prohibits their application in production setups.



**Figure 2. LIM framework.** (Left) Given multi-view images on 2 timesteps  $k$  and  $k + 1$ , LIM interpolates any intermediate 3D representation at  $k + \alpha, \alpha \in [0, 1]$ . It achieves this notably via cross-attention with the latest intermediate features of LRM on keyframe  $k$ . In practice, our LIM architecture has 6 blocks and LRM 12 blocks. (Right) Block structure of LRM and LIM. We include layer normalization before each module in blocks.

161

### 3. Method

162 In Sec. 3.1, we review the LRM [15] architecture which  
 163 our method is based on; in Sec. 3.2, we introduce LIM,  
 164 our large interpolator model, for efficient 3D interpolation  
 165 and; in Sec. 3.4, we show how LIM can be used for fast 4D  
 166 reconstruction and mesh tracking.

167

#### 3.1. Preliminaries

168 Our Large Interpolation Model (LIM) is built upon the  
 169 multi-view version of Large Reconstruction Model (LRM)  
 170 [15]. We first review LRM and its multi-view version.

171 **LRM.** The LRM [15] is a single-view reconstructor. Given  
 172 a source image  $I_{\text{src}}$  and its camera  $\pi_{\text{src}}$ , LRM reconstructs  
 173 a triplane [7] representation  $\mathcal{T} := \text{LRM}_{\theta}(I, \pi)$  of the  
 174 depicted scene. The triplane may be rendered from any target  
 175 view  $\pi_{\text{tgt}}$  using Emission-Absorption raymarching yielding  
 176 an RGB render  $R(\pi_{\text{tgt}}, \mathcal{T})$ , depth render  $R(\pi_{\text{tgt}}, \mathcal{T}_D)$  and  
 177 alpha-mask render  $R(\pi_{\text{tgt}}, \mathcal{T})_{\alpha}$ . In practice, we use the  
 178 Lightplane renderer [5] to implement  $R$ .

179 In a single-view setting, 3D reconstruction is highly am-

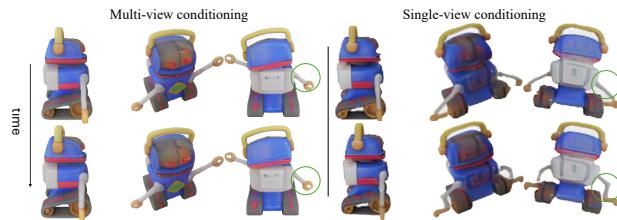


Figure 3. LRM conditioned on a single-view [50] is sensitive to small changes on the input image, which gives inconsistent result from one video frame to another. The multi-view LRM prevents this instability. For each model, left shows an input-view, right shows two target views. Each line is a different timestep.

biguous. Indeed, as depicted in Fig. 3, when applied to reconstruct monocular-video frames, LRM outputs triplanes with significantly time-inconsistent shape and texture.

**Multi-view LRM setup.** Hence, in order to minimize reconstruction ambiguity, we leverage a few-view conditioned version of LRM. Formally, given a set  $\mathcal{I}^{\text{src}} := \{I_{\text{src}}^i\}_{i=1}^{N^{\text{src}}}$  of  $N^{\text{src}}$  source images with corresponding cameras  $\Pi^{\text{src}} := \{\pi_{\text{src}}^i\}_{i=1}^{N^{\text{src}}}$  we predict a triplane  $\mathcal{T} := \text{LRM}_{\theta}(\mathcal{I}^{\text{src}}, \Pi^{\text{src}})$ , where we overload the same symbol for the multi-view and single-view versions for compactness. The architecture follows [27, 59] – the pixels of each source image  $I^{\text{src}}$  are first concatenated with the Plucker ray coordinates encoding the corresponding camera pose  $\pi^{\text{src}}$  and then fed to DinoV2 [6] yielding image tokens. Then, these tokens enter cross-attention layers inside a large 12-layer transformer that refines a set of fixed shape tokens into the final triplane representation  $\mathcal{T}$  of the reconstructed scene.

**Multi-view LRM training.** We train LRM in a fully-supervised manner on a large dataset of artist-created meshes, similar to Objaverse [10]. We render each mesh from a set of pre-defined camera viewpoints  $\Pi$ . The latter rendering, besides the RGB image  $I$ , also provides the ground-truth depth map  $D$  and the alpha mask  $M$ . For each training scene, we sample  $N_{\text{src}} = 4$  random images as input views, and render into  $N_{\text{tgt}} = 4$  randomly sampled held-out target views where losses are optimized.

We optimize three losses. (i) The photometric loss  $\mathcal{L}_{\text{photo}} := \sum_{i=1}^{N_{\text{tgt}}} \|I^i - R(\pi^i, \mathcal{T})\|^2 + \text{LPIPS}(I^i, R(\pi^i, \mathcal{T}))$ ; (ii) mask loss  $\mathcal{L}_{\text{mask}} := \sum_{i=1}^{N_{\text{tgt}}} \text{BCE}(M^i, R(\pi^i, \mathcal{T})_{\alpha})$ , where BCE is binary cross-entropy; and (iii) depth loss  $\mathcal{L}_{\text{depth}} := \sum_{i=1}^{N_{\text{tgt}}} \|D^i - R(\pi^i, \mathcal{T})_D\|$ . Recall that  $\mathcal{T} := \text{LRM}_{\theta}(\mathcal{I}^{\text{src}}, \Pi^{\text{src}})$  is the triplane output by LRM given the 4 source views. The total loss  $\mathcal{L}_{\text{photo}} + \mathcal{L}_{\text{depth}} + \mathcal{L}_{\text{mask}}$  is minimized with the Adam optimizer [24] with a learning rate of

214     $10^{-4}$  until convergence.

### 215    3.2. LIM: Large Interpolator Model

216    Given a monocular video, our aim is to predict the 3D representation of the scene at any continuous timestep. Furthermore, we aim to achieve this in a feed-forward manner, and we require the ability to trace the 3D representation in time, which eventually yields a practically applicable animated mesh with a shared UV texture.

222    **Multi-view LIM.** As mentioned in Sec. 3.1, reconstructing monocular videos is a highly ambiguous task and, hence, we first focus on the simpler multi-view version with access to multiple views at each timestep. At the end of this section, we describe how to tackle the harder monocular task by converting it to the multi-view setting described here.

228    Formally, we are given a multi-view RGB video  $\{\mathcal{I}_k\}_{k \in (1, 2, \dots, N_f)}$  composed of  $N_f$  timesteps where, for 229 each integer timestep  $k$ , we have a set  $\mathcal{I}_k = \{I_k^i\}_{i=1}^{N_v}$  of 230  $N_v$  view-points with cameras  $\Pi_k = \{\pi_k^i\}_{i=1}^{N_v}$ . In order to 231 4D-reconstruct the latter we can, in principle, use LRM to 232 predict a set  $\{\mathcal{T}_k\}_{k \in (1, 2, \dots, N_f)}$  containing a triplane for each 233 keyframe in the video. However, the latter remains **discrete** 234 in time and, hence, we cannot obtain a 3D representation 235 at any intermediate continuous timestep  $k + \alpha, \alpha \in [0, 1]$ . 236 Furthermore, such frame-specific triplanes encode implicit 237 shapes disconnected across different timesteps. This 238 prevents us from converting the time-series of reconstructions 239 into a time-varying mesh.

241    Thus, to achieve continuous reconstruction in time, and 242 to enable surface tracking, we introduce our Large Interpolator 243 Model (LIM). Given 2 keyframe sets  $\mathcal{I}_k, \mathcal{I}_{k+1}$  at discrete 244 timesteps  $k$  and  $k+1$ , LIM predicts an interpolated triplane 245  $\hat{\mathcal{T}}_{k+\alpha}$  at any continuous timestep  $t = k+\alpha, \alpha \in [0, 1]$ :

$$246 \quad \hat{\mathcal{T}}_{k+\alpha} := \text{LIM}_\psi(\mathcal{F}_k(\mathcal{I}_k, \Pi_k), \mathcal{I}_{k+1}, \alpha). \quad (1)$$

247    The architecture of LIM, illustrated in Fig. 2, takes advantage of the pretrained multiview LRM model from Sec. 3.1. 248 More specifically, we begin by calculating the intermediate 249 features  $\mathcal{F}_k$  as predicted by LRM from the frame set  $\mathcal{I}_k$  at 250 the start timestep  $k$ . These features are extracted after each 251 of the last  $L = 6$  transformer blocks of LRM. Then, we 252 broadcast and concatenate a positional encoding of the 253 interpolation time  $\alpha$  to  $\mathcal{F}_k$  and feed the result to LIM. This 254 input is then refined by series of cross-attentions with the 255 image tokens of the next keyframes  $\mathcal{I}_{k+1}$  to predict the final 256 interpolated triplane  $\hat{\mathcal{T}}_{k+\alpha}$ .

### 258    3.3. Training LIM

259    We train LIM on a large dataset of artist-created meshes 260 animated with a range of motions. For each scene, we render 261 the asset from several random viewpoints at each key-frame 262 of the animation.

263    In order to train LIM, for each scene, we first sample a 264 pair of keyframe interpolation endpoints at timesteps  $k_{\text{src}}$  265 and  $k_{\text{tgt}}$  such that  $k_{\text{tgt}} - k_{\text{src}} \in \{2, 3, 4\}$ . Then, we 266 additionally sample a middle keyframe  $k_m$  such that  $k_{\text{src}} \leq 267 k_m \leq k_{\text{tgt}}$ . We then task LIM to predict the interpolated 268 triplane  $\hat{\mathcal{T}}_{k_{\text{src}}+\alpha_m} := \text{LIM}(\mathcal{F}_{k_{\text{src}}}, \mathcal{I}_{k_{\text{tgt}}}, \alpha_m)$  at an intermediate 269 keyframe  $k_m$  given the source and target conditioning 270  $\mathcal{F}_{k_{\text{src}}}, \mathcal{I}_{k_{\text{tgt}}}$  and the interpolation time  $\alpha_m = \frac{k_m - k_{\text{src}}}{k_{\text{tgt}} - k_{\text{src}}}$ , 271 which converts the discrete timestep  $k_m$  into a continuous 272 interpolation time  $\alpha_m \in [0, 1]$ . The interpolated triplane 273  $\hat{\mathcal{T}}_{k_{\text{src}}+\alpha_m}$  is then compared to the pseudo-ground-truth triplane 274  $\mathcal{T}_{k_m} = \text{LRM}(\mathcal{I}_{k_m}, \Pi_{k_m})$  output by LRM at the interpolated 275 keyframe  $k_m$  with the following MSE loss:

$$276 \quad \mathcal{L}_{\mathcal{T}} := \|\hat{\mathcal{T}}_{k_{\text{src}}+\alpha_m} - \mathcal{T}_{k_m}\|^2, \alpha_m = \frac{k_m - k_{\text{src}}}{k_{\text{tgt}} - k_{\text{src}}}. \quad (2)$$

#### 277    Causal consistency for continuous-time interpolation.

278    The loss  $\mathcal{L}_{\mathcal{T}}$  provides a basic supervisory signal which, 279 however, only supervises LIM at keyframe times  $k_m$  that 280 are discrete. The latter prevents the model from becoming 281 a truly temporally-smooth interpolator because, during 282 training, it is never exposed to arbitrary interpolation times 283  $\alpha$  spanning the whole continuous range of  $[0, 1]$ .

284    To address this, we introduce a causal consistency loss 285  $\mathcal{L}_{\text{causal}}$ . In a nutshell, the loss enforces that a triplane 286 interpolated directly from time  $k_{\text{src}}$  to  $k_{\text{src}} + \delta, \delta \in [0, 1]$  has to 287 match a triplane that is first interpolated to an arbitrary 288 intermediate timestep  $k_{\text{src}} + \alpha_{\text{rand}}, \alpha_{\text{rand}} \in \mathcal{U}(0, \delta)$  and then 289 further interpolated to the target timestep  $k_{\text{src}} + \delta$ .

290    More formally, we define the causal consistency loss as:

$$291 \quad \mathcal{L}_{\text{causal}} := \left\| \text{LIM} \left( \hat{\mathcal{F}}_{k_{\text{src}}+\alpha_{\text{rand}}}, \mathcal{I}_{k_{\text{src}}+\delta}, \frac{\delta - \alpha_{\text{rand}}}{1 - \alpha_{\text{rand}}} \right) - \hat{\mathcal{T}}_{k_{\text{src}}+\delta} \right\|^2, \quad (3)$$

292    where  $\hat{\mathcal{F}}_{k_{\text{src}}+\alpha_{\text{rand}}}$  stands for the intermediate features 293 predicted by LIM when interpolating from  $k_{\text{src}}$  to  $k_{\text{src}} + \alpha_{\text{rand}}$ . 294 Note that we feed into the second LIM pass the intermediate 295 features  $\hat{\mathcal{F}}_{k_{\text{src}}+\alpha_{\text{rand}}}$  output by LIM as opposed to features 296  $\mathcal{F}$  output by LRM as prescribed by the original LIM 297 formulation in (1). We empirically observed that this works 298 as we can assume that the intermediate features of LIM follow 299 approximately the same distribution as the intermediate 300 features of LRM. Hence, LIM can, in a recurrent manner, 301 accept its own intermediate features to ground the interpolation 302 of the next timesteps. As demonstrated in Sec. 4.1 303 (Tab. 4), the causal loss  $\mathcal{L}_{\text{causal}}$  significantly improves the 304 temporal consistency and the quality of the interpolations.

305    Note that for LIM training, LRM model weights  $\theta$  are 306 already optimized and we keep them frozen. We minimize 307 the total loss  $\mathcal{L}_{\mathcal{T}} + \mathcal{L}_{\text{causal}}$  using the Adam optimizer [24] 308 with a learning rate of  $10^{-4}$  until convergence.

309    **Monocular 4D reconstruction.** Given the trained LIM 310 and LRM models, we can now predict the 3D scene rep-

311 representation at any continuous timestep. However, as mentioned above, the LIM model relies on multiple views at 312 each timestep, which are not always available in practice. 313

314 Our method can, however, also be used in the monocular- 315 video to 4D setting. Here, we leverage a pretrained diffu- 316 sion model [58] to recover 3 videos at different viewpoints, 317 consistent in shape and motion with the monocular source. 318 We then reconstruct 3D per timestep with LRM, and add in- 319 between timesteps (depending on the user frame-rate need), 320 by interpolating with LIM. This replaces the optimization 321 of a 4D representation from the multi-view videos [2, 46], 322 which in practice takes minutes to hours for a single scene.

### 323 3.4. Tracing shapes with LIM

324 In Secs. 3.2 and 3.3, we have described how LIM, together 325 with LRM, can be trained to predict a continuous-time 3D 326 representation of a scene given a set of multi-view images at 327 discrete timesteps. As mentioned before, a key goal of our 328 model is to also trace the deformable shape through time. 329 Next, we describe how to extend LIM and LRM to output 330 canonical surface coordinates to enable surface tracing, and 331 how to use these coordinates to output a time-deforming 332 mesh with fixed topology and texture.

333 **Interpolating canonical surface coordinates.** To simplify 334 the surface tracing task, we extend LIM and LRM to label 335 the interpolated implicit surface with intrinsic coordinates 336 defined in the canonical coordinate of the object to be 337 interpolated. More specifically, aside from tasking the 338 interpolated triplanes  $\mathcal{T}$  with representing the RGB color and 339 geometry of the implicit shape, we also task them with sup- 340 porting a volumetric function  $f : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  that maps each 341 point in the 3D space to its canonical surface coordinate. 342 Without loss of generality, we set the canonical coordinates 343 of time-deforming shape to the XYZ coordinates of the cor- 344 responding surface points in the start timestep  $k_{\text{src}}$ . Since 345 we have a dataset of artist-created meshes with a known de- 346 formation of each vertex in time, we can easily calculate the 347 canonical coordinate function  $f_{k_{\text{src}}}$  of each vertex at the start 348 timestep  $k_{\text{src}}$  and then transport those using the known ani- 349 mation to any other deformation time  $k_{\text{tgt}}$  yielding a func- 350 tion  $f_{k_{\text{tgt}}}$ . Importantly,  $f$  can be rendered at any point in 351 time from a viewpoint  $\pi^i$  yielding a 3-channel canonical 352 coordinate render  $C^i$ .

353 Given the above, we train a second LRM, dubbed  $\overline{\text{LRM}}$ , 354 which shares the same architecture, but predicts triplanes 355  $\mathcal{T}^C := \overline{\text{LRM}}(\mathcal{C}_{\text{src}}, \Pi_{\text{src}})$  supporting the coordinate function 356  $f$ , by accepting a set  $\mathcal{C}_{\text{src}} := \{C_{\text{src}}^i\}_{i=1}^{N_{\text{src}}}$  of multi-view source 357 canonical renders  $C_{\text{src}}^i$ . This canonical-coordinate  $\overline{\text{LRM}}$  is 358 supervised with the following canonical loss:

$$359 \mathcal{L}_{\text{can}} := \|C_{\text{src}}^i - R(\pi^i, \mathcal{T}_{\text{src}}^C)\|^2, \quad (4)$$

360 and with the same depth, and mask losses as the LRM 361 in Sec. 3.1. Similarly, we train  $\overline{\text{LIM}}$  to interpolate the

362 canonical-coordinate triplane as,

$$363 \hat{\mathcal{T}}_{\text{src}}^C := \overline{\text{LIM}}(\mathcal{F}_{k_{\text{src}}}^C, \mathcal{I}_{k_{\text{src}}}, \mathcal{I}_{k_{\text{tgt}}}, \alpha). \quad (5)$$

364 Note that this  $\overline{\text{LIM}}$ , analogous to LIM, is conditioned on 365 the features  $\mathcal{F}_{k_{\text{src}}}^C$  of the canonical-coordinate  $\overline{\text{LRM}}$ . How- 366 ever, differently from LIM,  $\overline{\text{LIM}}$  accepts target and source 367 RGB frames  $\mathcal{I}_{k_{\text{src}}}$  and  $\mathcal{I}_{k_{\text{tgt}}}$  instead of the target-time canon- 368 ical image  $\mathcal{C}_{k_{\text{tgt}}}$ . This is because the canonical coordinates 369 can only be carried forward in time and, as such, are not 370 available for the target timestep  $k_{\text{tgt}}$ . Hence,  $\overline{\text{LIM}}$  instead 371 learns how to propagate the coordinates by analyzing the 372 RGB frames that are available in both timesteps. We super- 373 vise  $\overline{\text{LIM}}$  with the MSE loss  $\mathcal{L}_{\text{f}}^C$  and the causal consistency 374 loss  $\mathcal{L}_{\text{causal}}^C$  that are defined analogously to (2) and (3) but 375 with the canonical-coordinate triplanes  $\mathcal{T}^C$  and images  $\mathcal{C}$ .

376 **Mesh tracing.** Given the RGB and canonical-coordinate 377 versions of LIM and LRM, we can trace a mesh multi-view 378 frames  $\mathcal{I}_{\text{src}}$  and  $\mathcal{I}_{\text{tgt}}$  (recall that, using an image diffusion 379 model, this is also possible for a monocular video).

380 We start by extracting the color triplane  $\mathcal{T}_{k_{\text{src}}}$  at timestep 381  $k_{\text{src}}$  with LRM. We then render  $\mathcal{T}_{k_{\text{src}}}$  to obtain a depth map 382  $D_{k_{\text{src}}}$  that we unproject to form 3D points yielding the multi- 383 view canonical coordinates  $\mathcal{C}_{k_{\text{src}}}$ . Given  $\mathcal{C}_{k_{\text{src}}}$ , we can predict 384 the canonical-coordinate triplane  $\mathcal{T}_{k_{\text{src}}}^C$  with  $\overline{\text{LRM}}$ .

385 Then, for a series of monotonic time offsets 386  $\alpha_0, \dots, \alpha_N; |\alpha_{j+1} - \alpha_j| \rightarrow 0$  the canonical coor- 387 dinate triplane  $\mathcal{T}_{k_{\text{src}}}^C$ , together with  $\mathcal{I}_{\text{src}}$  and  $\mathcal{I}_{\text{tgt}}$ , is fed to  $\overline{\text{LIM}}$  388 to interpolate the canonical-coordinate triplane  $\hat{\mathcal{T}}_{k_{\text{src}}+\alpha_j}^C$  at 389 all continuous timesteps  $k_{\text{src}} + \alpha_j$ .

390 The series of resulting canonical-coordinate triplanes 391  $\hat{\mathcal{T}}_{k_{\text{src}}+\alpha_j}^C$  provides a series of implicit shapes annotated with 392 surface coordinates. To obtain a time deforming mesh, we 393 first run Marching Cubes (MC) [33] on the first triplane 394  $\hat{\mathcal{T}}_{k_{\text{src}}+\alpha_0}^C$  resulting in a mesh  $\mathcal{M}_{k_{\text{src}}+\alpha_0}(V_{k_{\text{src}}+\alpha_0}, F)$  with 395 time-dependent vertices  $V_{k_{\text{src}}+\alpha}$  and time-invariant faces  $F$ . 396 We then run MC on the next triplane  $\hat{\mathcal{T}}_{k_{\text{src}}+\alpha_1}^C$  and match 397 the vertices of the previous mesh to the surface on the 398 next mesh using nearest neighbor search in the space of 399 canonical coordinates defined by the triplanes  $\hat{\mathcal{T}}_{k_{\text{src}}+\alpha_0}^C$  and 400  $\hat{\mathcal{T}}_{k_{\text{src}}+\alpha_1}^C$ , respectively. Afterwards, we replace the vertices 401  $V_{k_{\text{src}}+\alpha_0}$  with the corresponding nearest neighbors from the 402 next time  $k_{\text{src}} + \alpha_1$ , and repeat the process for all the re- 403 maining timesteps  $\alpha_2, \dots, \alpha_N$

## 4. Experiments

### 4.1. Feed-forward Triplane Interpolation

404 In this section, we evaluate the ability of LIM to interpolate 405 triplanes so that the renders of the latter match the ground- 406 truth views extracted at the target interpolation timestep. 407 More specifically, given a multi-view video of a deforma- 408 ting object, which contains the frame set  $\{\mathcal{I}_k\}_{k=1}^{N_f}$  at each 409 410

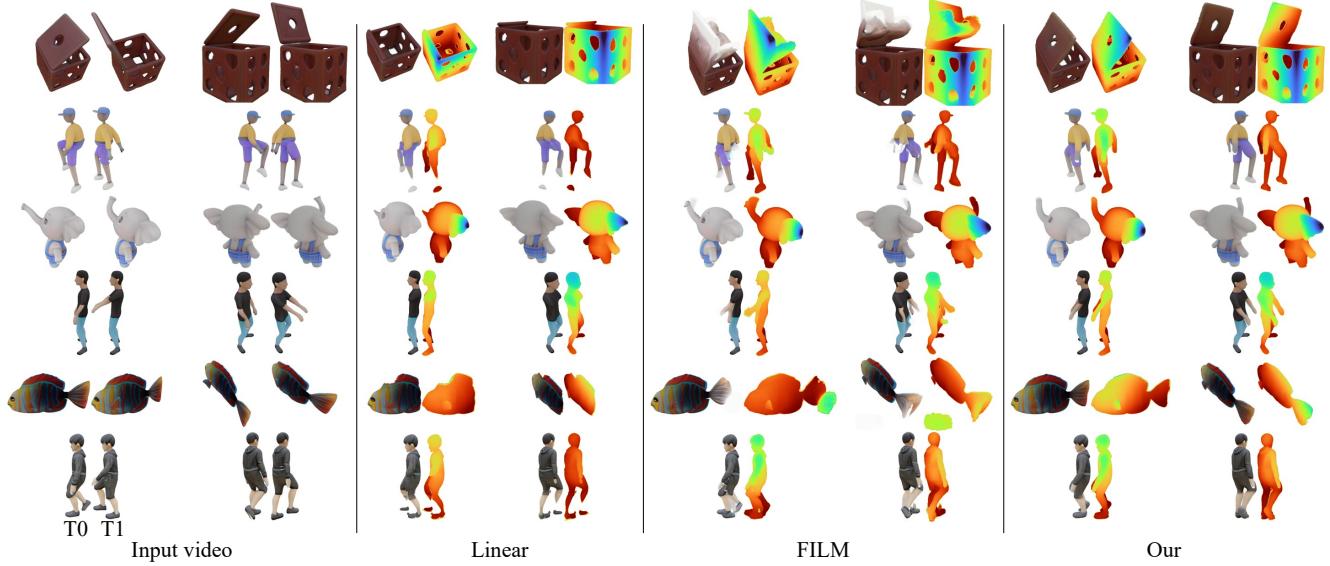


Figure 4. **Interpolation results** comparing (i) linear interpolation in triplane space, which fails on dynamic parts; (ii) image-based interpolator [41] (FILM), yielding view-consistent frame interpolations leading to defective reconstructions (ghosting around dynamic parts; for example, the tip of the elephant’s trunk or fish’s tail); and (iii) our LIM-based interpolation, which yields the most plausible results.

411 timestep  $k$ , we first split the frame sets into adjacent triplets  
 412  $\mathcal{I}_k, \mathcal{I}_{k+1}, \mathcal{I}_{k+2}$  for a  $k$ -th triplet. In each triplet, we then  
 413 evaluate the ability of a method to interpolate the 3D repre-  
 414 sentation  $\hat{\mathcal{T}}_{k+1}$  at the mid-point  $k + 1$  given the boundary  
 415 frames  $\mathcal{I}_k, \mathcal{I}_{k+2}$ . Note that, since the frames are sampled at  
 416 uniform time intervals, the continuous interpolation index  
 417 is  $\alpha = 0.5$ . For evaluation, we select 32 scenes from our  
 418 heldout dataset of animated objects.

419 Given the interpolant  $\hat{\mathcal{T}}_{k+1}$ , we then evaluate its quality  
 420 by rendering into the set of (novel) evaluation views, and re-  
 421 porting three photometric errors measuring the discrepancy  
 422 between the renders and the corresponding ground-truth im-  
 423 ages: (i) peak signal-to-noise ratio **PSNR**; (ii) perceptual  
 424 loss **LPIPS** [65]; and (iii) **PSNR<sub>FG</sub>** calculating PSNR only  
 425 over the foreground pixels.

426 **Baselines.** We compare our LIM interpolation with two  
 427 baselines. The first baseline (Linear) is a simple linear in-  
 428 terpolation, which defines the interpolated triplane  $\hat{\mathcal{T}}_{k+1}^{\text{linear}} =$   
 429  $(1 - \alpha)\mathcal{T}_k + \alpha\mathcal{T}_{k+2}$  as a linear combination of the two tri-

430 planes predicted by LRM for each set of boundary frames.  
 431 The second baseline (FILM) is image-based. Specifically,  
 432 we first interpolate in the image-space using a pre-trained  
 433 deep image interpolator *FILM* [41] which, given the bound-  
 434 ary views  $I_k^i, I_{k+2}^i$ , generates the interpolant  $\hat{I}_{k+1}^i$  followed  
 435 by multi-view LRM reconstruction yielding the interpo-  
 436 lated triplane  $\hat{\mathcal{T}}_{k+1}^{\text{FILM}}$ . We also report results from an Or-  
 437 acle method, which has access to the ground-truth images  
 438  $\mathcal{I}_{k+1}$  and reconstructs the corresponding triplane  $\hat{\mathcal{T}}_{k+1}^{\text{Oracle}}$   
 439 with LRM. The latter provides an upper performance limit.

440 **Results.** Tab. 1 presents the results. We also provide Fig. 4  
 441 and a supplemental video for visual evaluation. LIM out-  
 442 performs linear interpolation and image-based interpolation  
 443 on all three metrics. Here, we notice that linear interpo-  
 444 lation in the triplane space fails to correctly represent dy-  
 445 namic elements, which often disappear after being interpo-  
 446 lated. Furthermore, the image-based interpolation often re-  
 447 sults in artifacts in the color and opacity fields, which is due  
 448 to view-inconsistencies between the images interpolated at  
 449 the same timestep. LIM scores closest to the Oracle bound.

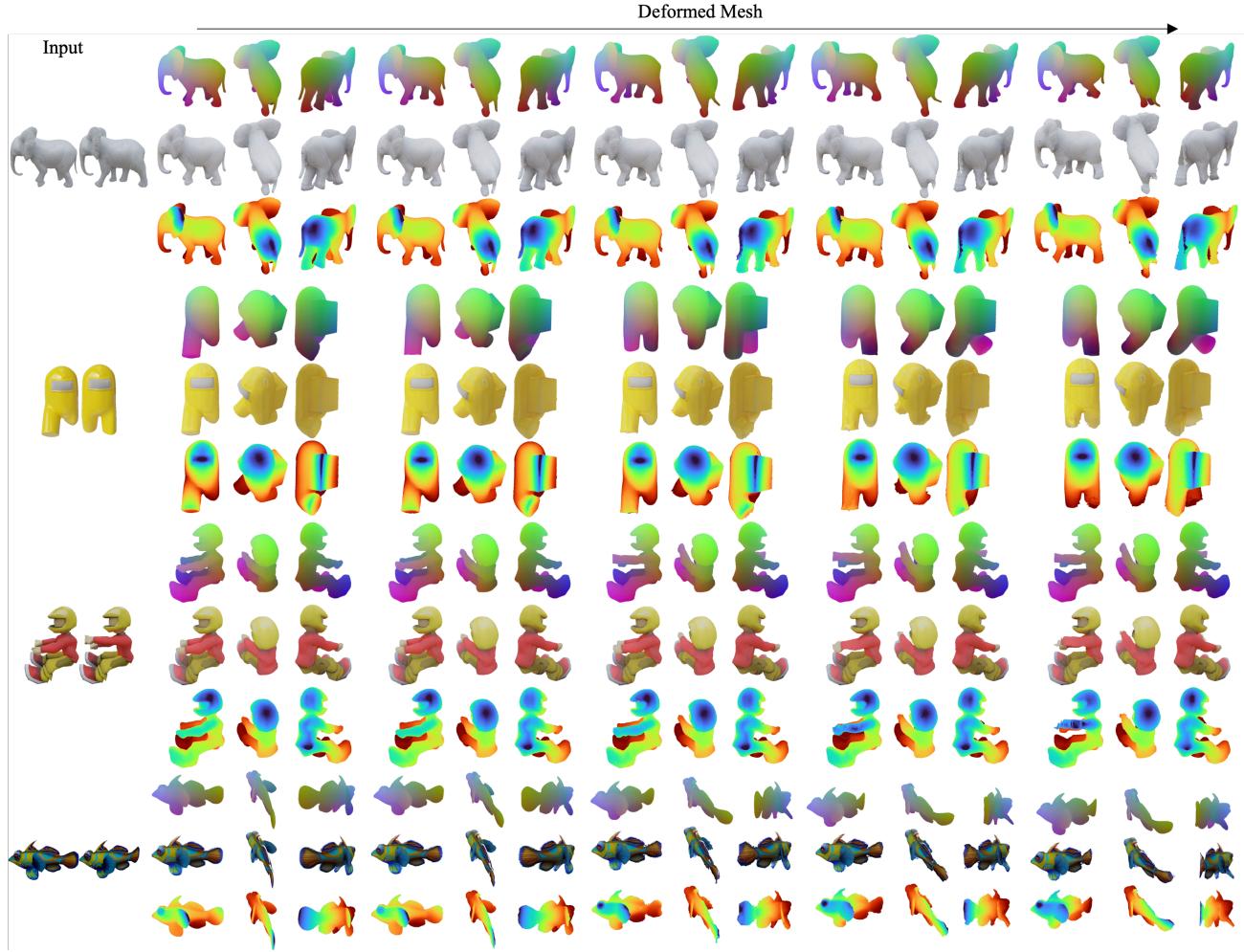
450 **Ablating causal consistency.** Table 2 reports the per-  
 451 formance of a LIM trained without the causal consistency  
 452 loss  $\mathcal{L}_{\text{causal}}$  on the aforementioned benchmark. The evalua-  
 453 tion reveals a significant drop in performance, confirming  
 454 the merit of the causal loss. See supplementary material.

## 4.2. Deformable Mesh Reconstruction

455 In this section, we evaluate the quality of the dynamic  
 456 mesh reconstructions output by LIM’s mesh-tracing method

Table 1. **Interpolation results** comparing LIM to a linear triplane interpolation, and to an image-based interpolation implemented with the FiLM image interpolator [41].

	PSNR $\uparrow$	PSNR <sub>FG</sub> $\uparrow$	LPIPS $\downarrow$
Linear	21.36	13.26	0.096
FILM [41]	22.61	15.44	0.091
LIM (Our)	<b>23.58</b>	<b>16.44</b>	<b>0.083</b>
Oracle	24.80	18.06	0.070



**Figure 5. Mesh Tracking results.** Given two implicit 3D representations, LIM can interpolate densely in time and hence can track a source mesh to produce a deforming mesh sequence. For each scene, we show (top to bottom) canonical-coordinate interpolation, RGB rendering of the tracked mesh and depth. See supplemental video.

**Table 2. Ablating  $\mathcal{L}_{\text{causal}}$ .** Interpolation accuracy comparing our LIM with its ablation removing the causal consistency loss  $\mathcal{L}_{\text{causal}}$ .

	PSNR $\uparrow$	PSNR <sub>FG</sub> $\uparrow$	LPIPS $\downarrow$
LIM- wo/ $\mathcal{L}_{\text{causal}}$	22.51	15.03	0.091
LIM	<b>23.58</b>	<b>16.44</b>	<b>0.083</b>

$\mathcal{I}_k, k \in [2, \dots, 8]$ . Note that we keep the topology of the first-frame mesh, as well as its texture shared across all 8 timesteps. After the mesh is traced, we render it at timesteps  $\{4, 6, 8\}$  to several heldout views and again evaluate the PSNR, LPIPS, and PSNR<sub>FG</sub>. We compare our LIM mesh tracing, described in Sec. 3.4, with a baseline (Nearest Matching) that iteratively deform the vertices of the mesh at

465

466

467

468

469

470

471

458 (Sec. 3.4). More specifically, we create a dataset of 8-  
459 step test sequences heldout from the train set. As before,  
460 the  $k$ -th timestep of the 8 timesteps contains a frame set  
461  $\mathcal{I}_k$ . First, we reconstruct a mesh in the canonical pose,  
462 defined as the shape at the first frame ( $k = 1$ ). We then  
463 use our mesh-tracing method to deform the vertices of the  
464 mesh so they follow the motion observed in the frame sets

**Table 3. Evaluation of deformable mesh tracking** comparing our LIM with Nearest-Neighbor tracing

	PSNR $\uparrow$	PSNR <sub>FG</sub> $\uparrow$	LPIPS $\downarrow$
NN-tracing	20.80	16.34	0.125
LIM (Our)	<b>21.92</b>	<b>17.46</b>	<b>0.103</b>
Oracle	23.40	17.80	0.096

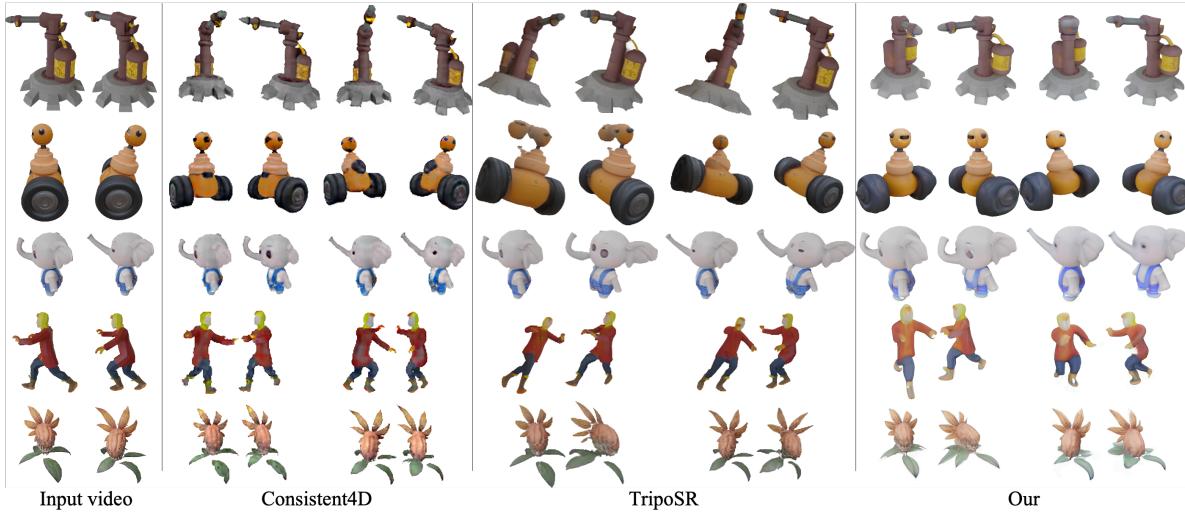


Figure 6. **Monocular 4D Reconstruction** comparing LIM with Consistent4D and TripoSR applied to each input frame separately.

472 timestep  $i$  to the nearest match on the surface of the mesh  
473 at timestep  $i + 1$ , where the matching is performed over  
474 distance and RGB-features.

475 **Results.** Tab. 3 present the quantitative results. Our main  
476 observation is that LIM’s ability to densely and accurately  
477 interpolate RGB and XYZ values over time, allows one  
478 to avoid explicitly solving the challenging correspondence  
479 problem between keyframed poses, separated by non-trivial  
480 deformations. Our evaluations show that LIM works well  
481 overall, although the results degrade around thin structures.  
482 Note that additional qualitative results are in the supplemen-  
483 tary material.

### 4.3. 4D Reconstruction

485 Finally, we evaluate LIM on 4D reconstruction from a  
486 monocular video. For each evaluation scene, we extract  
487 a video sequence of 16 frames  $\{I_k^1\}_{k=1}^{16}$ . We then lever-  
488 age an external diffusion-based model to generate 3 addi-  
489 tional views of the scene  $\{\hat{I}_k^i\}_{k \in [1, \dots, 16], i \in \{2, 3, 4\}}$  for each  
490 timestep. We then reconstruct the 3D representation on the  
491 odd frames with LRM and interpolate with our LIM to pre-  
492 dict the representation on the even frames. We evaluate on  
493 all the even frames, on a set of four random views, outside  
494 the training views. We report two metrics: the LPIPS er-  
495 rror between the ground-truth images and the renders of the  
496 reconstructed triplanes, and the FVD [51] measuring com-  
497 paring generated new-view sequence against ground-truth  
498 temporal sequence of renders.

499 **Baselines.** We compare two baselines: (i) Consis-  
500 tent4D [18], an optimization-based model, conditioned on  
501 monocular video and supervised via SDS loss; (ii) Tri-  
502 poSR [50], an open-source LRM conditioned on a single  
503 image, *i.e.*, the model is at a disadvantage.

504 **Results.** Evaluation results are in Tab. 4. The combination  
505 of LIM with multi-view diffusion model outperforms the  
506 competing methods by a significant margin.

Table 4. **Monocular video reconstruction** results comparing our LIM to Consistent4D [18] and to TripoSR [50] applied independently to each frame of the input video.

	Feed-fwd.	Inf. Time	LPIPS ↓	FVD ↓
Consistent4D [18]	✗	~1.5hours	0.484	1289.0
TripoSR [50]	✓	~30secs	0.515	1553.6
LIM (Ours)	✓	~3min	<b>0.119</b>	<b>720.6</b>

### 5. Conclusion

507 We have proposed LIM, a novel method paired with multi-  
508 view LRM to enable continuous and feed-forward render-  
509 ing in both space and time. As opposed to image-based  
510 interpolators or direct triplane baselines, we demonstrated  
511 that LIM results in high-quality and consistent 4D inter-  
512 polations, realistically capturing deformations, and can sup-  
513 port different type of modalities. A key advantage of ours is  
514 that we can interpolate in RGB and (canonical) XYZ, which  
515 in turn allows to directly output consistently-textured dy-  
516 namic mesh assets that applicable in production workflows.

517 A key limitation of our approach is it being trained on  
518 synthetic data. Since LIM learns to interpolate deforma-  
519 tion/motion, rather than appearance, we expect the results  
520 to carry over to real data. However, we would need an  
521 LRM model trained on real-world data to test this hypoth-  
522 esis. We leave this for later exploration. Also, in the future,  
523 we would like to extend our framework to handle extrap-  
524 olation, instead of interpolation. The challenge would be to  
525 effectively use video data and video generators for training.

527

## References

- [1] Sherwin Bahmani, Xian Liu, Yifan Wang, Ivan Skorokhodov, Victor Rong, Ziwei Liu, Xihui Liu, Jeong Joon Park, Sergey Tulyakov, Gordon Wetzstein, Andrea Tagliasacchi, and David B. Lindell. TC4d: Trajectory-conditioned text-to-4d generation, . 2
- [2] Sherwin Bahmani, Ivan Skorokhodov, Victor Rong, Gordon Wetzstein, Leonidas Guibas, Peter Wonka, Sergey Tulyakov, Jeong Joon Park, Andrea Tagliasacchi, and David B. Lindell. 4d-fy: Text-to-4d generation using hybrid score distillation sampling, . 2, 5
- [3] Benjamin Biggs, Thomas Roddick, Andrew W. Fitzgibbon, and Roberto Cipolla. Creatures great and SMAL: recovering the shape and motion of animals from video. In *Proc. ECCV*, 2018. 1
- [4] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. *arXiv.cs*, abs/2301.09632, 2023. 2
- [5] Ang Cao, Justin Johnson, Andrea Vedaldi, and David Novotny. Lightplane: Highly-scalable components for neural 3d fields. 3
- [6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proc. ICCV*, 2021. 3
- [7] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J. Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *Proc. CVPR*, 2022. 2, 3
- [8] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensoRF: Tensorial radiance fields. In *arXiv*, 2022. 2
- [9] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. Objaverse-XL: A universe of 10M+ 3D objects. *CoRR*, abs/2307.05663, 2023. 2
- [10] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3D objects. In *Proc. CVPR*, 2023. 2, 3
- [11] Sara Fridovich-Keil, Giacomo Meanti, Frederik Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. *arXiv.cs*, abs/2301.10241, 2023. 2
- [12] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. 2
- [13] Shubham Goel, Angjoo Kanazawa, and Jitendra Malik. Shape and viewpoint without keypoints. 2
- [14] Paul Henderson and Vittorio Ferrari. Learning single-image 3d reconstruction by generative modelling of shape, pose and shading. 2
- [15] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. LRM: Large reconstruction model for single image to 3d. 1, 2, 3
- [16] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. LRM: Large reconstruction model for single image to 3D. In *Proc. ICLR*, 2024. 1
- [17] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. SC-GS: Sparse-controlled gaussian splatting for editable dynamic scenes. 2
- [18] Yanqin Jiang, Li Zhang, Jin Gao, Weimin Hu, and Yao Yao. Consistent4d: Consistent 360° dynamic object generation from monocular video. 2, 8
- [19] Angjoo Kanazawa, Shubham Tulsiani, Alexei A. Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. 2
- [20] Hiroharu Kato and Tatsuya Harada. Learning view priors for single-view 3d reconstruction.
- [21] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. 2
- [22] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. 2
- [23] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3D gaussian splatting for real-time radiance field rendering. *Proc. SIGGRAPH*, 42(4), 2023. 1
- [24] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, San Diego, CA, USA, 2015. 3, 4
- [25] Filippos Kokkinos and Iasonas Kokkinos. Learning monocular 3d reconstruction of articulated categories from motion. 2
- [26] Filippos Kokkinos and Iasonas Kokkinos. To the point: Correspondence-driven monocular 3d category reconstruction, 2021. 2
- [27] Jiahao Li, Hao Tan, Kai Zhang, Zexiang Xu, Fujun Luan, Yinghao Xu, Yicong Hong, Kalyan Sunkavalli, Greg Shakhnarovich, and Sai Bi. Instant3D: Fast text-to-3D with sparse-view generation and large reconstruction model. *Proc. ICLR*, 2024. 3
- [28] Xueteng Li, Sifei Liu, Shalini De Mello, Kihwan Kim, Xiaolong Wang, Ming-Hsuan Yang, and Jan Kautz. Online adaptation for consistent mesh reconstruction in the wild. 2
- [29] Xueteng Li, Sifei Liu, Kihwan Kim, Shalini De Mello, Varun Jampani, Ming-Hsuan Yang, and Jan Kautz. Self-supervised single-view 3D reconstruction via semantic consistency. In *Proc. ECCV*, 2020. 2
- [30] Huan Ling, Seung Wook Kim, Antonio Torralba, Sanja Fidler, and Karsten Kreis. Align your gaussians: Text-to-4d with dynamic 3d gaussians and composed diffusion models. 2
- [31] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. SyncDreamer: Generating multiview-consistent images from a single-view image. *arXiv*, (2309.03453), 2023. 2

- 641 [32] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard  
642 Pons-Moll, and Michael J. Black. SMPL: a skinned multi-  
643 person linear model. *ACM Trans. on Graphics (TOG)*, 2015.  
644 1
- 645 [33] W. Lorensen and H. Cline. Marching cubes: A high reso-  
646 lution 3D surface construction algorithm. *ACM Computer  
647 Graphocs*, 21(24), 1987. 5
- 648 [34] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and  
649 Deva Ramanan. Dynamic 3d gaussians: Tracking by per-  
650 sistent dynamic view synthesis. 2
- 651 [35] Luke Melas-Kyriazi, Christian Rupprecht, Iro Laina, and  
652 Andrea Vedaldi. RealFusion: 360 reconstruction of any ob-  
653 ject from a single image. In *Proceedings of the IEEE Confer-  
654 ence on Computer Vision and Pattern Recognition (CVPR)*,  
655 2023. 2
- 656 [36] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik,  
657 Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF:  
658 Representing scenes as neural radiance fields for view syn-  
659 thesis. In *Proc. ECCV*, 2020. 2
- 660 [37] Thomas Müller, Alex Evans, Christoph Schied, and Alex-  
661 ander Keller. Instant neural graphics primitives with a multires-  
662 olution hash encoding. In *Proc. SIGGRAPH*, 2022. 2
- 663 [38] Zijie Pan, Zeyu Yang, Xiatian Zhu, and Li Zhang. Fast dy-  
664 namic 3d object generation from a single-view video. 2
- 665 [39] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Milden-  
666 hall. DreamFusion: Text-to-3D using 2D diffusion. In *Proc.  
667 ICLR*, 2023. 2
- 668 [40] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren,  
669 Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Sko-  
670 rokhodov, Peter Wonka, Sergey Tulyakov, and Bernard  
671 Ghanem. Magic123: One image to high-quality 3D object  
672 generation using both 2D and 3D diffusion priors. *arXiv.cs*,  
673 abs/2306.17843, 2023. 2
- 674 [41] Fitsum Reda, Janne Kontkanen, Eric Tabellion, Deqing Sun,  
675 Caroline Pantofaru, and Brian Curless. Film: Frame interpo-  
676 lation for large motion, 2022. 1, 6
- 677 [42] Jiawei Ren, Liang Pan, Jiaxiang Tang, Chi Zhang, Ang Cao,  
678 Gang Zeng, and Ziwei Liu. DreamGaussian4d: Generative  
679 4d gaussian splatting. 1, 2
- 680 [43] Jiawei Ren, Kevin Xie, Ashkan Mirzaei, Hanxue Liang, Xi-  
681 aohui Zeng, Karsten Kreis, Ziwei Liu, Antonio Torralba,  
682 Sanja Fidler, Seung Wook Kim, and Huan Ling. L4gm:  
683 Large 4d gaussian reconstruction model, 2024. 1, 2
- 684 [44] Remy Sabathier, Niloy J. Mitra, and David Novotny. Animal  
685 avatars: Reconstructing animatable 3d animals from casual  
686 videos. In *Computer Vision – ECCV 2024*, pages 270–287,  
687 Cham, 2025. Springer Nature Switzerland. 1
- 688 [45] Ruoxi Shi, Hansheng Chen, Zhuoyang Zhang, Minghua Liu,  
689 Chao Xu, Xinyue Wei, Linghao Chen, Chong Zeng, and Hao  
690 Su. Zero123++: a single image to consistent multi-view dif-  
691 fusion base model. *arXiv.cs*, abs/2310.15110, 2023. 2
- 692 [46] Uriel Singer, Shelly Sheynin, Adam Polyak, Oron Ashual,  
693 Iurii Makarov, Filippos Kokkinos, Naman Goyal, Andrea  
694 Vedaldi, Devi Parikh, Justin Johnson, and Yaniv Taigman.  
695 Text-to-4d dynamic scene generation. 2, 5
- 696 [47] Samarth Sinha, Roman Shapovalov, Jeremy Reizenstein, Ig-  
697 nacio Rocco, Natalia Neverova, Andrea Vedaldi, and David  
698 Novotný. Common pets in 3D: Dynamic new-view synthe-  
699 sis of real-life deformable categories. In *Proceedings of the  
700 IEEE Conference on Computer Vision and Pattern Recog-  
701 nition (CVPR)*, 2023. 2
- 702 [48] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi,  
703 Lizhuang Ma, and Dong Chen. Make-It-3D: High-fidelity 3d  
704 creation from A single image with diffusion prior. *arXiv.cs*,  
705 abs/2303.14184, 2023. 2
- 706 [49] Fengrui Tian, Shaoyi Du, and Yueqi Duan. MonoN-  
707 eRF: Learning a generalizable dynamic radiance field from  
708 monocular videos. 2
- 709 [50] Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan  
710 Huang, Adam Letts, Yangguang Li, Ding Liang, Christian  
711 Laforte, Varun Jampani, and Yan-Pei Cao. TripoSR: Fast 3d  
712 object reconstruction from a single image. 3, 8
- 713 [51] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach,  
714 Raphaël Marinier, Marcin Michalski, and Sylvain Gelly.  
715 Fvd: A new metric for video generation. In *ICLR*, 2019.  
716 8
- 717 [52] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei  
718 Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh  
719 models from single RGB images. 2
- 720 [53] Xinyue Wei, Kai Zhang, Sai Bi, Hao Tan, Fujun Luan,  
721 Valentin Deschaintre, Kalyan Sunkavalli, Hao Su, and Zex-  
722 iang Xu. MeshLRM: large reconstruction model for high-  
723 quality mesh. *arXiv*, 2404.12385, 2024. 2
- 724 [54] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng  
725 Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang.  
726 4d gaussian splatting for real-time dynamic scene rendering,  
727 . 2
- 728 [55] Shangzhe Wu, Tomas Jakab, Christian Rupprecht, and An-  
729 drea Vedaldi. DOVE: Learning deformable 3d objects by  
730 watching videos, . 2
- 731 [56] Shangzhe Wu, Ruining Li, Tomas Jakab, Christian Rup-  
732 precht, and Andrea Vedaldi. MagicPony: Learning articu-  
733 lated 3D animals in the wild. In *Proceedings of the IEEE  
734 Conference on Computer Vision and Pattern Recognition  
735 (CVPR)*, 2023. 2
- 736 [57] Zijie Wu, Chaohui Yu, Yanqin Jiang, Chenjie Cao, Fan  
737 Wang, and Xiang Bai. SC4d: Sparse-controlled video-to-4d  
738 generation and motion transfer, . 2
- 739 [58] Yiming Xie, Chun-Han Yao, Vikram Voleti, Huaiyu Jiang,  
740 and Varun Jampani. Sv4d: Dynamic 3d content generation  
741 with multi-frame and multi-view consistency, 2024. 5
- 742 [59] Yinghao Xu, Hao Tan, Fujun Luan, Sai Bi, Peng Wang, Ji-  
743 ahao Li, Zifan Shi, Kalyan Sunkavalli, Gordon Wetzstein,  
744 Zexiang Xu, and Kai Zhang. DMV3D: Denoising multi-  
745 view diffusion using 3D large reconstruction model. In *Proc.  
746 ICLR*, 2024. 3
- 747 [60] Gengshan Yang, Minh Vo, Natalia Neverova, Deva Ra-  
748 manan, Andrea Vedaldi, and Hanbyul Joo. BANMo: Build-  
749 ing animatable 3d neural models from many casual videos.  
750 2
- 751 [61] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman.  
752 Volume rendering of neural implicit surfaces. *arXiv.cs*,  
753 abs/2106.12052, 2021. 2

- 754 [62] Yuyang Yin, Dejia Xu, Zhangyang Wang, Yao Zhao, and  
755 Yunchao Wei. 4dgen: Grounded 4d content generation with  
756 spatial-temporal consistency. 2
- 757 [63] Yuyang Yin1, Dejia Xu2, Zhangyang Wang, Yao Zhao, and  
758 Yunchao Wei. 4DGen: Grounded 4D content generation with  
759 spatial-temporal consistency. *arXiv.cs*, 2023. 1
- 760 [64] Yifei Zeng, Yanqin Jiang, Siyu Zhu, Yuanxun Lu, Youtian  
761 Lin, Hao Zhu, Weiming Hu, Xun Cao, and Yao Yao.  
762 STAG4d: Spatial-temporal anchored generative 4d gaus-  
763 sians. 2
- 764 [65] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman,  
765 and Oliver Wang. The unreasonable effectiveness of deep  
766 features as a perceptual metric. In *CVPR*, 2018. 6