

# Dynamic Scene Novel View Synthesis via Deferred Spatio-temporal Consistency

Beatrix-Emőke Fülöp-Balogh<sup>a</sup>, Eleanor Tursman<sup>b</sup>, James Tompkin<sup>b</sup>, Julie Digne<sup>c</sup>, Nicolas Bonneel<sup>c,\*</sup>

<sup>a</sup>Univ. Lyon, UCBL, France

<sup>b</sup>Brown University, Rhode Island, USA

<sup>c</sup>Univ. Lyon, CNRS, France

## ARTICLE INFO

### Article history:

Received 27 April 2022

Received in final form 5 July 2022

Accepted 18 July 2022

Available online 25 July 2022

2000 MSC: 68U05, 68T45

## ABSTRACT

Structure from motion (SfM) enables us to reconstruct a scene via casual capture from cameras at different viewpoints, and novel view synthesis (NVS) allows us to render a captured scene from a new viewpoint. Both are hard with casual capture and dynamic scenes: SfM produces noisy and spatio-temporally sparse reconstructed point clouds, resulting in NVS with spatio-temporally inconsistent effects. We consider SfM and NVS parts together to ease the challenge. First, for SfM, we recover stable camera poses, then we *defer* the requirement for temporally-consistent points across the scene and reconstruct only a sparse point cloud per timestep that is noisy in space-time. Second, for NVS, we present a variational diffusion formulation on depths and colors that lets us robustly cope with the noise by enforcing spatio-temporal consistency via per-pixel reprojection weights derived from the input views. Together, this deferred approach lets us generate novel views for dynamic scenes without requiring challenging spatio-temporally consistent reconstructions nor training complex models on large datasets. We demonstrate our algorithm on real-world dynamic scenes against classic and more recent learning-based baseline approaches.

© 2022 Elsevier B.V. All rights reserved.

## 1. Introduction

Novel-view synthesis (NVS) creates a new view of a scene by combining existing images captured from different viewpoints. Much progress in NVS has been made over the past two decades to tackle its two core problems: 1) how to build a proxy scene geometry to aid in rendering, such as constructing simplified sparse depth points or a piecewise planar mesh via structure from motion (SfM), and 2) how to interpolate or extrapolate an image via the reprojected proxy given the existing captured imagery. NVS increases in difficulty across many axes [1, 2]: as the cameras become farther apart (wide baseline [3]), as their number decreases (few cameras [4]), as they become handheld (casual capture [5]), as the scene itself contains motion (dynamic scene

[6]), as the scene phenomena become more visually complex (geometry, materials, and motion [7]), and as the time given to generate the result decreases (compute cost [8]).

We consider dynamic scenes captured by a small number of cameras (5–12) over baselines of around  $60^\circ$  ( $\approx 1\text{m}$  for close scenes up to  $\approx 3\text{m}$  for far scenes), as might occur with a crowd of people capturing an event (Figure 1). We assume that these videos are synchronized. Within this scenario, we include sequences with casual handheld cameras. This is a relatively rare and challenging setting because both the cameras and the scene objects move simultaneously, and because sequences with only a small number of casual cameras makes robustness hard to obtain. This complicates camera pose estimation and depth estimation in SfM [9] and, if the proxy geometry is not perfect, causes ghosting, bleeding, and flickering artifacts across views and time during NVS in both moving objects and the background [10]. Thus, one key component of any algorithm is a way to enforce

\*Corresponding author: Nicolas Bonneel

e-mail: [nicolas.bonneel@liris.cnrs.fr](mailto:nicolas.bonneel@liris.cnrs.fr) (Nicolas Bonneel)

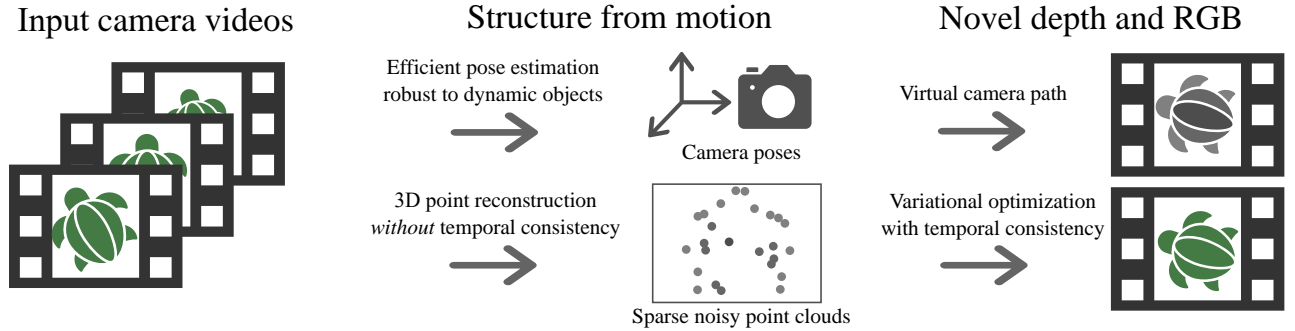


Fig. 1: **Overview.** Given a small set of video sequences of a performance, our method computes camera poses and sparse points, then optimizes those points into a novel video sequence following a user-defined camera path. Our space-time SfM intentionally does not compute temporal consistency for points on dynamic objects and instead defers spatio-temporal consistency in both depth and RGB reconstructions to the novel view synthesis stage via our variational formulation.

spatio-temporal consistency in both the SfM and the NVS.

We propose to address these challenges by *deferring* the difficult problem of reconstructing dynamic objects in time via SfM, and instead using a NVS approach to enforce temporal consistency. To ease the task of reconstructing dynamic scenes via SfM, many approaches first segment out moving objects or feature points and process the static background and the dynamic foreground separately [9, 11, 12]. Instead, we first recover camera poses for all views without any explicit dynamic object segmentation. Then, we recover scene points on both static and dynamic objects without temporal consistency and perform per-frame SfM across views only. This is easier to solve, but leads to significantly noisy reconstructions temporally.

Next, we turn our sparse (and noisy) reconstructed point clouds into novel views. This is commonly completed by densifying points into a depth map [13] for each view in a consistent way, and using the depths to reproject and merge input RGB views into a virtual view. We present a formulation which only densifies a depth map in the virtual camera’s view, rather than for all input views, which leads to a more efficient solve. For this, we take a coarse-to-fine variational approach and solve a diffusion-based formulation. Importantly, this formulation lets us enforce robust temporal consistency in the output depth to overcome the initial noisy reconstructions from the SfM. To determine our final RGB values, we also solve for the output color within the coarse-to-fine variational formulation.

We perform comparisons to recently-proposed approaches in point densification and view interpolation, using both optimization and learning-based approaches. Further, we show results on a synthetic dataset in an ablation study. In a nutshell, we show that considering SfM and NVS together allows us to ease the difficult temporally-consistent reconstruction problem and instead cope with it at the rendering stage. Overall, our work takes another step forward in improving digital content creation for scenes captured by multiple video cameras.

## 2. Related work

Rendering a novel viewpoint of a real-world scene captured with photographs is a problem that has received much attention over the past 30 years [1].

*Static scene IBR.* Image-based rendering (IBR) has initially attempted to render static scenes either from set of images or videos. This can be achieved either via warping input views using optical flow [14], using coarse geometric proxies [15] or via deep learning approaches [16]. In complex environments, IBR techniques often need some 3D proxy reconstruction. For example, the Lumigraph [15, 17] uses planar or coarse geometric proxies; Shade et al. [18] used multiple planar sprites; and Debevec et al. [19] employed photogrammetric reconstructions of buildings. Others have used 3D meshes from multi-view stereo reconstructions [20, 21]. For instance, Chaurasia et al. [22] proposed a depth-based synthesis using planar superpixel patches [23]. Matzen et al. [24] used two spherical cameras to synthesize an omni-directional stereo panorama. Recently, Riegler and Koltun [25] synthesized new views via neural textures atop a Delaunay reconstruction of sparse points obtained from video of static scenes. Beyond surface geometry, NeRF [26] performs an expensive optimization to create a volumetric function that is then rendered to synthesize new views.

Solving problems in the gradient domain can help too; for instance, to achieve smoother interpolations [7] or to densify sparse scene points. Holynski and Kopf spatio-temporally propagate sparse depth samples in a single view by solving a Poisson problem [13]. Their method relies on camera motion to detect depth edges, which limits it to static scenes. Inspired by gradient domain approaches, we formulate a variational approach that jointly enforces depth smoothness and consistency, color smoothness and consistency, as well as temporal consistency. Our approach additionally works with multiple potentially-dynamic cameras, and introduces a view-consistency term to ensure geometric consistency between views.

Deep learning can also be employed for static scene IBR. This includes plane sweep volumes [27] and multi-plane images to interpolate between two static narrow-baseline views [28] or between multiple views at once [29, 16], appearance flows to generate novel views from a single image of isolated objects [30], and light-field view interpolation [31]. Hedman et al. [5] use a geometric proxy and learn blending weights between view reprojections using a CNN. To improve the quality around depth discontinuities, Choi et al. [4] use a 3D uncertainty volume as a proxy and neural network-based patch refinement. Srinivasa et al. [32] train a CNN to predict a light field from a single image

for small-baseline view synthesis. Similarly, Song et al. [33] synthesize new views from a single image of a static scene.

While these techniques were not designed for videos and so neither explicitly maintain temporal consistency nor are constrained by speed, we nevertheless compare our approach to relevant methods for static scenes taken frame by frame.

*Dynamic scene VBR.* For dynamic scenes, please see dos Anjos et al. [2] for an exhaustive survey on video-based rendering (VBR) techniques. The need for a controlled capture setting is shared by many methods. Zitnick et al [8] use a specific system of 8 cameras combined with segmentation based stereo to extract the geometry. Similarly Wilburn et al. [34] use an array of 100 tightly-packed cameras. Broxton et al. [35] describe a custom camera array of 46 synchronized cameras mounted on a dome used to capture 6 degrees of freedom (DoF) wide-baseline light field videos. Guo et al. [36] relight videos with a set up of 331 light sources and 90 cameras, while Collet et al. [37] require 106 cameras. In a less constrained way, Pozo et al. [38] create a 16-camera rig to reconstruct 360 panoramic videos and synthesize new views. Penner and Zhang [39] use a soft volumetric representation for narrow baseline IBR to enforce smooth reconstructions. This method can handle motion, but has trouble handling unstructured data and works best from camera arrays. Our method also works with handheld cameras.

Casually-captured videos have also been considered. Ballan et al. [3] allow for quick transitions between handheld video sequences. Their method segments a single dynamic foreground subject approximated by a planar proxy, and creates a 3D reconstructed static background. To cope with dynamic background objects reprojecting incorrectly, the method blurs background transitions between captured viewpoints. Our method assumes no segmentation nor planarity assumptions for dynamic objects. Lipski et al. [10] use dense correspondence fields to interpolate views between videos. They disambiguate matches in difficult cases by manually drawing correspondence lines on image pairs to use as priors in their matching algorithm. Mustafa et al. [11, 12] reconstruct isolated moving objects after segmenting them out from the initial video. These methods focus on specific object meshes, and so do not provide re-rendering of an entire scene from a novel viewpoint.

Recently, Luo et al. [40] introduced a consistency term by fine tuning a neural network to improve the estimated depth per point. This works for a single camera with no or limited dynamic motion. Bansal et al. [6] use foreground and background extraction together with a self-supervised CNN based composition operator, and Yoon et al. [41] use deep learning to extrapolate new views from a single monocular video camera; we compare our approach to this method. While deep learning techniques have shown progress in this area, and neural methods show significant gains, they still often need hours of computation for a single static scenes (e.g., NeRF [26]) or use voxel-based acceleration structures that are not obvious to extend to video due to memory constraints (e.g., DirectVoxGo [42] or InstantNGP [43]).

Outside of NVS, other video reconstruction tasks raise consistency questions. Vo et al. [9] used a spatio-temporal bundle adjustment technique and human motion priors to reconstruct actor performances by temporally aligning videos at sub-frame

precision. Bao et al. [44] used deep learning for consistent video super resolution. Finally, Davis et al. [45] recovered depth in dynamic scenes by unifying structured light and laser scanning into a space-time stereo framework.

Given how challenging consistency can be for dynamic scenes with just RGB cameras, our approach considers how to defer temporal consistency from the reconstruction step to the novel view synthesis step.

### 3. Method

Our algorithm takes as input a set of casually-captured synchronized videos. We also provide the focal lengths for a pair of cameras (required by OpenMVG [46]), while the remaining focal lengths are estimated automatically by our algorithm. Our method proceeds in two steps (Figure 1):

1. **Camera pose estimation and 3D scene points.** We perform a three-step structure from motion reconstruction to provide both the set of camera poses and a set of sparse 3D points for each time step (Section 3.1).
2. **Novel depth and novel view rendering.** We densify the sparse points into a depth map and render a new virtual camera frame by optimizing a coarse-to-fine variational formulation while enforcing spatio-temporal consistency (Section 3.2).

#### 3.1. Camera pose estimation and 3D scene points.

Let us consider a set of  $S$  synchronized video views of a dynamic scene, each composed of  $T$  frames. We call  $I = \{I_{s,t} | s = 1, \dots, S; t = 1, \dots, T\}$  the set of all frames indexed by  $s$  (camera index) and  $t$  (time step). At each frame, via SfM, we recover the camera parameters  $\mathbf{C}_{s,t}$  consisting of the intrinsic matrix and extrinsic rotation and translation matrices, and a set of sparse 3D points for each time step. First, we efficiently recover a set of camera poses for all frames. In contrast to other methods [3, 11, 12], we estimate poses without an explicit dynamic object segmentation step. Second, we recover 3D points by solving a per-timestep SfM problem without a complex temporal reconstruction. We solve each SfM problem with an *a contrario* algorithm [47]. This automatically adapts thresholds to the input data instead of using global thresholds, which is more flexible to different inputs.

*Efficient camera pose estimation.* A straightforward approach for accurate SfM is to solve a problem across all frames simultaneously, but this can be expensive and memory prohibitive. A second approach might consider solving only between consecutive time steps, but this is known to produce camera position drift [48]. Instead, we take a coarse-to-fine approach.

We begin by computing SfM across keyframes at every  $\kappa$  time steps of each video. We detect and match SIFT keypoints [49] within this subset and then simultaneously solve for all camera poses and 3D points. Then, we refine our estimate with a second SfM that only matches keypoints between successive frames of the same camera view, with previously-estimated camera poses held fixed. This considers every frame of every video, but we

only match  $I_{s,t}$  to  $I_{s,t+1}$ , and not to  $I_{s+1,t}$  or  $I_{s+1,t+1}$ . To recover smooth camera paths per view, we add two additional penalty terms to the bundle adjustment:

$$w(t-t') \|C_{s,t} - C_{s,t'}\|^2, \quad t-3 \leq t' \leq t+3 \quad (1)$$

and

$$w(t-t') \|A_{s,t} - A_{s,t'}\|^2, \quad t-3 \leq t' \leq t+3, \quad (2)$$

where  $w(t-t')$  is a Gaussian weight function with a standard deviation of 1.16,  $C_{s,t}$  is the center of each camera pose, and  $A_{s,t}$  is the angle-axis representation of the rotation matrix  $R_{s,t}$ , where the angle is expressed in radians. This second SfM reduces computation time over all-pairs matching while still reducing drift by constraining the frame-to-frame pose estimates by the keyframe pose estimates. For hyperparameters, smaller  $\kappa$  will increase processing time, while larger  $\kappa$  may make it more difficult to match fast camera motion. We found  $\kappa = 20$  to be a good compromise in our test sequences.

**3D scene points.** To recover 3D points across the scene, we solve a keypoint reconstruction problem that is independent per time step. Taking as fixed the recovered camera poses for each video frame, we match 2D keypoints only between frames with the *same timestamp*, then reconstruct a set of sparse 3D points per time step. This is our key to handling dynamic scenes: Motion often makes it difficult to match dynamic objects over time, but as 2D keypoints are *not* matched in time in this last SfM step, moving objects are correctly recovered at least in space. This knowingly produces temporal inconsistencies; we will recover from these errors in the novel view synthesis stage where it is easier to enforce consistency (Sec. 3.2).

**Post processing.** Finally, we increase the density of our point matches using PatchMatch [50], as proposed in the OpenMVS<sup>1</sup> and COLMAP<sup>2</sup> [51] frameworks. This process splats points to each view and assigns colors to the 3D point cloud.

### 3.2. Novel depth and novel view rendering

Our SfM recovers a set of camera poses and an RGB 3D point cloud per time step. However, at this stage of our algorithm, projecting these points to a novel view still leaves large regions of empty space. To synthesize more realistic views, we diffuse these points in depth and RGB in the new view in image space while enforcing spatio-temporal constraints.

**Notation.** We will often warp the content of a frame  $I_{s,t}$  into the domain of the novel view  $I_t$ : this reprojection is computed using the extrinsic and intrinsic parameters of both reprojected frames and virtual camera, as well as the depth map  $D_t$ . We denote the projected frame as  $I_{s,t}^{proj}(x) = I_{s,t}(C_{s,t}C_t^{-1}(x, D_t(x)))$ , where  $C^{-1}(x, D_t(x))$  is the image plane to world coordinate system transformation of the pixel location  $x$  given its depth value  $d$ . We also denote a sparse map by  $\hat{\cdot}$ . The sparse depth map obtained by projecting the sparse point cloud into frame  $t$  of the new virtual camera path is then  $\hat{D}_t$  and its corresponding sparse color image is  $\hat{I}_t$ .

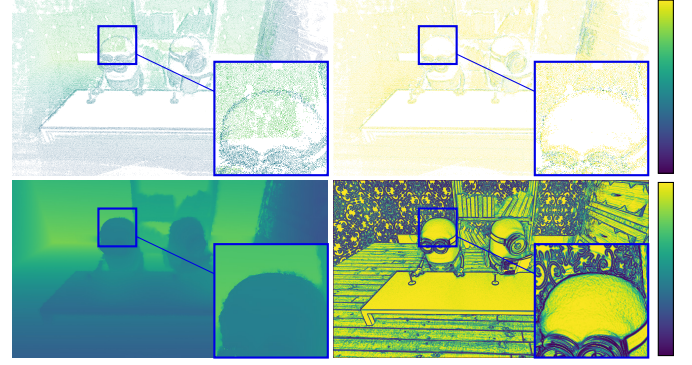


Fig. 2: **3D point weights preserve color edges and occlusions.** *Top row:* Sparse reconstructed 3D points (left) and their weights  $w_{\hat{D}}$  (right) projected into the virtual view. White indicates areas of empty space; depth map is bright green in far depth regions. *Bottom row:* Points diffused into a full depth map  $D$  (left) according to the weight map  $w_D$  (right). Note how the color edges are correctly identified via Eq. 5, and how the occluded points from behind the head of the character on the left are given no weight by Eq. 6 (top right) and so do not corrupt the depth.

**Algorithm progression.** We wish to warp a frame  $I_{s,t}$  to the novel view  $I_t$  to be blended into a final novel view. For this, we need both the estimated camera poses and the dense depth maps  $D_t$ , which are yet to be computed. But, to properly constrain the diffusion of the sparse depth values  $\hat{D}_t$ , recovered in Sec. 3.1, we need RGB information from the virtual camera’s point of view. Thus, we jointly solve for the depth maps  $D_t$  and color images  $I_t$  by minimizing the energy functional:

$$E(D_t, I_t) = E_D(D_t) + E_I(I_t). \quad (3)$$

The functional relates terms constraining the depth map ( $E_D$ ) to terms constraining the color image ( $E_I$ ) by weights that guide the diffusion process. We solve  $E$  iteratively: we first solve for the depth map  $D_t$  while fixing the color values  $I_t$ , and then conversely we fix the depth values and solve for color. This avoids having to solve a nonlinear system of equations, and lets us use slightly-improved depth values to warp the input frames at each step. This improves the estimate of the rendered RGB image, which in turn constrains the diffusion of the depth.

**Depth diffusion.** We project the sparse point cloud into the novel view image plane approximated to integer pixel locations, creating the sparse depth map  $\hat{D}_t$  as an initialization. When several sparse points project onto the same pixel, we keep only the depth value corresponding to the sparse point that is closest to the camera, disregarding possible outliers. Then, we densify the sparse depth map by minimizing the following energy:

$$E_D(D_t) = \int_{x \in \Omega} w_D(x, t) |\nabla D_t(x)|^2 dx + \lambda_{PC} \int_{x \in \Omega} w_D(x, t) |D_t(x) - \hat{D}_t(x)|^2 dx. \quad (4)$$

The first integral is a smoothness term controlled by weight  $w_D$ . We wish diffusion to decrease around color edges to produce sharp results. We also wish diffusion of depth values to

<sup>1</sup><https://github.com/cdcseacave/openMVS>

<sup>2</sup><https://colmap.github.io/>



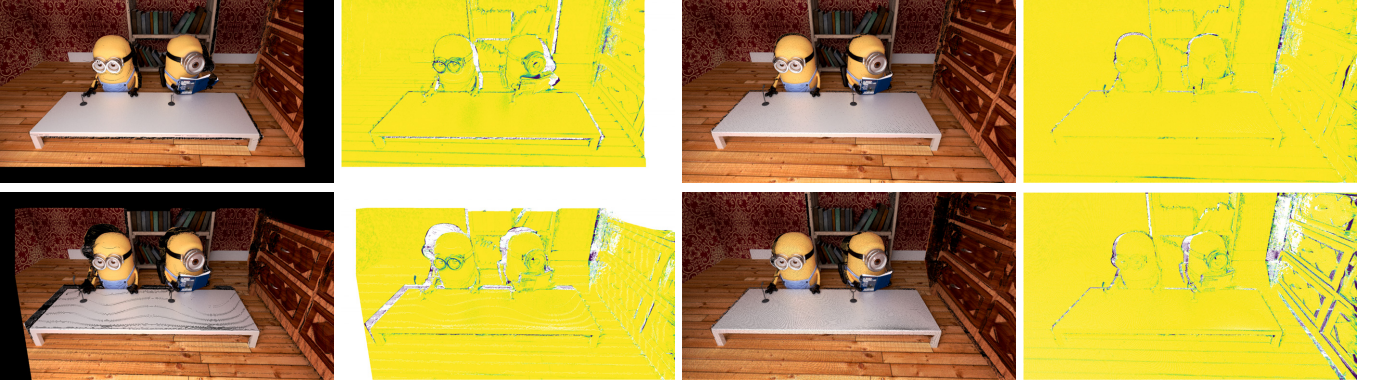


Fig. 3: **Input image contributions.** Four closest input images  $I_{s,t}$  projected onto the virtual camera's view point alongside their corresponding weight maps  $w_P$  (Eq. 8). We see that occluded regions are given little weight.

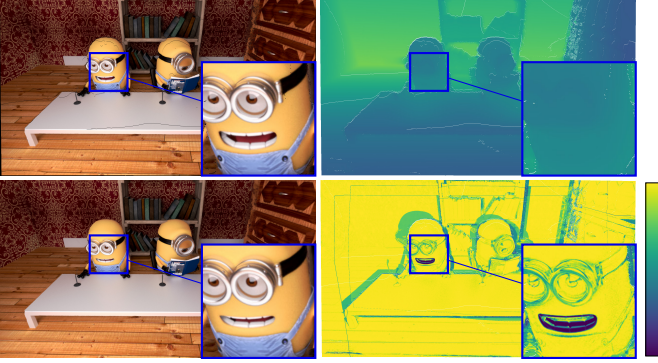


Fig. 4: **Effect of temporal weight.** Top row: Color image  $I_{t-1}$  and depth map  $D_{t-1}$  of a previous time step warped to the camera view at time step  $t$ . Bottom left: Current color image  $I_t$ . Bottom right: Weight map  $w_T$  (Eq. 11) modulates consistency, notably around the moving mouth of the character on the left.

increase when the colors from reprojected input views are similar, as they are likely correct. As such, we define  $w_D$  as:

$$w_D(x, t) = \frac{1}{\|\nabla I_t(x)\|^2 \sum_{s=1}^n \sigma_{vis}^s(x, t)} \sum_{s=1}^n w_P^s(x, t), \quad (5)$$

where  $1/\|\nabla I_t\|^2$  modulates depth diffusion around color edges, and  $1/\sum_{s=1}^n \sigma_{vis}^s(x, t)$  is a normalization factor that accounts for each pixel's visibility in the novel view. As both the visibility term  $\sigma_{vis}$  and the projection weight  $w_P^s$  pertain more to the color diffusion process, we will define them later on in Eq. 8.

The second integral in Eq. 4 reduces the weight of sparse 3D points that are occluded from the point of view of the virtual camera or are erroneously reconstructed. For this, we relax the constraint of  $D_t$  where it exactly matches the projected sparse point cloud:

$$w_{\hat{D}}(x, t) = \exp\left(-\frac{\|\hat{I}_t(x) - I_t(x)\|^2}{2\sigma^2}\right). \quad (6)$$

In Figure 2, we show example weight maps  $w_{\hat{D}}$  and  $w_D$  that govern the depth diffusion process, as defined in Eqs. 5 and 6.

There are two parameters in this diffusion process: 1)  $\sigma$  controls the soft occlusion tolerance, and we set  $\sigma = 0.075$  in all our experiments; 2) the sparse point cloud attachment weight

$\lambda_{PC}$  controls the influence of the reconstructed points, and we set it in the range  $\lambda_{PC} \in [0.25, 2]$ .

**Color diffusion.** Given depth map  $D_t$ , we initialize the RGB image to a projection of the color in the input point cloud. Then, we densify it by minimizing the following diffusion energy on each color channel independently:

$$\begin{aligned} E_I(I_t) = & \int_{x \in \Omega} \|\nabla I_t(x)\|^2 \\ & + \sum_{s=1}^n \int_{x \in \Omega} \lambda_P w_P^s(x, t) \|I_t(x) - I_{s,t}^{proj}(x)\|^2 dx \\ & + \sum_{s=1}^n \int_{x \in \Omega} \lambda_G w_P^s(x, t) \|\nabla I_t(x) - \nabla I_{s,t}^{proj}(x)\|^2 dx \end{aligned} \quad (7)$$

The first integral encourages smooth gradients over the intensity of the novel view, which aids blending of the projected input images especially along their borders. The second integral constrains the RGB intensities and  $I_t$  to be close to the intensities of  $I_{s,t}^{proj}$ , and the third integral constrains the RGB gradients similarly. They are both modulated by the weight

$$w_P^s(x, t) = \sigma_{vis}^s(x, t) \exp\left(-\frac{\|I_{s,t}^{proj}(x) - I_t(x)\|^2}{2\sigma^2}\right), \quad (8)$$

that measures the agreement of each warped input frame with the novel view and is held constant at each iteration. Figure 3 shows warped input frames along with their weight maps  $w_P^s$ .

$w_P$  incorporates visibility term  $\sigma_{vis}^s(x, t)$  (Fig. 5) that is 1 for a given pixel  $x$  of the novel view  $I_t$  only if, out of every pixel that is projected to the same pixel location in an input image  $I_{s,t}$ ,  $x$  has the smallest depth value  $d$  in the input image's coordinate frame.

Color diffusion relies on two new parameters:  $\lambda_P$

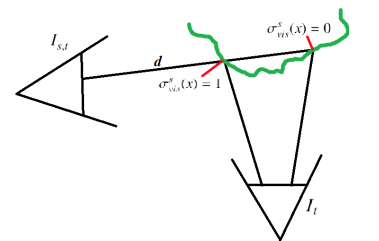


Fig. 5: Visibility term  $\sigma_{vis}^s$  determines whether a projected input pixel is seen in the novel view.

controls the influence of the data constraint and  $\lambda_G$  controls the influence of the gradient equality constraints. We set them both in the range [5, 20].  $\sigma$  serves the same function and values as in the depth map diffusion.

### 3.3. Temporal consistency

We enforce temporal consistency within novel views by additional terms in  $E_D$  and  $E_I$ . With slight abuse of notation, this term is added to Eqs. 4 and 7:

$$E_D(D_t) = \dots + \lambda_T \int_{x \in \Omega} w_T(x, t) \|D_t(x) - D_{t-1}^{proj}(x)\|^2 dx, \quad (9)$$

$$E_I(I_t) = \dots + \lambda_T \int_{x \in \Omega} w_T(x, t) \|I_t(x) - I_{t-1}^{proj}(x)\|^2 dx. \quad (10)$$

These terms constrain depth  $D_t$  to remain similar to the warped previous depth  $D_{t-1}$  projected to the current camera location, and for color  $I_t$  similarly. This constraint is relaxed by a weight

$$w_T(x, t) = \frac{1}{n} \sum_{s=1}^n \exp \left( -\frac{\|I_{t-1}^{proj}(x) - I_{s,t}^{proj}(x)\|^2}{2\sigma^2} \right) \quad (11)$$

for pixels for which an agreement in color was not reached. This is expected in regions containing motion because the depth values of frame  $t-1$  may be invalid, as is the case around the mouth of the character on the left in Figure 4.  $w_T$  allows the computation of depth and color values of these pixels to rely more freely on the other terms of the functional, like the data term of the depth or the color of the warped input images.

The parameter controlling the strength of the temporal consistency  $\lambda_T$  is set in the range of [0.01, 0.1].

### 3.4. Implementation details

To avoid using input frames that are far away from the novel camera’s view, we rank each input camera based on its distance from the novel camera according to the following formula:

$$r_F(s) = \frac{1}{\|C_t - C_{s,t}\|^2} \exp \left( -\frac{\arccos((\text{tr}(R_t R_{s,t}^T) - 1)/2)}{2\pi\sigma^2} \right) \quad (12)$$

This penalizes frames with camera poses that are either far in position or in viewing direction from the novel view. Then, we use the first  $n = 4$  ranked input frames to minimize the functional of Eq. 3.

We approximate every partial derivative as central differences and use Dirichlet boundary conditions.

For efficiency, we also proceed in a multiscale fashion: we solve for depth and color at a coarse resolution  $l$ , and then use these to initialize a finer resolution  $l-1$  at twice the previous resolution—our lowest level is 1/64 of the original frame size. At each level we iterate  $10 \cdot 2^l$  times. Finally, we also proceed in a streaming manner: we reproject the previous frame’s depth and color (denoted as  $D_{t-1}^{proj}$  and  $I_{t-1}^{proj}$ ) into the current virtual camera pose for use within the temporal consistency constraint.

## 4. Experiments and results

### 4.1. Dataset sequences

*Real-world existing dataset.* We exploit existing datasets used in the context of novel view synthesis, all of them captured using camera arrays:

- *Jumping* [41]: A group of four people jump (12 cameras).
- *Skating* [41]: A person rides a skateboard (12 cameras).
- *Playground* [41]: A person flies a dinosaur balloon (12 cameras).
- *Umbrella* [41]: A person opens and rotates an umbrella (12 cameras).
- *DynamicFace* [41]: A person of making faces (12 cameras).
- *Breakdancers* [8]: A person break dancing in front of 4 people (8 cameras).

*Custom dataset.* We test our algorithm on three 100-frame real world sequences that we acquired each with five Canon Rebel EOS T7i cameras at  $1920 \times 1080$  resolution. The cameras were set up with a mix of hand held or tripod capture. The videos were synchronized based on the audio track. Our sequences are:

- *Cat and Dog*: Two pet animatronics.
- *Elephant Wiggle*: A puppet hanging by a wire.
- *Drone*: A drone hanging by a wire.

We additionally generate a synthetic 100 frames long *Minions* sequence using 11 input cameras to compare to ground truth RGB and depth estimation from a 12th camera. It contains a rendering of two characters laughing behind a table. In this sequence, all cameras are moving.

### 4.2. Metrics

To evaluate reconstruction results, we compare camera pose position and orientation error in world space, and 3D point reprojection error in pixels. To evaluate novel view synthesis results, we quantitatively compare methods in terms of PSNR (higher is better), more perceptually-motivated SSIM [52] (higher is better), and video temporal consistency measures SRRED and TRRED [53] (lower is better).

### 4.3. Results and ablation study

To show results in this paper, we extract frames from output videos to highlight the comparisons; please see our accompanying video to better evaluate the results and comparisons. To begin with our method, Figure 6 shows rendered frames from novel views and corresponding depth maps for the *Cat and Dog* and the *Elephant Wiggle* sequences. While some artifacts remain in the depth video, the generation of the final novel view RGB rendered sequence is robust to these and has fewer artifacts. Note that the borders of the view partially appear blurry when there is insufficient field of view overlap between input videos.

We ablate our SfM method using the synthetic *Minions* dataset with moving objects, where points are known to be either static or dynamic (Table 1). We compared the recovered pose over 50

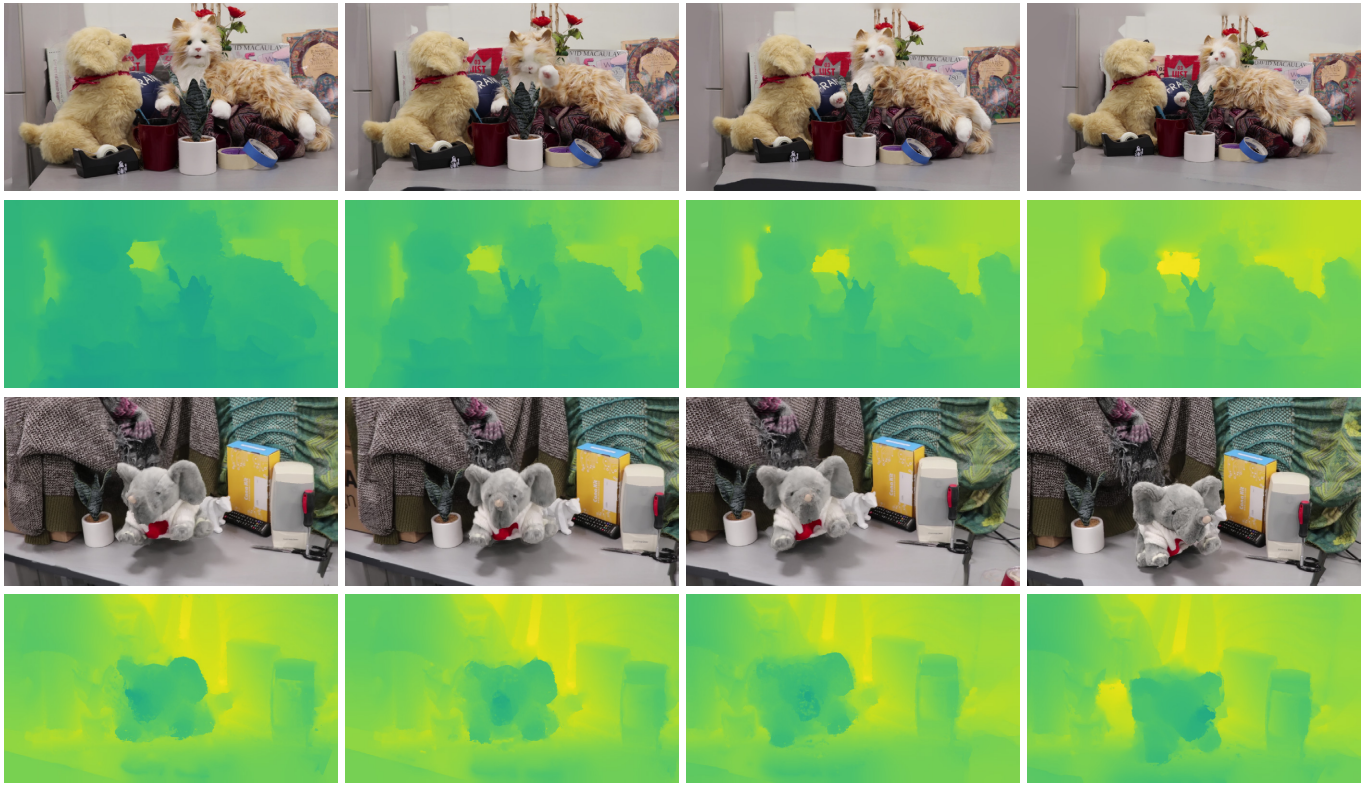


Fig. 6: **Results.** Color and depth outputs for *Cat and Dog* and *Elephant Wiggle* scenes.

Method	Pos. error (mm)	Orient. error (°)	Median reproj. error (pix)	Mean reproj. error (pix)
Naive SfM	0.0016	0.1088	0.0553	0.0972
Naive SfM from static objects only	0.0017	0.1482	0.0512	0.0901
Our SfM without path smoothing	0.0016	0.0675	0.0513	0.0940
Our SfM	0.0015	0.0740	0.0511	0.0937
SfM with Ground Truth Poses	0	0	0.0552	0.0975

Table 1: **Ablation study—quantitative camera pose and 3D point reconstruction.** On the synthetic *Minions* scene, we compare estimated camera pose accuracy for naive SfM, naive SfM using a ground truth mask for the static parts of the scene, and an ablated version of our space time SfM without camera smoothing. While adding smoothing slightly increases the orientation error, it reduces positional errors. Our approach also minimizes the median reprojection error of the feature points. SfM minimizes reprojection error by construction: using ground truth poses results in zero position and orientation error, but can still lead to inaccurate point locations and so increased reprojection error.



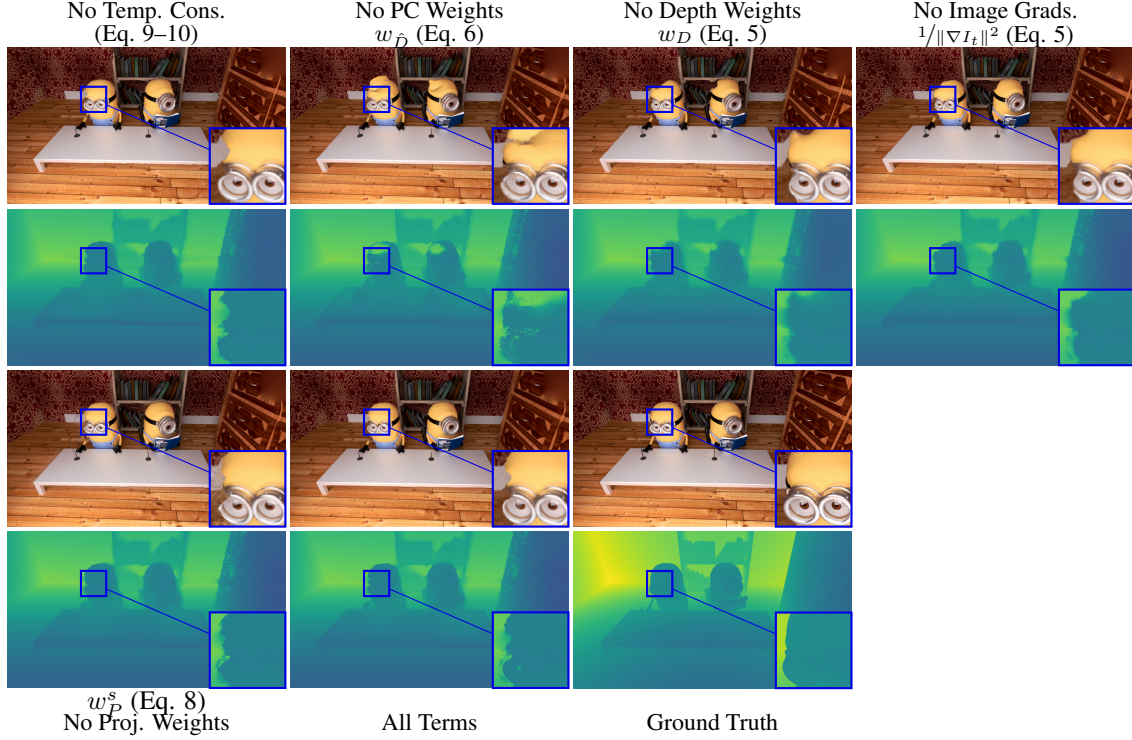


Fig. 7: **Ablation study—qualitative novel views.** Our rendered result without temporal consistency (1st column), without weights on the projected sparse 3D points (2nd column), without depth weights and inverse image gradients (3rd column), only without inverse image gradients (4th column), without weights on the projected input images (5th column), together with the full result (6th column) and the ground truth (7th column), for both the depth (first row) and the color view (2nd row).

Metric	PSNR↑	SSIM↑	SRRED↓	TRRED↓
No Temp. Cons.	26.878	0.846	11.82	26.60
No PC Weights	26.875	0.848	11.65	19.62
No Depth Weights	25.256	0.833	15.62	32.60
No Image Grads.	26.207	0.847	12.88	21.34
No Proj. Weights	27.302	0.846	11.30	19.20
<b>All Terms</b>	<b>27.427</b>	<b>0.853</b>	<b>10.35</b>	<b>17.77</b>

Table 2: **Ablation study—quantitative novel views.** Evaluated by PSNR, perceptual quality measure SSIM [52], and video consistency measures SRRED/TRRED [53]. For PSNR and SSIM, higher values are better; for SRRED/TRRED, lower values are better. All terms help improve quality.

timestamps and 5 cameras. First, we compare against a naive SfM approach that solves for all frames simultaneously without consideration of dynamic objects. Next, some methods rely on segmenting out moving objects to cope with dynamic scenes [11, 12]. To compare to this idea, we created a segmentation-based SfM baseline from the naive SfM by performing reconstruction only from points that are known to be static using perfect ground truth masks. While the segmentation slightly improves the 3D reconstruction, its positive effect is not clear on the recovery of camera positions, even though the dynamic object segmentation is a pixel accurate ground truth. Against both baselines, our method makes better use of dynamic points to more accurately recover camera paths and its median reprojection error is the smallest even compared to the reconstruction of static points only. Finally, we compare against the non-smoothed camera path version of our approach. While the positional error decreases, the

rotation error slightly increases. Our camera path regularization also reduced reprojection errors. Overall, we found smoothing to provide better final results.

We also test our rendering method in an ablation study by disabling the temporal consistency term and the weight functions in the diffusion process one at a time. The results are shown in Figure 7 and Table 2. For all metrics, our full model achieves the best performance. The effect of the temporal consistency term on our rendering process can be best seen in our accompanying video, where the rendering without it results in jitters mostly along object boundaries. If we forego the weights  $w_D$  (Eq. 6) on the projected sparse 3D points, the background points that should be occluded by foreground objects gain out-sized influence over the depth diffusion process and ultimately over the color rendering. The depth weights  $w_D$  (Eq. 5) boost the propagation of high-confidence depth values and reduce the influence of incorrect ones and so prevent oversmoothing of the depth maps. The inverse image gradients applied to modulate the depth diffusion give additional sharpness to the depth images along object boundaries. Finally, the weights  $w_P^s$  (Eq. 8) on the projected input images prevent ghosting artifacts caused by occlusions and erroneous reprojections.

#### 4.4. Baseline method comparisons

We compare our method to four recent methods, including deep-learning-based methods that require external training databases: Deep Blending [5], Local Light Field Fusion [29], Extreme View Synthesis [4], and MonoCam [41]. Furthermore, we use the *Breakdancers* scene to compare to results provided by



Method	Scenes	Min. views	Training time	Preprocess per frame	Render time	Fig. #
Deep Blending [5]	Static	4	37 h	8 h	Real time	8
LLFF [29]	Static	6	?	10 min	Real time	11
EVS [4]	Static	2	?	10 min	98 sec	10
MonoCam [41]	Dynamic	1	?	?	?	12
VI [8]	Dynamic	8	N/A	?	Real time	13
VVC [10]	Dynamic	5	N/A	<i>Manual fix</i>	Real time	13
<b>Ours</b>	<b>Dynamic</b>	<b>4</b>	<b>N/A</b>	<b>2 min</b>	<b>6.8 sec</b>	<b>7–11</b>

Table 3: **Related work comparisons.** Considering scene type, minimum number of input views, and speed. LLFF and EVS use pre-trained networks, so we did not re-train them. For MonoCam, we used the author’s results for comparison. ‘?’ denotes where no information is available.

Method	PSNR $\uparrow$	SSIM $\uparrow$	SRRED $\downarrow$	TRRED $\downarrow$
<b>Ours</b>	26.22 (25.59)	0.80 (0.79)	7.12 (7.82)	11.13 (10.50)
Deep Blending [5]	9.8 (21.59)	0.34 (0.63)	12.51 (11.87)	57.39 (50.82)
LLFF [29]	21.88 (21.09)	0.57 (0.56)	22.28 (21.67)	35.61 (31.48)

Table 4: **Results—quantitative metrics.** The metrics were computed on the first 34 frames of the synthetic *Minions* scene. In parenthesis, we also compute metrics on the first 12 frames only, after which Deep Blending [5] starts to occasionally produce entirely black frames.



Fig. 8: **Results—qualitative.** Comparison with Deep Blending [5] on the *Cat and Dog* and *Elephant Wiggle* sequences.

two older methods that best match our intended setup: the View Interpolation (VI) method of Zitnick et al. [8] and the Virtual Video Camera (VVC) method of Lipski et al. [10]. Each of these methods work with different numbers of input views and require different amounts of processing time. Some of these methods are only intended for static scenes, and so we would expect them to produce temporally inconsistent results. Table 3 summarizes these properties.

**Static—Deep Blending [5].** We compare our rendering method with Deep Blending [5]. This learns optimized weights to blend four layers of mosaic images, where the first layer is composed of the best fitting pixels, the second layer holds the second best, and so on, with ‘best’ determined by a heuristic. For the comparison, we first reconstructed each scene separately for each time step as described in their method. Afterwards, to use the same camera path as for our results, we registered each time step to our full space-time reconstruction based on

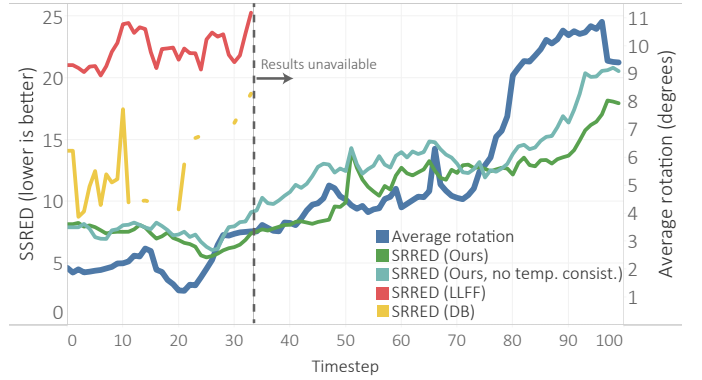


Fig. 9: **Results—quantitative temporal consistency.** We show the evolution of SRRED values (lower is better) over time on the *Minions* scene for our method (green), our method without temporal consistency (cyan), and for Deep Blending (DB) [5] and LLFF [29]. At the same time, we display the average rotation angle between the target view and the four nearest views used for reconstruction (see Sec. 3.4). While our method outperforms DB and LLFF, for this sequence, the novel view rotation with respect to the four nearest views increases over time, which correlates with increasing errors. As expected, our temporal consistency reduces this drift. DB produces several frames with large black areas that result in *NaN* SRRED values, and thus, gaps in plots.

the camera positions. Finally, we used the pre-trained network provided by the authors to render each frame.

Figure 8 shows that Deep Blending cannot always reconstruct marginal parts of the scene, and is often blurrier than our result. Table 4 quantitatively shows that our method outperforms it in every metric. Figure 9 shows that SRRED values correlate with the angular distance between the target view and the four nearest input views used for reconstruction, and that this distance tends to increase over time for our particular test sequence.

**Static—Extreme View Synthesis [4].** Figure 10 shows a comparisons with Extreme View Synthesis (EVS) [4]. As input, EVS receives our SfM results. As expected, it exhibits flickering

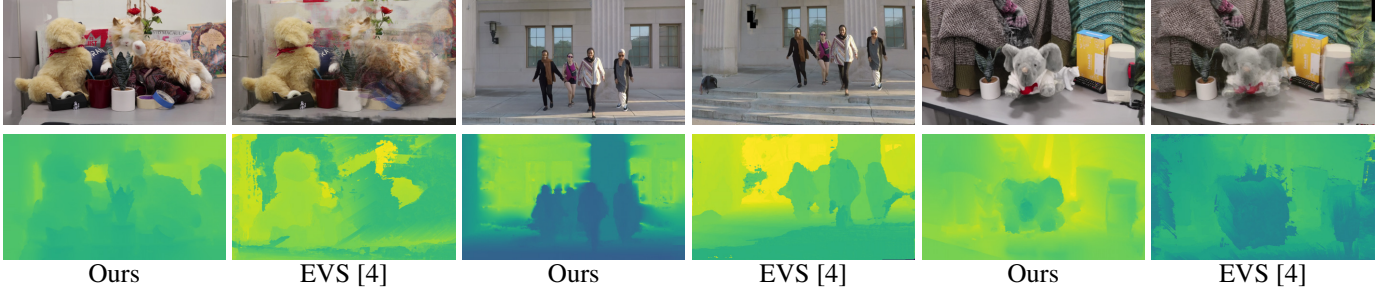


Fig. 10: **Results—qualitative comparison with Extreme View Synthesis (EVS) [4].** On the *Cat and Dog*, *Jumping*, and *Elephant Wiggle* sequences. Our method produces fewer artifacts than EVS.

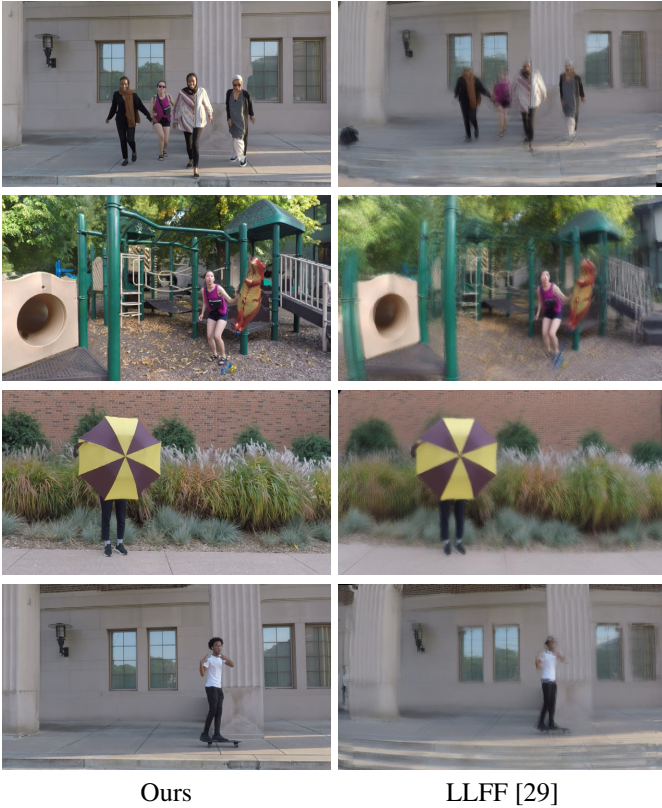


Fig. 11: **Results—qualitative comparison with Local Light Field Fusion (LLFF) [29].** On the *Jumping*, *Playground*, *Umbrella*, and *Skating* sequences.

since this method is designed for static scenes and does not enforce temporal consistency. In addition, EVS cannot handle high resolution input because of its memory requirements, so we lowered the resolution of the input video from  $1920 \times 1080$  to  $1280 \times 720$ . Similarly, we could not increase the depth resolution of its scene reconstruction step, which leads to inaccurate depth maps and thus severe ghosting in the affected areas.

**Static—Local Light Field Fusion [29].** Figure 11 shows a comparison with Local Light Field Fusion (LLFF) [29]. Since LLFF requires at least 6 cameras to work, we could only compare on the 12-camera dataset sequences. To improve stability, we provide LLFF with our temporally consistent reconstruction results instead of their COLMAP reconstruction. Since our 3D reconstruction remains noisy by design, which is not expected by this method, we fixed the minimum and maximum depths to

known correct values—this lets LLFF correctly set space bounds for its Multi-Plane Image computations. Without this, extreme flickering occurs; with, flickering is reduced but is not eliminated completely. Our result also appears sharper and with reduced ghosting artifacts, which is also reflected in Table 4.

**Dynamic—Monocam [41].** We conduct a comparison with Monocam [41], using the results given by the authors. It is important to note that the input sequences provided by the authors as a dataset differ slightly from the sequences used as results in the corresponding paper [41], which makes direct comparison impossible. For instance, in the *Skating* sequence, the skater is performing hand gestures in the input sequence, but not in the result sequence.

In the MonoCam results, Figure 12 shows that dynamic background objects like the plants in the *Umbrella* sequence incorrectly appear static if the virtual camera is static, and are inconsistent if the virtual camera is dynamic. MonoCam results also exhibit temporal coherence artifacts. For instance, the reflections in the jumping and skating sequences jump back and forth based on which view was used to render them. Please see these artifacts in the accompanying video.

**Dynamic—View Interpolation, Virtual Video Camera [8, 10].** Figure 13 shows a comparison with View Interpolation (VI) [8] and Virtual Video Camera (VVC) [10] methods. We approximately reproduced the camera path of the video provided by the authors for the Breakdancers scene for this comparison. While VI requires a fixed and calibrated camera grid, our method can use hand-held devices. VVC eliminates these restrictions, but relies on user input to correct correspondence matches—a user-interaction that our methods avoids.

#### 4.5. Challenging sequence—Drone

This sequence shows a quadcopter drone (Figure 14). The drone has many thin features: the chassis, the fan blades, and exposed wires between battery and motors. Here, our SfM reconstruction fails to find feature points on or nearby thin features at the correct depth, therefore our consistent propagation cannot provide the correct depth. As such, we see ghosting effects.

#### 4.6. Computational resources

We implemented our system in C++ on a Intel(R) Xeon(R) CPU ES-2630 v3 @2.4GHz computer. We used the OpenMVG library to compute Structure from Motion and sparse depth maps;



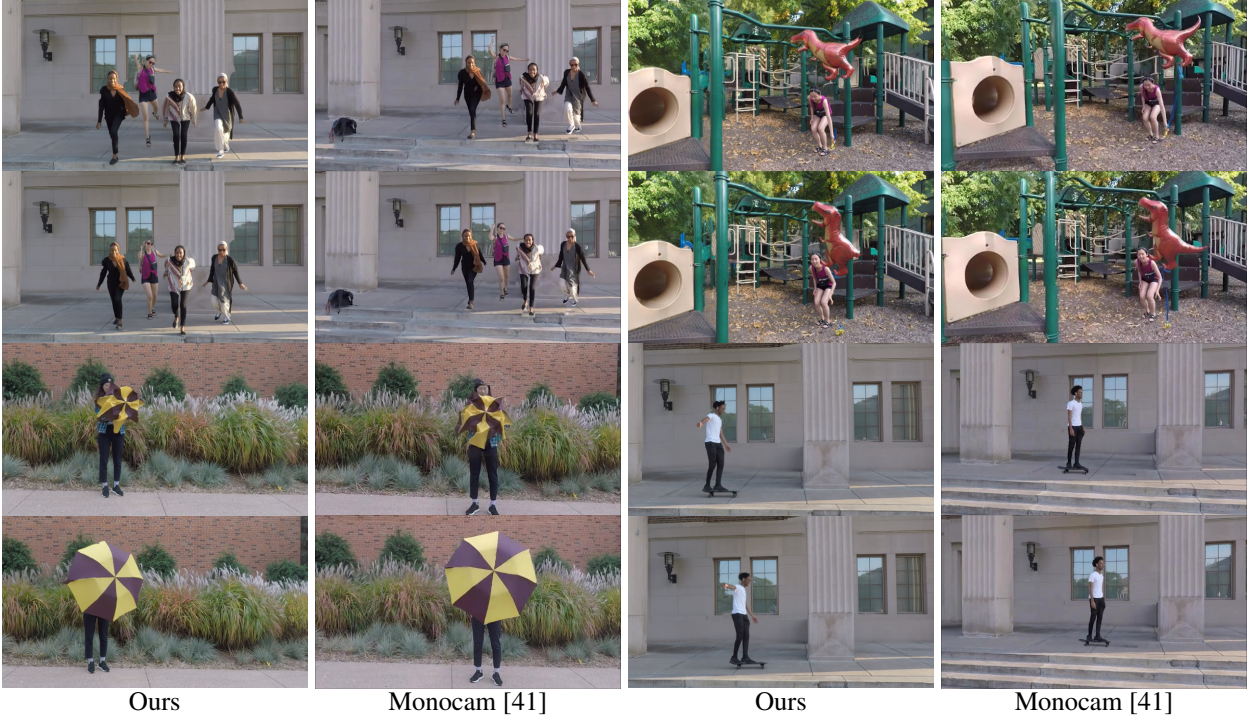


Fig. 12: **Results—qualitative comparison with MonoCam [41].** On the *Jumping*, *Playground*, *Umbrella*, and *Skating* sequences. Stronger temporal inconsistencies in results computed with MonoCam can be seen in the accompanying video. Camera poses and small scene details differ between the input data and MonoCam results, but were provided as is by the authors.



Fig. 13: **Results—qualitative comparison with VI [8] and VVC [10].** On the *Breakdancers* sequence, using results provided by the authors.

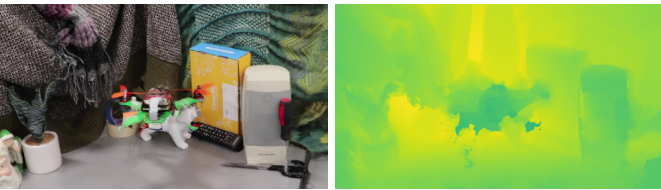


Fig. 14: **Limitations on thin features.** In the *Drone* sequence, the scene has thin features which makes geometry reconstruction difficult. Here, our consistent propagation has trouble correcting for missing sparse feature points.

and both OpenMVS and COLMAP to compute the PatchMatch-based sparse depth map post processing (Section 3.1). We parallelize the code using OpenMP and run on 32 cores; the rendering algorithm loads up to 2GBs of data per frame. As an example of wall-clock time, it took 2.2 hours to process the *Elephant Wiggle* sequence (5 cameras, 100 Full HD 1080p frames per camera). The computation time breaks down to camera and sparse depth estimation (2 hours), and the rendering itself (6.8s per frame at half of the input resolution, 11.3 minutes for the whole video).

## 5. Discussion

Our approach solves a joint optimization for color and depth via alternating descent, where each of color and depth is optimized in turn via Eq. 3 through the coarse-to-fine iterations. Solving for both color and depth simultaneously struggled to converge under the same scheme, producing poorer results.

Concerning limitations, first, our choice of the OpenMVG library [46] for computing the SfM has the drawback that we must provide focal lengths for a pair of cameras to initiate the reconstruction process. Our method also requires that the video sequences should have enough texture on the objects and in the background such that enough SIFT keypoints can be detected and matched. This also relates to the amount of motion in the frame: conceptually, if SIFT keypoints are only detected on moving objects, then camera pose estimation will fail. In practice, we did not find this to be a problem. Furthermore, if the baseline is too wide, then not enough points will be obtained on moving objects and the depth propagation will fail. Similarly, in the case of objects not properly supported by 3D points, our propagation technique will produce ghosting artifacts around depth discontinuities. This effect is most noticeable around (dynamic) foreground object boundaries and the background, where the difference in depth values tends to be larger. Finally, our optimization has parameters that can be tuned for each sequence; we provide reasonable initial values (Sec. 3.2).

Finally, our current implementation is unoptimized C++ running on a CPU. Even if we optimize the implementation, one bottleneck is that keypoints from several images must be matched, and this is time consuming. If we consider SfM as an offline task

to be performed once per scene, then the view rendering part currently takes 7 seconds per frame. Given the fixed grid, GPU-based diffusion optimizers are possible, which would produce a much more application-friendly render time.

## 6. Conclusion

We introduce a novel view synthesis method which can handle dynamic scenes. It is based around the key insight that reconstructing temporally-consistent 3D points on dynamic objects is hard, yet a structure-from-motion reconstruction method need not be temporally consistent if temporal consistency can be enforced in the rendering algorithm. We show that this can be accomplished by deferring consistency to a variational screen-space formulation, which makes it easy to robustly enforce spatio-temporal consistency via reprojection constraints weighted by confidences. While our setting has some restrictions, we show competitive results against existing baselines for video-based rendering without resorting to powerful but time-consuming deep learning techniques. In the future, we hope to reduce constraints in camera motions and time with asynchronous videos.

## Acknowledgments

This work was funded in part by ANR project CALiTrOp (ANR-16-CE33-0026).

## References

- [1] Zhang, C, Chen, T. A survey on image-based rendering—representation, sampling and compression. *Signal Processing: Image Communication* 2004;19(1):1–28.
- [2] dos Anjos, RK, Pereira, J, Gaspar, J. A navigation paradigm driven classification for video-based rendering techniques. *Computers & Graphics* 2018;77:205–216.
- [3] Ballan, L, Brostow, GJ, Puwein, J, Pollefeys, M. Unstructured video-based rendering: Interactive exploration of casually captured videos. *ACM Transactions on Graphics (TOG)* 2010;29(4):87.
- [4] Choi, I, Gallo, O, Troccoli, A, Kim, MH, Kautz, J. Extreme view synthesis. In: *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. 2019, p. 7780–7789.
- [5] Hedman, P, Philip, J, Price, T, Frahm, JM, Drettakis, G, Brostow, G. Deep blending for free-viewpoint image-based rendering. *ACM Trans Graph* 2018;37(6). URL: <https://doi.org/10.1145/3272127.3275084>. doi:10.1145/3272127.3275084.
- [6] Bansal, A, Vo, M, Sheikh, Y, Ramanan, D, Narasimhan, S. 4d visualization of dynamic events from unconstrained multi-view videos. *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 2020;.
- [7] Kopf, J, Langguth, F, Scharstein, D, Szeliski, R, Goesele, M. Image-based rendering in the gradient domain. *ACM Transactions on Graphics (TOG)* 2013;32(6):199.
- [8] Zitnick, CL, Kang, SB, Uyttendaele, M, Winder, S, Szeliski, R. High-quality video view interpolation using a layered representation. *ACM Trans Graph* 2004;23(3):600–608. URL: <https://doi.org/10.1145/1015706.1015766>.
- [9] Vo, M, Narasimhan, SG, Sheikh, Y. Spatiotemporal bundle adjustment for dynamic 3d reconstruction. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016, p. 1710–1718.
- [10] Lipski, C, Linz, C, Berger, K, Sellent, A, Magnor, M. Virtual video camera: Image-based viewpoint navigation through space and time. *Computer Graphics Forum* 2010;29(8):2555–2568.
- [11] Mustafa, A, Kim, H, Guillemaut, J, Hilton, A. Temporally coherent 4d reconstruction of complex dynamic scenes. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, p. 4660–4669. doi:10.1109/CVPR.2016.504.
- [12] Mustafa, A, Volino, M, Kim, H, Guillemaut, J, Hilton, A. Temporally coherent general dynamic scene reconstruction. *CoRR* 2019;abs/1907.08195. URL: <http://arxiv.org/abs/1907.08195>. arXiv:1907.08195.
- [13] Holynski, A, Kopf, J. Fast depth densification for occlusion-aware augmented reality. In: *SIGGRAPH Asia 2018 Technical Papers*. ACM; 2018, p. 194.
- [14] Chen, SE, Williams, L. View interpolation for image synthesis. In: *Proceedings of the 20th annual conference on Computer graphics and interactive techniques*. ACM; 1993, p. 279–288.
- [15] Gortler, SJ, Grzeszczuk, R, Szeliski, R, Cohen, MF. The lumigraph. In: *Siggraph*; vol. 96. 1996, p. 43–54.
- [16] Flynn, J, Broxton, M, Debevec, P, DuVall, M, Fyffe, G, Overbeck, R, et al. Deepview: View synthesis with learned gradient descent. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2019, p. 2367–2376.
- [17] Buehler, C, Bosse, M, McMillan, L, Gortler, S, Cohen, M. Unstructured lumigraph rendering. In: *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*. ACM; 2001, p. 425–432.
- [18] Shade, J, Gortler, S, He, Lw, Szeliski, R. Layered depth images. In: *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*. 1998, p. 231–242.
- [19] Debevec, PE, Taylor, CJ, Malik, J. Modeling and rendering architecture from photographs: A hybrid geometry-and image-based approach. In: *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*. 1996, p. 11–20.
- [20] Snavely, N, Seitz, SM, Szeliski, R. Photo tourism: exploring photo collections in 3d. In: *ACM transactions on graphics (TOG)*; vol. 25. ACM; 2006, p. 835–846.
- [21] Hedman, P, Ritschel, T, Drettakis, G, Brostow, G. Scalable inside-out image-based rendering. *ACM Transactions on Graphics (TOG)* 2016;35(6):231.
- [22] Chaurasia, G, Duchene, S, Sorkine-Hornung, O, Drettakis, G. Depth synthesis and local warps for plausible image-based navigation. *ACM Transactions on Graphics (TOG)* 2013;32(3):30.
- [23] Achanta, R, Shaji, A, Smith, K, Lucchi, A, Fua, P, Süsstrunk, S. Slic superpixels. *Tech. Rep.*; 2010.
- [24] Matzen, K, Cohen, MF, Evans, B, Kopf, J, Szeliski, R. Low-cost 360 stereo photography and video capture. *ACM Trans Graph* 2017;36(4).
- [25] Riegler, G, Koltun, V. Free view synthesis. In: *European Conference on Computer Vision*. 2020;.
- [26] Mildenhall, B, Srinivasan, PP, Tancik, M, Barron, JT, Ramamoorthi, R, Ng, R. Nerf: Representing scenes as neural radiance fields for view synthesis. 2020. arXiv:2003.08934.
- [27] Flynn, J, Neulander, I, Philbin, J, Snavely, N. Deepstereo: Learning to predict new views from the world’s imagery. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016, p. 5515–5524.
- [28] Zhou, T, Tucker, R, Flynn, J, Fyffe, G, Snavely, N. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:180509817* 2018;.
- [29] Mildenhall, B, Srinivasan, PP, Ortiz-Cayon, R, Kalantari, NK, Ramamoorthi, R, Ng, R, et al. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Trans Graph* 2019;38(4). URL: <https://doi.org/10.1145/3306346.3322980>. doi:10.1145/3306346.3322980.
- [30] Zhou, T, Tulsiani, S, Sun, W, Malik, J, Efros, AA. View synthesis by appearance flow. In: *European conference on computer vision*. Springer; 2016, p. 286–301.
- [31] Kalantari, NK, Wang, TC, Ramamoorthi, R. Learning-based view synthesis for light field cameras. *ACM Transactions on Graphics (TOG)* 2016;35(6):193.
- [32] Srinivasan, PP, Wang, T, Sreelal, A, Ramamoorthi, R, Ng, R. Learning to synthesize a 4d RGBD light field from a single image. *International Conference on Computer Vision (ICCV)* 2017 2017;.
- [33] Song, J, Chen, X, Hilliges, O. Monocular neural image based rendering with continuous view control. In: *ICCV* 2019. 2019, p. 4089–4099.
- [34] Wilburn, B, Joshi, N, Vaish, V, Talvala, EV, Antunez, E, Barth, A, et al.



- High performance imaging using large camera arrays. *ACM Trans Graph* 2005;24(3).
- [35] Broxton, M, Flynn, J, Overbeck, R, Erickson, D, Hedman, P, DuVall, M, et al. Immersive light field video with a layered mesh representation. *ACM Transactions on Graphics (Proc SIGGRAPH)* 2020;39(4):86:1–86:15.
  - [36] Guo, K, Lincoln, P, Davidson, P, Busch, J, Yu, X, Whalen, M, et al. The relightables: Volumetric performance capture of humans with realistic relighting. *ACM Trans Graph* 2019;38(6).
  - [37] Collet, A, Chuang, M, Sweeney, P, Gillett, D, Evseev, D, Calabrese, D, et al. High-quality streamable free-viewpoint video. *ACM Trans Graph* 2015;34(4).
  - [38] Pozo, AP, Toksvig, M, Schrager, TF, Hsu, J, Mathur, U, Sorkine-Hornung, A, et al. An integrated 6dof video camera and system design. *ACM Trans Graph* 2019;38(6).
  - [39] Penner, E, Zhang, L. Soft 3d reconstruction for view synthesis. *ACM Transactions on Graphics (TOG)* 2017;36(6):235.
  - [40] Luo, X, Huang, J, Szeliski, R, Matzen, K, Kopf, J. Consistent video depth estimation. *ACM Transactions on Graphics (Proceedings of ACM SIGGRAPH)* 2020;39(4).
  - [41] Yoon, JS, Kim, K, Gallo, O, Park, HS, Kautz, J. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 2020;.
  - [42] Sun, C, Sun, M, Chen, HT. Direct voxel grid optimization: Superfast convergence for radiance fields reconstruction. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2022, URL: <https://arxiv.org/abs/2111.11215>. doi:10.48550/ARXIV.2111.11215.
  - [43] Müller, T, Evans, A, Schied, C, Keller, A. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:220105989* 2022;.
  - [44] Bao, W, Lai, WS, Ma, C, Zhang, X, Gao, Z, Yang, MH. Depth-aware video frame interpolation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2019, p. 3703–3712.
  - [45] Davis, J, Ramamoorthi, R, Rusinkiewicz, S. Spacetime stereo: A unifying framework for depth from triangulation. In: *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.*; vol. 2. IEEE; 2003, p. II–359.
  - [46] Moulon, P, Monasse, P, Perrot, R, Marlet, R. Openmvg: Open multiple view geometry. In: *International Workshop on Reproducible Research in Pattern Recognition*. Springer; 2016, p. 60–74.
  - [47] Moulon, P, Monasse, P, Marlet, R. Adaptive structure from motion with a contrario model estimation. In: *Proceedings of the Asian Computer Vision Conference (ACCV 2012)*. Springer Berlin Heidelberg; 2012, p. 257–270. doi:10.1007/978-3-642-37447-0\_20.
  - [48] Cornelis, K, Verbiest, F, Van Gool, L. Drift detection and removal for sequential structure from motion algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 2004;26(10):1249–1259. doi:10.1109/TPAMI.2004.85.
  - [49] Lowe, DG. Distinctive image features from scale-invariant keypoints. *International journal of computer vision* 2004;60(2):91–110.
  - [50] Barnes, C, Shechtman, E, Finkelstein, A, Goldman, DB. Patchmatch: A randomized correspondence algorithm for structural image editing. In: *ACM SIGGRAPH 2009 Papers*. SIGGRAPH '09; New York, NY, USA: Association for Computing Machinery. ISBN 9781605587264; 2009, URL: <https://doi.org/10.1145/1576246.1531330>. doi:10.1145/1576246.1531330.
  - [51] Schonberger, JL, Frahm, JM. Structure-from-motion revisited. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, p. 4104–4113.
  - [52] Wang, Z, Bovik, A, Sheikh, H, Simoncelli, E. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing* 2004;13(4):600–612. doi:10.1109/TIP.2003.819861.
  - [53] Soundararajan, R, Bovik, AC. Video quality assessment by reduced reference spatio-temporal entropic differencing. *IEEE Transactions on Circuits and Systems for Video Technology* 2013;23(4):684–694. doi:10.1109/TCSVT.2012.2214933.