

# DS-Hybrid GS: auto-reconstruction and decomposition from monocular video

Anonymous CVPR submission

Paper ID 4792

## Abstract

001     *The dynamic nature of our environment, paired with the*  
002     *prevalence of monocular videos as a medium for capturing*  
003     *reality, necessitates efficient methods for reconstructing*  
004     *high-dimensional representations from low-dimensional*  
005     *video data while tackling dynamic and static elements sep-*  
006     *arately for further editing. However, existing methods of-*  
007     *ten rely heavily on costly and ambiguous priors, such as*  
008     *manually labeled masks, optical flow, and depth informa-*  
009     *tion, which can introduce inefficiencies and result in sub-*  
010     *optimal reconstructions. In this work, we propose a novel*  
011     *reconstruction approach that not only facilitates the re-*  
012     *construction of dynamic scenes from monocular videos but*  
013     *also effectively decouples dynamic and static components.*  
014     *To achieve this, we introduce an automatic segmentation*  
015     *pipeline that distinguishes between dynamic and static ob-*  
016     *jects. Our method leverages large pre-trained models for*  
017     *generating high-quality masks and employs feature point*  
018     *registration to enhance generalization beyond traditional*  
019     *optical flow techniques. Additionally, we incorporate a se-*  
020     *mantic filter to further refine the segmentation results. The*  
021     *second stage of our process focuses on the reconstruction*  
022     *of dynamic scenes, relying solely on the dynamic masks*  
023     *obtained in the first step. This approach results in a hy-*  
024     *brid dynamic scene representation that enables effective*  
025     *dynamic decomposition called **DS-Hybrid GS**(Dynamic-*  
026     *Static-Hybrid Gaussian Splatting). By using only the dy-*  
027     *namic masks as prompt inputs, our method becomes robust*  
028     *and applicable to a wider range of datasets. This work*  
029     *makes three key contributions: (1) an automatic monocular*  
030     *video reconstruction method that facilitates the decoupling*  
031     *of dynamic components, (2) an innovative dynamic ele-*  
032     *ment classifier based on point density matching, (3) In street*  
033     *scenes, it significantly improves the efficiency and reduces*  
034     *the cost of reconstruction, and (4) the insight that fewer*  
035     *constraints in the reconstruction process lead to greater ro-*  
036     *bustness in monocular scene reconstruction.*



Figure 1. **Video with extreme camera motion** Previous methods, such as OmnimatteRF, address video with parallax effect well. However, when facing this kind of monocular video, previous methods often tend to treat the entire scene as part of the foreground due to the influence of optical flow. In contrast, our method, which uses only the mask as input, is not affected by the ambiguities of optical flow, enabling more accurate segmentation results.

## 1. Introduction

The world we inhabit can be viewed as a dynamic scene characterized by both temporal and spatial consistency. At the same time, video, particularly monocular video capturing everyday life in a causal manner, has become the primary medium for recording and representing the real world. Recovering high-dimensional representations of the real world from low-dimensional video data has long been a significant challenge and remains crucial for scaling up machine intelligence’s ability to perceive and understand the world. Furthermore, the decoupling of dynamic and static components in the reconstructed scene is essential for subsequent applications and editing, playing a key role in the practical viability of reconstruction techniques.

Recent works [7, 12, 14, 15, 20, 22, 42] have demonstrated promising results in reconstructing dynamic scenes from monocular videos, with the majority treating dynamic

054 and static scenes separately. To achieve robust outcomes,  
055 these studies often incorporate various priors and supervi-  
056 sion, such as optical flow, depth information, and manual  
057 annotations to assist in reconstruction. However, methods  
058 relying on additional prior constraints can be costly and may  
059 encounter unavoidable ambiguities. Thus, obtaining eco-  
060 nomically efficient and robust information suitable for dy-  
061 namic components has become a critical issue in this way  
062 of monocular video reconstruction. Meanwhile, some ap-  
063 proaches [26, 39, 41] treat all scenes as dynamic for re-  
064 construction. However, these methods are often limited  
065 by their rendering techniques and storage formats. Adding  
066 temporal information to static scenes can introduce sig-  
067 nificant redundancy, resulting in inefficient reconstruction  
068 processes unsuited for subsequent editing tasks. There-  
069 fore, optimizing these methods for cost reduction and effi-  
070 ciency, while ensuring their usability, remains an important  
071 challenge. In prior work on dynamic-static decomposition,  
072 methods [2] have demonstrated reasonable performance  
073 in simple cases, effectively separating dynamic and static  
074 components. Since the concept of Omnimatte [21] pro-  
075 posed as a general of co-effects such as shadows in dynamic  
076 scenes, related and following works [12, 16, 21, 33, 39] set-  
077 ting a higher standard for decomposition tasks. However,  
078 these methods struggle when applied to more complex dy-  
079 namic scenarios. As illustrated in Fig. 1, significant camera  
080 motion can introduce optical flow ambiguities, which cause  
081 previous methods to fail and lack robustness. Addressing  
082 such cases, particularly more extreme dynamic situations  
083 remains an open problem.

084 Therefore, we first propose an efficient and rapid solu-  
085 tion for obtaining dynamic scene information by designing  
086 a pipeline for the automatic segmentation of dynamic and  
087 static objects in videos. We leverage the priors from large  
088 pre-trained models to achieve high-quality masks and inno-  
089 vatively employ feature point registration to filter dynamic  
090 objects, offering greater generalization compared to similar  
091 optical flow methods. Additionally, we utilize a semantic  
092 filter to refine the results. Furthermore, addressing the de-  
093 mands for efficient storage and rendering while considering  
094 the characteristics of monocular videos, we introduce a hy-  
095 brid dynamic scene representation that enables dynamism  
096 decomposition with only dynamic masks as additional in-  
097 put. This approach significantly reduces redundant infor-  
098 mation compared to previous dynamic methods. Further-  
099 more, our research has substantiated that a reduction in the  
100 number of constraints leads to enhanced robustness in per-  
101 formance.

102 1. A fully automatic monocular video reconstruction  
103 method while enabling dynamic components decoupling.

104 2. Novel dynamic elements classifier base match points  
105 density.

106 3. In street scene reconstruction, we eliminated the re-

liance on bounding boxes, significantly reducing costs and  
enhancing reconstruction efficiency. 107  
108

4. An insight that fewer constraints provide more robust-  
ness in monocular reconstruction. 109  
110

## 2. Related work 111

**3D reconstruction** Throughout the development of com-  
puter vision, recovering spatial information from images  
has remained a challenging problem. Traditional meth-  
ods [1, 29, 30, 32] have primarily focused on reconstruct-  
ing geometric information. In recent years, however, novel  
view synthesis approaches have emerged, such as Neu-  
ral Radiance Fields (NeRF) [24] and its subsequent exten-  
sions [3, 4, 9, 35, 48] which are capable of capturing view-  
dependent effects. However, vanilla NeRF requires query-  
ing the MLP for hundreds of points each ray, significantly  
constraining its training and rendering speed. Although  
some works [10, 17, 19, 25, 28, 44] have attempted to im-  
prove the training or rendering speed, these methods re-  
main confined to the nuances of differentiable volume ren-  
dering until 3D Gaussian Splatting(3DGS) [13] proposed.  
3DGS utilizes rasterization to achieve real-time rendering  
of high-quality results in complex scenes. While numer-  
ous subsequent works have made advancements in geome-  
try reconstruction [11], large-scale representation [31], and  
anti-aliasing [45], we argue that the inherent effectiveness  
of Gaussian primitives, coupled with the theoretical founda-  
tion of Gaussian Mixture Models (GMM) for fitting arbi-  
trary shape probability distributions, renders 3DGS a more  
robust representation for static scenes in our work. 112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135

**Monocular video reconstruction** While the input of re-  
construction of a scene is diverse, monocular video is the  
most common and challenging set. With the emergence  
of NeRF and 3DGS, various works [8, 14, 15, 20, 22, 26,  
36, 37, 40, 41, 46, 47] have attempted to address this issue.  
Many of these approaches [14, 15, 37, 47] utilize prior in-  
formation, such as optical flow and depth information, to  
guide the reconstruction process. However, as dycheck [8]  
points out, most works focus on quasi-static scenes and are  
not generalizable for most videos in our lives. In our work,  
our primary focus is on utilizing robust prior knowledge to  
reconstruct causal monocular videos. 136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147

**Video dynamics decomposition** Video dynamic decom-  
position plays a fundamental and vital role in diverse video  
editing. Traditional methods have largely depended on  
green screens, multi-view observations, or rotoscoping.  
However, these approaches do not apply to typical monocu-  
lar videos encountered in everyday situations. Thus, several  
methods [2, 18] have attempted to address the decoupling  
of dynamic components in monocular videos, successfully  
148  
149  
150  
151  
152  
153  
154  
155

156 isolating the RGBA representation of both the foreground  
157 and background. However, prior methods primarily fo-  
158 cused on the main dynamic components or considered shad-  
159 dows in isolation [38], neglecting the overall associated ef-  
160 fects of dynamic elements, such as shadows and lighting.  
161 Omnimatte [21], was the first to propose a generic frame-  
162 work capable of learning all associated effects. In recent  
163 years, highly relevant improvements have emerged, whether  
164 by incorporating 3D information [16, 33], employing self-  
165 supervised techniques to obtain foregrounds [39, 42], or  
166 utilizing UV mapping to facilitate follow editing [12], all  
167 of which have shown promising results. However, similar  
168 to NeRF [24], these methods are limited by their render-  
169 ing techniques or over-reliance on priors like optical flow  
170 and depth estimation, resulting in insufficient robustness.  
171 Recent work on 4D Gaussian Splatting (4DGS) has also  
172 demonstrated some capabilities for dynamic reconstruction  
173 and static-dynamic decoupling. Building on this, we aim  
174 to enhance the robustness of the video dynamic Omni-  
175 matte decomposition framework by refining 4DGS, thereby  
176 broadening its applicability to causal videos.

### 177 3. Method

178 Given a monocular video, our task is to generate a high-  
179 quality reconstructed scene while decoupling dynamic and  
180 static elements. Monocular video reconstruction is ill-posed  
181 since the observation of dynamic objects is limited under  
182 one view of one frame and is always insufficient. Although  
183 the static scene usually has richer views and information  
184 to achieve a stable and reliable result, previous methods  
185 strongly rely on various priors to help reconstruct the dy-  
186 namic part. Our method can reconstruct a dynamic scene  
187 from a general monocular video taken freely and achieve the  
188 decomposition of dynamic and static scenes. It is fully self-  
189 supervised and does not require additional training data, no  
190 manual labeling, and no optical flow, but only input videos.

191 The overview of our method is divided into two stages  
192 as shown in Figure 2, stage 1 is to obtain a high-quality  
193 mask for the dynamic part including (b)Automask and  
194 (c)Dynamic classifier, and stage 2 is to decouple the static  
195 part and the dynamic part and reconstruct the dynamic  
196 scene. In stage 1 we will first generate enough numeral  
197 temporal consistent masks, then we will use a robust match-  
198 ing model [5] to gain matched points cross frames and use  
199 epipolar geometry to classify whether the masked object is  
200 dynamic or not. After that, we propose a semantic filter  
201 to avoid potential ambiguity in poor-feature regions. More  
202 details will be explained in Section 3.1.

203 As for stage 2 in Section 3.2, we will introduce the  
204 static dynamic hybrid representation based Gaussian Splat-  
205 ting which largely refers to 4d-gaussian splatting [41] and  
206 the training and optimization details will be discussed in  
207 Section 3.3.

### 3.1. Dynamic Mask Estimation

208 Previous methods for obtaining dynamic masks typically  
209 rely on manual labeling or optical flow techniques [23].  
210 However, manual labeling is labor-intensive and costly,  
211 while optical flow often fails in areas with limited features  
212 or when camera motion exceeds object movement, as seen  
213 in street scenes. To address these limitations, we leverage  
214 large pre-trained models in the preprocessing phase, pro-  
215 viding a more robust and automated pipeline for dynamic  
216 mask generation. Specifically, we use SAM2 [27] to ini-  
217 tialize masks with general semantic priors and RoMa [5] to  
218 classify dynamism based on these priors.  
219

**Mask initialization** The process of obtaining the mask is  
220 always under-considered. The cost of annotation and ro-  
221 bustness challenge is often magnified or prioritized for res-  
222 olution in practical applications. To address this issue, we  
223 initialize masks through SAM2(Segment Anything model  
224 2) [27]. Leveraging the capabilities of SAM2, we can au-  
225 tomatically perform a comprehensive segmentation on an  
226 initial frame, generate the corresponding mask, and propa-  
227 gate through the whole video. This self-generated mask ap-  
228 proach is more labor-efficient and scalable compared to in-  
229 teractive methods. As a large model, it exhibits greater gen-  
230 eralization capability compared to fine-tuned, task-specific  
231 segmentation models. Additionally, it provides more re-  
232 fined edges and retains higher-frequency details than pre-  
233 vious methods.  
234

**Matching-based classifier** After gaining the masks, we  
235 need to classify their dynamism. Previous methods for de-  
236 termining physical dynamism predominantly rely on optical  
237 flow, with some approaches even deriving dynamic masks  
238 directly from optical flow. However, optical flow operates  
239 under a strong assumption of content consistency between  
240 frames, which conflicts with the incomplete observation of  
241 dynamic objects in monocular video sequences. Addition-  
242 ally, optical-flow-based methods [23] tend to fail when en-  
243 counter non-rigid dynamic objects or objects with less  
244 prominent image features. In short, optical flow represents  
245 the correspondence instead of motion itself, which results  
246 in misalignment.  
247

248 We propose a feature-matching-based approach to iden-  
249 tify dynamic objects and design a semantic filter to in-  
250 corporate commonsense knowledge, thereby mitigating un-  
251 avoidable errors arising from ambiguities in image features.  
252 Given two paired frames at  $t_i, t_{i+\Delta t}$ , where  $\Delta t$  is fixed  
253 frame intervals for equal frame rate video to ensure there  
254 will be significant dynamics. We use RoMa [5] to estimate  
255 a dense warp  $W^{t \rightarrow t_i+\Delta t}$  and a matchability score  $p(P_{t_i})$ ,  
256 where  $P_{t_i}$  means the matched key points in frame  $t_i$  pair.  
257 We sample key points paired according to the matchabil-

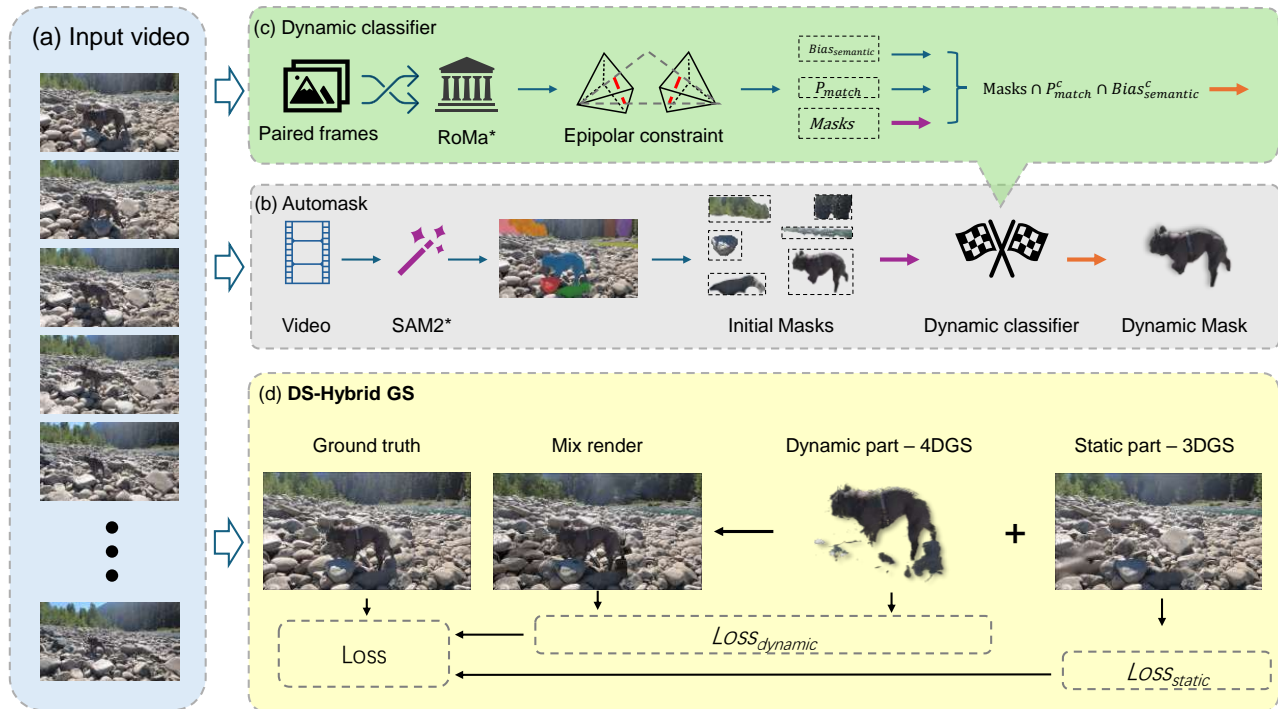


Figure 2. **Pipeline** (a) Our pipeline takes a causal video as input, enabling automatic segmentation and reconstruction with decomposition of dynamic components. (b) First, we preprocess the video using SAM2 [27] to generate a set of high-quality masks for the entire video. Based on these initial masks, we design a dynamic classifier to determine whether the masked object is dynamic. (c) With consistent masks across frames, we select paired frames at fixed intervals and use RoMa [5], a dense feature matching model, for coarse key points matching. Since dynamic objects often violate epipolar constraints, we apply epipolar constraint filtering to obtain refined match key points  $P_{match}^c$ , and the superscript  $^c$  means complementary set. Additionally, we introduce a semantic bias to avoid ambiguity and mismatches in low-feature areas, such as the sky or road. Consequently, the final criterion for a dynamic mask is that it falls within sparsely matched key points and is not within the semantic bias. (d) Finally, given the dynamic part masks obtained from the monocular video, we can initialize a more accurate point cloud via structure-from-motion [29], enabling scene reconstruction with the dynamic parts decoupled.

258 ity score. Since dynamic objects violate the epipolar con-  
 259 straint, we can conclude that the registration points within  
 260 the dynamic mask are likely to be sparser as shown in Fig 3.  
 261 Once we obtain accurately matched key points between two  
 262 frames, we can employ the RANSAC [6] method to esti-  
 263 mate the essential matrix to exclude points that do not con-  
 264 form to epipolar constraint. This condition can be utilized  
 265 to determine whether the objects within the mask are dy-  
 266 namic. Moreover, ambiguity always exists in some low-  
 267 texture parts, thus we propose a semantic bias set to avoid  
 268 some usual textures that may result in sparse match points,  
 269 like sky or road, which significantly work in street scenes.  
 270 Above all, we classify an object as dynamic or not by fol-  
 271 lowing the equation:

$$D = \left\{ x \in I \mid \left| \frac{\Delta t}{T} \sum_{i=1} \text{Density}(P_{t_i}) < \tau \text{ and } x \notin \text{Bias} \right. \right\} \quad (1)$$

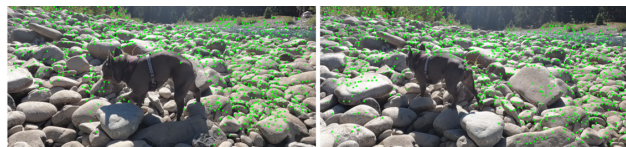


Figure 3. **Match points Example**

273 where  $x \in I$  represents points in the initial mask, 273  
 $\frac{\Delta t}{T} \sum_{i=1} \text{Density}(P_{t_i})$  is the average key point density 274  
 across frames with interval  $\Delta t$ , and  $\tau$  is the dynamic thresh- 275  
 old for identifying sparse matches.  $x \notin \text{Bias}$  excludes 276  
 points within the semantic bias set (e.g., sky or road). 277

### 3.2. Dynamic Scene Representation 278

279 As a scene representation method, 3D Gaussian Splat- 279  
 ting (3DGS) [13] benefits from a well-optimized rasteri- 280  
 zation system on GPUs, achieving high-quality real-time 281  
 novel view synthesis. And 4D Gaussian splotting(4DGS) 282

283 [41], build on 3DGS, achieve real-time photorealistic dy-  
284 namic novel view synthesis. Thus, we aim to combine the  
285 strengths of both approaches by proposing a hybrid scene  
286 representation method that effectively integrates dynamic  
287 and static components.

288 **Preliminary: 4D Gaussian Splatting** 3DGS represent  
289 the whole scene as a cloud of 3D Gaussians while each  
290 Gaussian has a theoretically infinite scope. Compare to nor-  
291 malized Gaussian function in origin 3DGS, 4DGS prove the  
292 the unnormalized Gaussian function of a multivariate Gaus-  
293 sian can be factorized as the production of the unnormal-  
294 ized Gaussian functions of its condition and margin distri-  
295 butions and hold the critical properties. Thus the influence  
296 of a Gaussian on a given spatial position  $x \in \mathbb{R}^3$  defined by  
297 an unnormalized Gaussian function:

$$298 \quad p(x|\mu, \Sigma) = e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}, \quad (2)$$

299 where  $\mu \in \mathbb{R}^3$  is its mean vector, and  $\Sigma \in \mathbb{R}^{3 \times 3}$  is an  
300 anisotropic covariance matrix. For the Mean vector  $\mu$  of a  
301 3D Gaussian is parameterized as  $\mu = (\mu_x, \mu_y, \mu_z)$  in static  
302 scene and  $\mu = (\mu_x, \mu_y, \mu_z, \mu_t)$  for a 4D Gaussian dynamic  
303 scene. And the covariance matrix  $\Sigma$  both Gaussian is fac-  
304 torized same into a scaling matrix  $S$  and a rotation matrix  
305  $R$  as:

$$306 \quad \Sigma = RSS^T R^T, \quad (3)$$

307 where  $S$  is summarized by its diagonal elements  $S =$   
308  $\text{diag}(s_x, s_y, s_z)$ , whilst  $R$  is constructed from a unit quater-  
309 nion  $q$  for static Gaussian. In contrast, 4D Gaussian extent  
310 time dimension equally to a space dimension in scale matrix  
311 and rotation matrix.

312 Moreover, a 3D Gaussian also includes a set of coeffi-  
313 cients of spherical harmonics (SH) for representing view-  
314 dependent color  $c_i(d)$ , where  $c_i$  denotes the color of the  $i$ -  
315 th visible Gaussian from the viewing direction  $d_i$ , along  
316 with an opacity  $\alpha$ . 4DGS proposes to represent  $c_i(d, t)$  as  
317 the combination of a series of 4D spherindrical harmonics  
318 (4DSH) which are constructed by merging SH with differ-  
319 ent 1D-basis functions.

320 In 4DGS rendering, given a pixel with spatial coordi-  
321 nates  $(u, v)$  and timestamp  $t$  in view  $\mathcal{I}$ , its color  $\mathcal{I}(u, v, t)$ ,  
322 after being further factorized as a product of a condi-  
323 tional probability  $p_i(u, v|t)$  and a marginal probability  
324  $p_i(t)$  at time  $t$ , can be computed by blending visible 4D  
325 Gaussians  $p_i(x, y, z, t)$ , that have been sorted according to  
326 their depth. Whereas, a 4D Gaussian can also be factorized  
327 into  $p_i(x, y, z|t)$  and  $t$ , where  $p_i(x, y, z|t)$  is a 3D Gaussian  
328 whose projection in view plane can be approximated by 2D  
329 Gaussian  $p_i(u, v|t)$ . Same as the linearize the perspective  
330 transformations in [13, 41, 49], mean of the derived 2D

Gaussian is obtained as:

$$331 \quad \mu_i^{2d} = \text{Proj}(\mu_i|E, K)_{1:2}, \quad (4) \quad 332$$

333 where  $\text{Proj}(\cdot|E, K)$  denotes the transformation from the  
334 world space to the image space given the intrinsic  $K$  and  
335 extrinsic  $E$ . The covariance matrix is given by

$$336 \quad \Sigma_i^{2d} = (JE\Sigma E^T J^T)_{1:2,1:2}, \quad (5) \quad 337$$

338 where  $J$  is the Jacobian matrix of the perspective projection.  
339 After get the 2D Gaussian  $p_i(u, v|t)$  for alpha blending, the  
340 rendering equation can be described as below:

$$341 \quad I(u, v, t) = \sum_{i=1}^N p_i(t) p_i(u, v|t) \alpha_i c_i(d, t) \quad 342$$

$$343 \quad \times \prod_{j=1}^{i-1} (1 - p_j(t) p_j(u, v|t) \alpha_j). \quad (6) \quad 344$$

345 where  $c_i(d, t)$  denotes the color of the  $i$ -th visible Gaussian  
346 from the viewing direction  $d_i$  at timestamp  $t$ ,  $\alpha_i$  represents  
347 its opacity.

348 **3D-4D hybrid representation** Although 4DGS has  
349 demonstrated remarkable results in reconstructing dynamic  
350 scenes, its performance in novel view synthesis is subopti-  
351 mal. This is primarily due to the tendency of 4DGS to  
352 overfit static scenes, resulting in issues with spatial consis-  
353 tency for static elements and significantly increasing both  
354 rendering and storage overhead. In other words, 4DGS can-  
355 not distinguish the view effect or time effect with only a  
356 monocular video as input. Inspired by techniques in video  
357 matting, we propose a hybrid 3D-4D hybrid scene represen-  
358 tation. In this approach, dynamic regions are first masked  
359 out, allowing separate optimization of the static scene be-  
360 fore refining the dynamic regions. It is noteworthy that  
361 certain "Omnimatte" elements, like shadows and reflections  
362 within the static scene, may also be learned as part of the dy-  
363 namic representation. To obtain a cleaner decoupled scene, a  
364 retraining strategy with background constraints can further  
365 enhance the decoupling of static and dynamic components.  
366 And since we majorly focus on monocular reconstruction,  
we simplify the general mix rendering formula into alpha  
blending which is enough for most monocular video cases  
shown below:

$$367 \quad I_{\text{blend}}(u, v, t) = \sum_{i=1}^N \alpha_i I_i(u, v, t) + \left(1 - \sum_{i=1}^N \alpha_i\right) I_b(u, v) \quad 368$$

$$369 \quad (7) \quad 370$$

371 where  $I_i(u, v, t)$  and  $\alpha_i$  means the  $i$ -th dynamic fore-  
ground color and its alpha, and  $I_b(u, v)$  is the time-invariant  
static background color. Our representation enables accu-  
rate scene reconstruction with improved spatial consistency,

372 making it efficient for subsequent editing tasks. Further-  
373 more, the internal representation is flexible and can be up-  
374 dated to incorporate any advanced image-based reconstruc-  
375 tion methods as they become available.

### 376 3.3. Optimization

377 Our training process is divided into two main stages: first,  
378 training the static components, followed by the optimiza-  
379 tion of the dynamic elements. Supervision is provided by  
380 a combination of three distinct loss functions. During the  
381 static training phase, dynamic regions are masked out, and  
382 only the real images of the static components are used to  
383 guide the reconstruction through a static loss term. Once  
384 the static scene converges, the static model is frozen, and  
385 we proceed with the reconstruction of the dynamic scene.  
386 In this stage, a background loss is introduced to preserve  
387 the clarity of the background in the dynamic regions after  
388 initializing the dynamic part, while the overall reconstruc-  
389 tion is supervised by comparing the rendered scene with the  
390 full real image using a reconstruction loss.

391 The total loss function used for training is as follows:

$$392 \text{Loss} = \lambda_{dssim}L_1 + (1 - \lambda_{dssim})L_{ssim} + \lambda_{stage}L_{bg} \quad (8)$$

393 Here,  $L_1$  is the L1 loss, which measures the absolute dif-  
394 ference between the predicted image and the ground truth  
395 image.  $L_{ssim}$  is the SSIM loss, which is derived from the  
396 Structural Similarity Index Measure and evaluates percep-  
397 tual similarity.  $L_{bg}$  represents the background loss, which  
398 contains the same helps to reduce artifacts in the dynamic  
399 layers. The hyperparameter  $\lambda_{dssim}$  determines the relative  
400 contribution of the L1 and SSIM losses, while  $\lambda_{stage}$  is a bi-  
401 nary parameter that ensures the background loss is applied  
402 only during the dynamic stage of training.

403 **Initialization** We first obtain the dynamic mask using  
404 our proposed method and preprocess the images with this  
405 mask. These preprocessed images are then used to derive  
406 the camera poses and initial point cloud required for Gaus-  
407 sian Splatting through structure-from-motion [29, 30]. Both  
408 the static and dynamic scene representations are initialized  
409 based on this point cloud.

410 **Implementation details** For the SAM2 hyperparameters,  
411 we set 64 points per side, 128 points per batch and only  
412 one crop number of layers. As for the RoMa matching  
413 model [5], we did not finetune or modify the model. The  
414 dynamic threshold in most cases is 0.01. As for the opti-  
415 mization part, we use Adam optimizer and we perform both  
416 10000 iterations in static and hybrid training stages for gen-  
417 eral scenes like in Omnimatte-wild datasets [21] and both  
418 30000 iterations for street scenes since it is harder to cov-  
419 erage with the same learning rate 0.00016. The  $\lambda_{dssim}$  is

0.2. Training a general scene usually takes 3.5 hours on a  
single RTX4090 graphic card. And all the preprocessing is  
also done on the same device. Our code and dataset will be  
made public and available.

## 4. Experiments

In this section, we present a comprehensive comparison  
with state-of-the-art methods via both qualitative and quan-  
titative evaluations. For qualitative results, we assess our  
method across various datasets, unlike previous methods  
that primarily focus on video reconstruction with limited  
camera motion dynamics, we also evaluate our approach on  
Waymo [34] datasets, which involve more complex and dy-  
namic scenarios. Additionally, our automatic mask genera-  
tion method is rigorously tested on these practical datasets  
to validate its applicability and robustness in real-world  
conditions. Meanwhile, quantitative results are presented  
on Waymo dataset for the reconstruction performance and  
Movies dataset to examine decomposition capability and  
visual fidelity. The Movies dataset, proposed by Omni-  
matteRF [16], includes ground truth backgrounds specifi-  
cally designed to evaluate model decoupling performance.  
The baseline setting and preprocessing are the same as the  
instruction.

### 4.1. Qualitative Evaluation

We present a qualitative comparison of methods in Fig. 4  
and Fig. 5. In Fig. 4, our method demonstrates strong  
performance on this dataset. By using only the dynamic  
car masks, our approach successfully incorporates the as-  
sociated shadows into the dynamic region. In contrast,  
the baseline method struggles with accurately decompos-  
ing the static background, as the ambiguity in optical flow  
often causes the network to mistakenly treat all scenes  
as dynamic. In Fig. 5, we also evaluate our method  
in iPhone dataset [26], dynamic scenes dataset [43],  
and Movies dataset with separation results.

### 4.2. Quantitative Evaluation

We select PSNR, SSIM, and LPIPS as the evaluation in-  
dex in both the reconstruction metric and decomposition  
metric. We first quantitatively evaluated our method on  
Waymo dataset, comparing the reconstruction performance  
against baseline methods as shown in Table 1. The results,  
along with several sampled visualizations in Fig 4, demon-  
strate that our approach outperforms current state-of-the-art  
methods of reconstruction from monocular video with de-  
composition ability. For further analysis, we also evalu-  
ated our method on the Movie dataset by comparing met-  
rics between our segmented backgrounds and the ground  
truth backgrounds, as presented in Table 2. (Some results  
cite from OmnimatteRF [16]) Our methods perform second  
best in this metric and have more details in some regions of

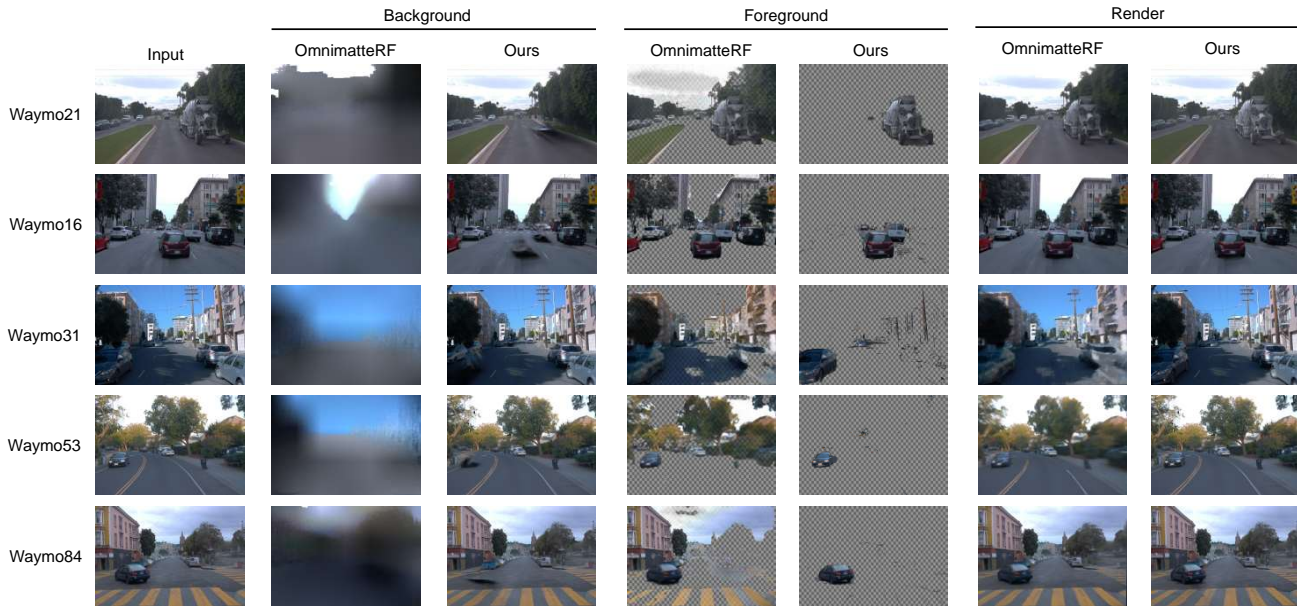


Figure 4. **Waymo qualitative evaluation** We evaluated our method and OmnimatteRF on several Waymo dataset cases. The baseline method fails to handle these cases with extreme camera motion. The input mask is generated by our automatic method instead of the dynamic masks projected from the labeled bounding box in the dataset. For the foreground render we apply an alpha threshold of 0.5 which is the same as in the baseline code.

Waymo	016			021			031			053			084		
	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑	PSNR↑
OmnimatteRF	0.271	0.854	30.60	0.288	0.877	30.82	0.435	0.715	22.47	0.393	0.742	24.04	0.198	0.902	33.31
Ours	<b>0.097</b>	<b>0.968</b>	<b>32.67</b>	<b>0.103</b>	<b>0.962</b>	<b>32.19</b>	<b>0.096</b>	<b>0.956</b>	<b>30.45</b>	<b>0.129</b>	<b>0.938</b>	<b>28.15</b>	<b>0.066</b>	<b>0.975</b>	<b>33.85</b>

Table 1. **Reconstruction quantitative evaluations.** We present the reconstruction comparison of our method and baselines on the waymo datasets. The better results are in **bold**.

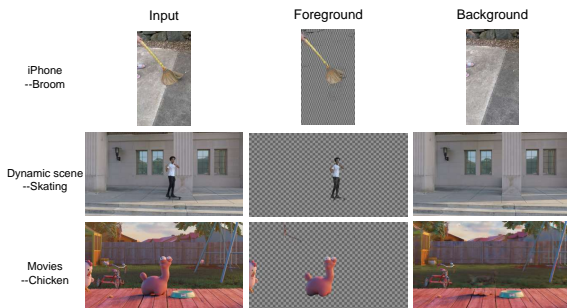


Figure 5. **Various datasets qualitative evaluation**

470 the background compared to the best method as shown in  
471 Fig. 6.

### 472 4.3. Ablation study

473 We demonstrate the effect of our background constraint.  
474 When directly reconstructing the scene, the static part's re-

construction is limited by the constrained observation from  
a monocular video, leading to unavoidable artifacts in some  
high-frequency regions. Meanwhile, in the dynamic part,  
4DGS tends to overfit the static scene components. To mit-  
igate this, we introduce a background loss that encourages  
the regions outside the dynamic masks to remain cleaner,  
helping to refine the reconstruction of the dynamic parts as  
we shown in Fig 7

## 5. Conclusion

We introduce an automated method for the complete recon-  
struction of scenes from monocular videos, with the capa-  
bility to automatically decouple dynamic and static compo-  
nents. Extensive experiments have demonstrated that our  
approach is not only comparable to existing methods on  
simpler datasets but also exhibits superior performance and  
robustness in more complex and rapidly changing scenar-  
ios. Additionally, our method for automatically obtaining  
dynamic masks is readily transferable to other techniques.

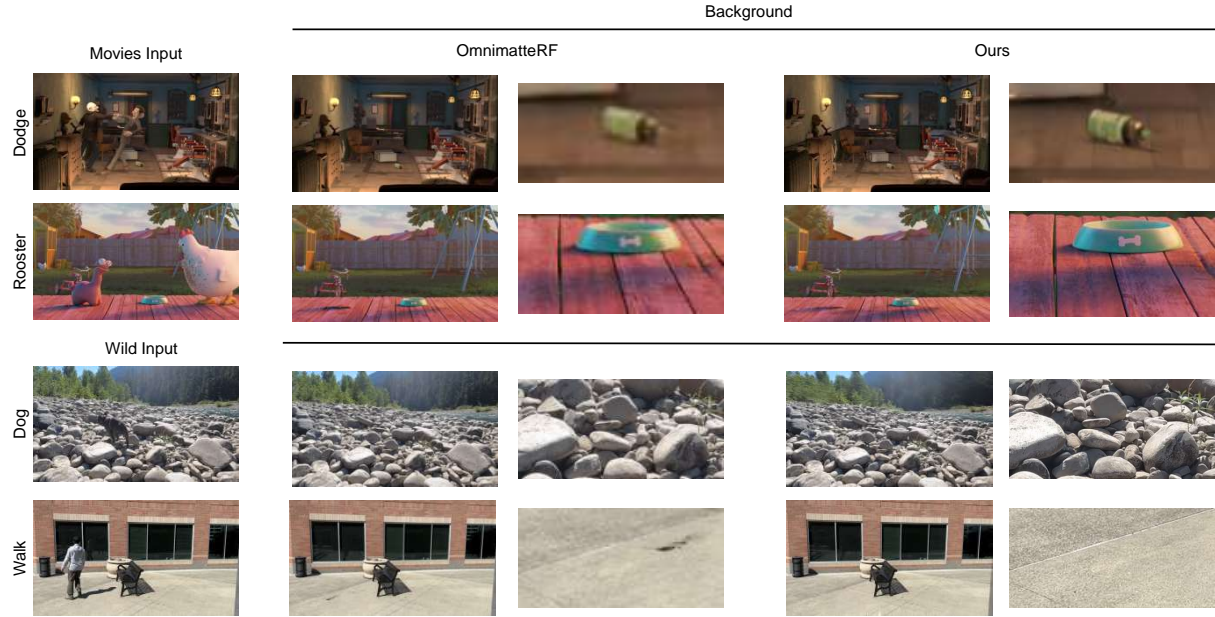


Figure 6. **Background visualizations** We also evaluate our background decomposition capabilities in comparison to OmnimatteRF on the datasets they proposed. Additionally, due to the advantages of Gaussian splatting, our method produces more precise results for the background in the near-field, particularly under similar training conditions.

Movies	Donkey			Dog			Chicken			Rooster			Dodge		
	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑	PSNR↑
D <sup>2</sup> NeRF	-	-	-	0.370	0.694	22.73	-	-	-	0.340	0.708	25.13	0.408	0.729	20.95
Omnimatte	0.315	0.653	19.11	0.279	0.706	21.74	0.312	0.704	20.95	0.220	0.741	23.14	0.067	0.879	23.88
LNA	<u>0.104</u>	<u>0.849</u>	18.79	<u>0.154</u>	0.828	26.08	0.190	0.818	19.22	<u>0.131</u>	0.804	<u>26.46</u>	0.068	0.937	24.94
4DGS	0.357	0.614	16.48	0.427	0.628	19.64	0.390	0.653	19.42	0.558	0.490	15.41	0.333	0.701	19.41
OmnimatteRF	<b>0.005</b>	<b>0.990</b>	<b>38.24</b>	<b>0.030</b>	<b>0.976</b>	<b>31.44</b>	<b>0.021</b>	<b>0.978</b>	<b>32.86</b>	<b>0.024</b>	<b>0.969</b>	<b>27.65</b>	<b>0.006</b>	<b>0.991</b>	<b>39.11</b>
Ours	0.242	0.783	<u>20.60</u>	<u>0.154</u>	<u>0.941</u>	<u>27.51</u>	<u>0.171</u>	<u>0.891</u>	<u>25.67</u>	0.189	<u>0.906</u>	23.95	<u>0.066</u>	<u>0.970</u>	<u>31.86</u>

Table 2. **Decomposition quantitative evaluations.** We present the background reconstruction comparison of our method and baselines on the *Movies* datasets. The best results are in **bold**. The second best results are in underline. Results marked - are the ones where the method failed to give good separations.

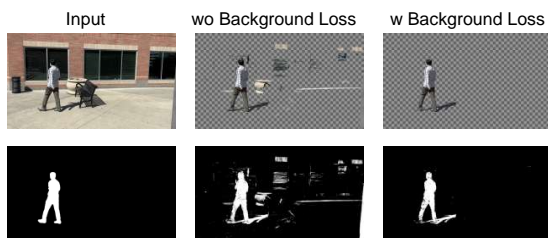


Figure 7. **Ablation** The input consists of the image and dynamic mask, while the output includes the dynamic rendering results and the alpha visualization. The alpha visualization provides a more intuitive way to assess the effectiveness of the background loss, as it clearly illustrates the separation between dynamic and static regions in the scene.

Moreover, our method has several possible further works. 493  
 First, our method can also be conveniently expanded to 494  
 multi-view reconstruction in street scenes, significantly re- 495  
 ducing annotation costs while providing high-quality recon- 496  
 struction results. Second, we could incorporate additional 497  
 regularization and constraints, such as depth, to improve 498  
 the method. But, in this work, we believe that fewer con- 499  
 straints and regularization lead to greater robustness in the 500  
 approach. 501

However, our method is not without limitations. The 502  
 automatic mask method requires threshold adjustments tai- 503  
 lored to specific kinds of scenes, and it may fail with rela- 504  
 tively stationary objects. To address this, we propose to 505  
 extend the video range in autonomous driving scenarios. 506

In summary, our method holds broad potential for appli- 507  
 cation across various domains. 508



509

## References

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

- [1] Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M. Seitz, and Rick Szeliski. Building rome in a day. *Communications of the ACM*, 54: 105–112, 2011. 2
- [2] Xue Bai, Jue Wang, David Simons, and Guillermo Sapiro. Video snapchat: robust video object cutout using localized classifiers. *ACM Transactions on Graphics (ToG)*, 28(3):1–11, 2009. 2
- [3] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. *arXiv e-prints*, 2021. 2
- [4] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. *IEEE*, 2023. 2
- [5] Johan Edstedt, Qiyu Sun, Georg Bökman, Mårten Wadenbäck, and Michael Felsberg. RoMa: Robust Dense Feature Matching. *IEEE Conference on Computer Vision and Pattern Recognition*, 2024. 3, 4, 6
- [6] Martin A Fischler and Robert C. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981. 4
- [7] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5712–5721, 2021. 1
- [8] Hang Gao, Ruilong Li, Shubham Tulsiani, Bryan Russell, and Angjoo Kanazawa. Monocular dynamic view synthesis: A reality check. *Advances in Neural Information Processing Systems*, 35:33768–33780, 2022. 2
- [9] Haoyu Guo, Sida Peng, Haotong Lin, Qianqian Wang, Guofeng Zhang, Hujun Bao, and Xiaowei Zhou. Neural 3d scene reconstruction with the manhattan-world assumption. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5511–5520, 2022. 2
- [10] Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5875–5884, 2021. 2
- [11] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for geometrically accurate radiance fields. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 2
- [12] Yoni Kasten, Dolev Ofri, Oliver Wang, and Tali Dekel. Layered neural atlases for consistent video editing. *ACM Transactions on Graphics (TOG)*, 40(6):1–12, 2021. 1, 2, 3
- [13] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 2, 4, 5
- [14] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2
- [15] Zhengqi Li, Qianqian Wang, Forrester Cole, Richard Tucker, and Noah Snavely. Dynibar: Neural dynamic image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4273–4284, 2023. 1, 2
- [16] Geng Lin, Chen Gao, Jia-Bin Huang, Changil Kim, Yipeng Wang, Matthias Zwicker, and Ayush Saraf. Omnimatterf: Robust omnimatte with 3d background modeling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23471–23480, 2023. 2, 3, 6
- [17] Haotong Lin, Sida Peng, Zhen Xu, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Efficient neural radiance fields for interactive free-viewpoint video. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. 2
- [18] Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian L Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Real-time high-resolution background matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8762–8771, 2021. 2
- [19] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020. 2
- [20] Yu-Lun Liu, Chen Gao, Andreas Meuleman, Hung-Yu Tseng, Ayush Saraf, Changil Kim, Yung-Yu Chuang, Johannes Kopf, and Jia-Bin Huang. Robust dynamic radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 1, 2
- [21] Erika Lu, Forrester Cole, Tali Dekel, Andrew Zisserman, William T Freeman, and Michael Rubinstein. Omnimatte: Associating objects and their effects in video. In *CVPR*, 2021. 2, 3, 6
- [22] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In *3DV*, 2024. 1, 2
- [23] Aravindh Mahendran, James Thewlis, and Andrea Vedaldi. Self-supervised segmentation by grouping optical-flow. In *Computer Vision – ECCV 2018 Workshops*, pages 528–534, Cham, 2019. Springer International Publishing. 3
- [24] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2, 3
- [25] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *ACM transactions on graphics (TOG)*, 41(4):1–15, 2022. 2
- [26] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021. 2, 6
- [27] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

- 623 Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feicht-  
624 enhofer. Sam 2: Segment anything in images and videos.  
625 *arXiv preprint arXiv:2408.00714*, 2024. 3, 4
- 626 [28] Christian Reiser, Rick Szeliski, Dor Verbin, Pratul Srinivasan,  
627 Ben Mildenhall, Andreas Geiger, Jon Barron, and Peter Hedman. Merf: Memory-efficient radiance fields for real-time view synthesis in unbounded scenes. *ACM Transactions on Graphics (TOG)*, 42(4):1–12, 2023. 2
- 628 [29] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 4,  
629 6
- 630 [30] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 6
- 631 [31] Qing Shuai, Haoyu Guo, Zhen Xu, Haotong Lin, Sida Peng, Hujun Bao, and Xiaowei Zhou. Real-time view synthesis for large scenes with millions of square meters. 2024. 2
- 632 [32] Noah Snavely, Steven M Seitz, and Richard Szeliski. Photo tourism: exploring photo collections in 3d. In *ACM siggraph 2006 papers*, pages 835–846. 2006. 2
- 633 [33] Mohammed Suhail, Erika Lu, Zhengqi Li, Noah Snavely, Leonid Sigal, and Forrester Cole. Omnimate3d: Associating objects and their effects in unconstrained monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 630–639, 2023. 2, 3
- 634 [34] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 6
- 635 [35] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022. 2
- 636 [36] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12959–12970, 2021. 2
- 637 [37] Qianqian Wang, Vickie Ye, Hang Gao, Jake Austin, Zhengqi Li, and Angjoo Kanazawa. Shape of motion: 4d reconstruction from a single video. *arXiv preprint arXiv:2407.13764*, 2024. 2
- 638 [38] Tianyu Wang, Xiaowei Hu, Qiong Wang, Pheng-Ann Heng, and Chi-Wing Fu. Instance shadow detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1880–1889, 2020. 3
- 639 [39] Tianhao Wu, Fangcheng Zhong, Andrea Tagliasacchi, Forrester Cole, and Cengiz Oztireli. D<sup>2</sup>nerf: Self-supervised decoupling of dynamic and static objects from a monocular video. *Advances in neural information processing systems*, 35:32653–32666, 2022. 2, 3
- 640 [40] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9421–9431, 2021. 2
- 641 [41] Zeyu Yang, Hongye Yang, Zijie Pan, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. In *International Conference on Learning Representations (ICLR)*, 2024. 2, 3, 5
- 642 [42] Vickie Ye, Zhengqi Li, Richard Tucker, Angjoo Kanazawa, and Noah Snavely. Deformable sprites for unsupervised video decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2657–2666, 2022. 1, 3
- 643 [43] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5336–5345, 2020. 6
- 644 [44] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenocets for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. 2
- 645 [45] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19447–19456, 2024. 2
- 646 [46] Wentao Yuan, Zhaoyang Lv, Tanner Schmidt, and Steven Lovegrove. Star: Self-supervised tracking and reconstruction of rigid objects in motion with neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13144–13152, 2021. 2
- 647 [47] Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, Deqing Sun, and Ming-Hsuan Yang. Monst3r: A simple approach for estimating geometry in the presence of motion. *arXiv preprint arxiv:2410.03825*, 2024. 2
- 648 [48] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv e-prints*, 2020. 2
- 649 [49] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa splatting. *IEEE TVCG*, 2002. 5