

Unsupervised Object-Centric Voxelization for Dynamic Scene Understanding

Siyu Gao* Yanpeng Zhao* Yunbo Wang[†] Xiaokang Yang
 MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University
 {siyu.gao, zhao-yan-peng, yunbow, xkyang}@sjtu.edu.cn
<https://sites.google.com/view/dynavol/>

Abstract

Understanding the compositional dynamics of multiple objects in unsupervised visual environments is challenging, and existing object-centric representation learning methods often ignore 3D consistency in scene decomposition. We propose DynaVol, an inverse graphics approach that learns object-centric volumetric representations in a neural rendering framework. DynaVol maintains time-varying 3D voxel grids that explicitly represent the probability of each spatial location belonging to different objects, and decouple temporal dynamics and spatial information by learning a canonical-space deformation field. To optimize the volumetric features, we embed them into a fully differentiable neural network, binding them to object-centric global features and then driving a compositional NeRF for scene reconstruction. DynaVol outperforms existing methods in novel view synthesis and unsupervised scene decomposition and allows for the editing of dynamic scenes, such as adding, deleting, replacing objects, and modifying their trajectories.

1 Introduction

Performing object-centric unsupervised learning in dynamic visual environments is of great importance but challenging due to the intricate entanglement between the spatial and temporal information [32; 26; 10]. Previous approaches primarily leverage the temporal cues across consecutive video frames but tend to ignore the 3D nature, resulting in a multi-view mismatch of 2D object segmentation [5; 27; 14]. In this paper, we provide an early study of unsupervised dynamic scene decomposition in 3D scenarios, where we believe that an effective object-centric representation should satisfy three conditions: (i) It should be able to accurately represent the *spatial structures* of the visual scene in a stereo-consistent manner and facilitate precise object localization. (ii) It should capture and decouple the *underlying dynamics* of each object from the visual appearance and spatial structures. (iii) It should obtain a *global understanding* of each object, which is crucial for downstream tasks such as scene editing and relational reasoning.

Accordingly, we propose to learn two sets of object-centric representations: one that represents the local spatial structures using 3D voxel grids that may vary over time, and another that represents the global understanding of each object, which is time-invariant. One may wonder about the advantages of introducing 3D voxelization. The answer is that if we can learn voxel grids that explicitly indicate the probability of each spatial location belonging to different objects, we can achieve 3D-consistent scene decomposition naturally. It is worth noting that in our approach, these two sets of global/local representations are interdependent during training, in the sense that well-trained and decoupled global features can guide the model to bind each spatial location to the corresponding object.

*Equal contribution.

[†]Corresponding author: Yunbo Wang.

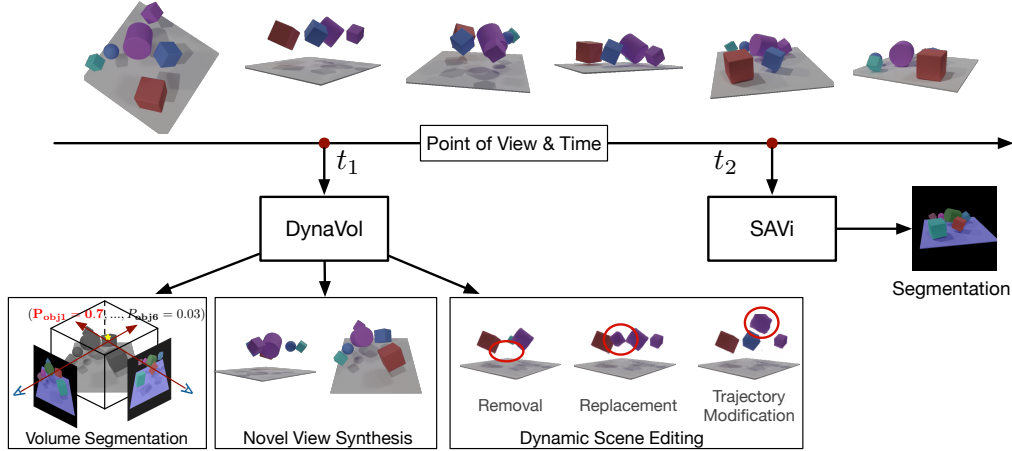


Figure 1: **Contributions:** DynaVol explores an unsupervised object-centric voxelization approach for dynamic scene decomposition. Unlike its 2D counterparts, such as SAVi [15], DynaVol ensures 3D consistency and provides additional capabilities, *e.g.*, novel view synthesis and scene editing.

Based on these intuitions, we propose DynaVol, an inverse graphics framework that learns to perform object-centric voxelization of 3D dynamic scenes. Our approach comprises three network components, including (i) a deformation network that learns the canonical-space transitions of volumetric representations over time, (ii) a 3D slot attention network that progressively refines the object-level, time-invariant global features by aggregating the volumetric representations, and (iii) a voxel-based, compositional neural radiance field (NeRF) for scene reconstruction, which introduces a strong geometric inductive bias that facilitates the learning of object-centric representations. As shown in Figure 1, our unsupervised voxelization approach provides three advantages for dynamic scene understanding. First, it allows for fine-grained separation of object-centric information. Second, it naturally ensures the 3D consistency of the decomposition results. Third, it enables direct scene editing that is not possible in existing dynamic scene decomposition methods [15; 25; 5].

DynaVol is trained on each individual dynamic scene using a two-stage learning scheme. In the first warmup stage, a sparse set of multi-view static images is used to provide strong geometric priors that facilitate the decoupling of spatial and temporal features. In the second stage, the entire model is optimized using sequential data from a monocular moving camera, allowing for a joint refinement of the initial voxel grid features, global slot-attention features, and the dynamics learning module.

We evaluate DynaVol on multiple 3D dynamic scenes with different numbers of objects, diverse motions, various shapes (such as cube, sphere, and real-world shapes), as well as different materials (such as rubber and metal). We demonstrate the effectiveness of our approach in three ways. First, DynaVol outperforms existing scene decomposition approaches (SAVi [15] and uORF [34]) by projecting the object-centric volumetric representations onto 2D planes. Second, it outperforms strong baseline models (D-NeRF [24], DeVRF [18]) for novel view synthesis. Finally, it also performs well for dynamic scene editing, such as object removal, replacement, and trajectory modification, by directly manipulating the voxel grids and the learned deformation function without further training.

2 Related Work

Unsupervised scene decomposition in 2D space. Most existing approaches in this area [9; 8; 6; 1] use latent features to represent objects in 2D scenarios like CLEVR [13]. The slot attention method [19] extracts object-centric latents through a cross-attention mechanism and repeatedly refines them using GRUs [3]. SAVi [15] extends slot attention into dynamic scenes by updating slots at each frame and uses optical flow data as the reconstruction target. STEVE [27] improves SAVi by simply replacing the spatial broadcast decoder with an autoregressive Transformer. SAVi++ [5] uses depth information to improve SAVi for modeling static objects and scenes with camera motion.

Unsupervised scene decomposition in 3D space. Recent methods [14; 2; 34; 28; 25] combine object-centric representations with view-dependent scene modeling techniques like neural radiance fields (NeRFs) [22]. ObSuRF [28] adopts the spatial broadcast decoder and takes depth information as training supervision. uORF [34] extracts the background latent and foreground latents from an

input static image to handle background and foreground objects separately. For dynamic scenes, Li *et al.* [17] proposed an auto-encoder framework that incorporates volume rendering to model dynamic scenes. However, it represents the whole scene using a single latent vector (instead of object-centric features), which can be insufficient for complex scenarios with multiple objects and various motions. Guan *et al.* [11] proposed to use a set of particle-based explicit representations in the NeRF-based inverse rendering framework, which is particularly designed for fluid physics modeling. Driess *et al.* [4] explored the combination of an object-centric auto-encoder and volume rendering for dynamic scenes, which is the most relevant work to our approach. However, unlike our approach which is totally unsupervised, it requires pre-prepared 2D object segments.

Dynamic scene rendering based on NeRFs. Besides those with object-centric representations, there is another line of work [24; 18; 33; 12] that models 3D dynamics using NeRF-based methods. D-NeRF [24] uses a deformation network to map the coordinates of the dynamic fields to the canonical space. DeVRF [18] models dynamic scenes with volume grid features [29] and voxel deformation fields. Recently, D²NeRF [33] presents a motion decoupling framework. However, unlike DynaVol, it cannot segment multiple moving objects.

3 Method

In this section, we first discuss the problem setup and the overall framework of DynaVol (Sect. 3.1). We then introduce a new set of representations for object-centric scene voxelization and present the network details of dynamics modeling, 3D slot attention, and object-centric neural rendering (Sect. 3.2–3.4). Finally, we describe the two-stage training procedure of DynaVol (Sect. 3.5).

3.1 Overview of DynaVol

Problem setup. We assume a set of RGB images of a dynamic scene $\{\mathbf{I}_t^v, \mathbf{T}_t^v\}_{t=1}^T$ collected with a moving monocular camera and a sparse set of views $\{\mathbf{I}_{t_0}^v, \mathbf{T}_{t_0}^v\}_{v=1}^V$ at the initial timestamp. $\mathbf{I}_t^v \in \mathbb{R}^{H \times W \times 3}$ are images acquired under camera poses $\mathbf{T}_t^v \in \mathbb{R}^{4 \times 4}$, T is the length of video frames, and V is the number of views at the starting moment t_0 . The goal is to understand the space-time structures of the visual scene from $\{\mathbf{I}_t^v\}_{t=1}^T$ and $\{\mathbf{I}_{t_0}^v\}_{v=1}^V$ without additional information.

Overall framework. DynaVol is trained in an inverse graphics learning framework to synthesize $\{\mathbf{I}_t^v\}_{t=1}^T$ and $\{\mathbf{I}_{t_0}^v\}_{v=1}^V$ without further supervision. Formally, the goal is to learn an object-centric projection of $(\mathbf{x}, \mathbf{d}, t) \rightarrow \{(\mathbf{c}_n, \sigma_n)\}_{n=1}^N$, where N is the assumed number of objects and $\mathbf{x} = (x, y, z)$ is a 3D point sampled by the neural renderer, which outputs the density and color for each object at view direction \mathbf{d} . By re-combining $\{(\mathbf{c}_n, \sigma_n)\}_{n=1}^N$ to approach true pixel values, the model is required to learn 3D-consistent object-centric representations. As shown in Figure 2, DynaVol maintains voxel grids $\mathcal{V}_{\text{density}}$ and a set of object-level slot features \mathbf{S} . The entire model consists of three network components: (i) Deformation networks f_ψ and f'_ξ that learn the canonical-space transitions of $\mathcal{V}_{\text{density}}$ over time. (ii) A volume encoder E_θ and a slot attention block Z_ω that progressively refine \mathbf{S} . (iii) A compositional NeRF³ N_ϕ that jointly uses $\mathcal{V}_{\text{density}}$ and \mathbf{S} to render the observed images. The training pipeline of DynaVol involves two stages, including a warmup stage that learns $(\mathcal{V}_{\text{density}}, f_\psi, f'_\xi, N_{\phi'})$ and a dynamic grounding stage that learns $(\mathcal{V}_{\text{density}}, \mathbf{S}, f_\psi, E_\theta, Z_\omega, N_\phi)$.

3.2 Object-Centric Volumetric Representation and Dynamics Modeling

Volumetric representation of object probability. Inspired by the idea of using a 3D voxel grid to maintain the volume density for neural rendering, we extend it with a 4D voxel grid, denoted as $\mathcal{V}_{\text{density}}$. The additional dimension is used to indicate the occurrence probabilities of each object within each grid cell. The occurrence probability $\{\sigma_n\}_{n=1}^N$ at an arbitrary 3D location can be efficiently queried through the trilinear interpolation sampling method:

$$\text{Interp}(\mathbf{x}, \mathcal{V}_{\text{density}}) : (\mathbb{R}^3, \mathbb{R}^{N \times N_x \times N_y \times N_z}) \rightarrow \mathbb{R}^N, \quad (1)$$

where (N_x, N_y, N_z) are the resolutions of $\mathcal{V}_{\text{density}}$. To achieve sharp decision boundaries during training, we apply the softplus activation function to the output of trilinear interpolation.

Canonical-space dynamics modeling. Inspired by D-NeRF [24], we use a dynamics module f_ψ to learn the deformation field from the voxel grid $\mathcal{V}_{\text{density}}^{t_0}$ at the initial timestamp to its canonical space

³It is important to note that we use a non-compositional NeRF denoted by $N_{\phi'}$ in the warmup stage.

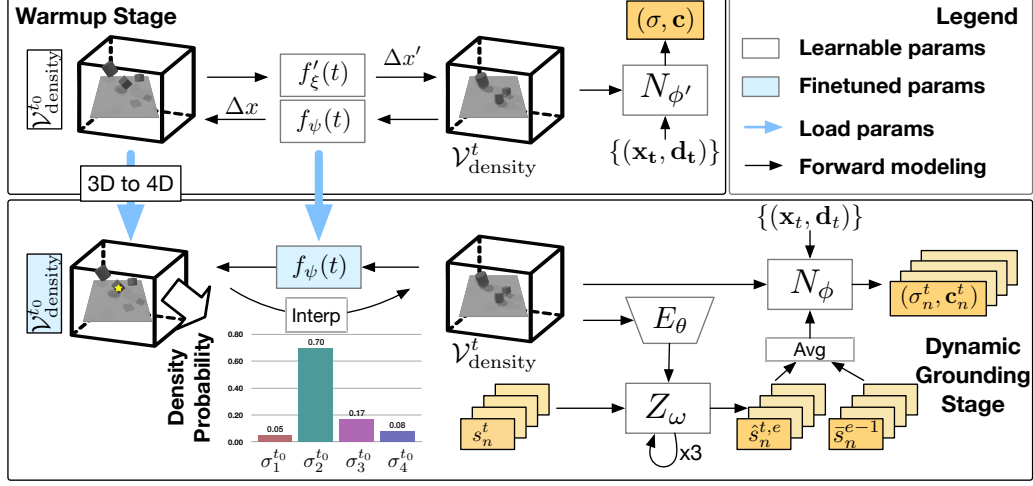


Figure 2: **An overview of DynaVol.** The model comprises three network components, including the voxel grid deformation module, the slot features refinement module, and the neural rendering module. The model has two training stages, including a warmup stage and a dynamic grounding stage.

variations over time. Given a 3D point $\forall \mathbf{x}_i \in \{\mathbf{x}\}$ at an arbitrary time, $f_\psi(\mathbf{x}_i, t)$ predicts a position movement $\Delta \mathbf{x}_i$, so that we can transform \mathbf{x}_i to the scene position at the first moment by $\mathbf{x}_i + \Delta \mathbf{x}_i$. We then query the occurrence probability from $\mathcal{V}_{\text{density}}^{t_0}$ by $\hat{\mathcal{V}}_{\text{density}}^t = \text{Interp}((\mathbf{x}_i + f_\psi(\mathbf{x}_i, t)), \mathcal{V}_{\text{density}}^{t_0})$. Notably, we encode \mathbf{x}_i and t into higher dimensions via positional embedding. Additionally, we use another dynamics module f'_ξ in the warmup stage to model the forward movement $\Delta x'$ from initial movement to timestamp t . This module enables the calculation of a cycle-consistency loss, ensuring the coherence of the estimated motion. More details are introduced in the Sect. 3.5.

3.3 Object-Centric Global Representation and 3D Slot Attention

Slot features. In addition to the time-varying volumetric representations, we further learn another set of object-centric features that are **time-invariant** and represent the global understanding of each object. Specifically, we use a set of latent codes referred to as “slots” to represent these object-level features. This terminology is in line with prior literature on 2D static scene decomposition [19]. The slots are randomly initialized from a normal distribution and progressively refined episode by episode throughout our second training stage. In this context, an *episode* refers to a training process that iterates through every moment of the data sequence. We denote the slots by $\mathbf{S}_{t,e} \in \mathbb{R}^{N \times D}$, where e is the episode index and D is the feature dimensionality. At the beginning of each episode, we have $\mathbf{S}_{t_0,e} = \bar{\mathbf{S}}_{e-1}$, where $\bar{\mathbf{S}}_{e-1}$ represents the average of $\{\mathbf{S}_{t,e-1}\}_{t=1}^T$ across all timestamps in the previous episode. Each slot feature captures the time-invariant properties such as the appearance of each object, which enables the manipulation of the scene’s content and relationships between objects.

Encoder. To bind the voxel grid representations to the corresponding object, at timestamp t , we pass $\mathcal{V}_{\text{density}}^t$ through a 3D CNN encoder E_θ , which consists of 3 convolutional layers with ReLU. It outputs N flattened features $\mathbf{h}_t \in \mathbb{R}^{M \times D}$, where M represents the size of the voxel grids that have been reduced in dimensionality by the encoder. From an optimization perspective, a set of well-decoupled global slot features can benefit the separation of the object-centric volumetric representations.

Slot attention. To refine the slot features, we employ the iterative attention block denoted by Z_ω to incorporate the flattened local features \mathbf{h}_t . In a single round of slot attention, we have:

$$\mathcal{A}_t = \text{softmax}_N \left(\frac{1}{\sqrt{D}} k(\mathbf{h}_t) \cdot q(\mathbf{S}_t)^T \right), \quad W_t^{i,j} = \frac{\mathcal{A}_t^{i,j}}{\sum_{l=1}^M \mathcal{A}_t^{l,j}}, \quad \mathcal{U}_t = W^T \cdot v(\mathbf{h}_t), \quad (2)$$

where (q, k, v) are learnable linear projections $\mathbb{R}^{D \rightarrow D}$ [20], such that $\mathcal{A}_t \in \mathbb{R}^{M \times N}$ and $\mathcal{U}_t \in \mathbb{R}^{N \times D}$, and \sqrt{D} is a fixed softmax temperature [30]. The resulted slots features are then updated by a GRU as $\hat{\mathbf{S}}_t = \text{GRU}(\mathcal{U}_t, \mathbf{S}_t)$. We repeat the attention computation 3 times at each timestamp.

3.4 Compositional Neural Renderer

Forward modeling. The renderer can be denoted by $N_\phi(\mathbf{x}, \mathbf{d} \mid \{\bar{\mathbf{s}}_n\})$. Previous compositional NeRF, like in uORF [34], typically uses an MLP to learn a continuous mapping g from sampling point \mathbf{x} , viewing direction \mathbf{d} , and slot features $\{\mathbf{s}_n\}$, to the emitted densities $\{\sigma_n\}$ and colors $\{\mathbf{c}_n\}$ of different slots. While in our renderer N_ϕ , we only adopt the MLP to learn the object-centric projections $g': (\mathbf{x}, \mathbf{d}, \{\bar{\mathbf{s}}_n\}) \rightarrow \{\mathbf{c}_n\}$ and query $\{\sigma_n\}$ directly from the voxel grid $\hat{\mathcal{V}}_{\text{density}}^t$ at the corresponding timestamp. At timestamp t , N_ϕ takes as inputs $\{\bar{\mathbf{s}}_n\}_t = \text{mean}(\hat{\mathbf{S}}_t, \bar{\mathbf{S}}_{e-1})$, where $\bar{\mathbf{S}}_{e-1}$ is the average of $\{\mathbf{S}_{t,e-1}\}_{t=1}^T$ across all timestamps in the previous episode. We use density-weighted mean to compose the predictions of \mathbf{c}_n and σ_n for different objects, such that:

$$w_n = \sigma_n / \sum_{n=1}^N \sigma_n, \quad \bar{\sigma} = \sum_{n=1}^N w_n \sigma_n, \quad \bar{\mathbf{c}} = \sum_{n=1}^N w_n \mathbf{c}_n, \quad (3)$$

where $\bar{\sigma}$ and $\bar{\mathbf{c}}$ is the output density and the color of a sampling point. We estimate the color $C(\mathbf{r})$ of a sampling ray with the quadrature rule [21]: $\hat{C}(\mathbf{r}) = \sum_{i=1}^P T_i (1 - \exp(-\bar{\sigma}_i \delta_i)) \bar{\mathbf{c}}_i$, where $T_i = \exp(-\sum_{j=1}^{i-1} \bar{\sigma}_j \delta_j)$, P is the number of sampling points in a certain ray, and δ_i is the distance between adjacent samples along the ray.

Objectives. At a specific timestamp, we take the rendering loss $\mathcal{L}_{\text{Render}}$ between the predicted and observed pixel colors, the background entropy loss \mathcal{L}_{Ent} , and the per-point RGB loss $\mathcal{L}_{\text{Point}}$ following DVGO [29] as basic objective terms. \mathcal{L}_{Ent} can be viewed as a regularization to encourage the renderer to concentrate on either foreground or background. To enhance dynamics learning in the warmup stage, we design a novel cycle loss between f_ψ and f'_ξ :

$$\begin{aligned} \mathcal{L}_{\text{Render}} &= \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \left\| \hat{C}(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2, & \mathcal{L}_{\text{Ent}} &= \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} -\hat{w}_l^r \log(\hat{w}_l^r) - (1 - \hat{w}_l^r) \log(1 - \hat{w}_l^r), \\ \mathcal{L}_{\text{Point}} &= \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \left(\frac{1}{P_r} \sum_{i=0}^{P_r} \left\| \bar{\mathbf{c}}_i - C(\mathbf{r}) \right\|_2^2 \right), & \mathcal{L}_{\text{Cyc}} &= \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \left(\frac{1}{P_r} \sum_{i=0}^{P_r} \left\| f_\psi(x_i, t) + f'_\xi(x'_i, t) \right\|_2^2 \right), \end{aligned} \quad (4)$$

where \mathcal{R} is the set of sampled rays in a batch, P_r is the number of sampling points along ray r , $x'_i = x_i + f_\psi(x_i, t)$, $i \in [0, P_r]$, and \hat{w}_l^r is the color contribution of the last sampling point along r . It is obtained by $\hat{w}_l^r = T_{P_r} (1 - \exp(-\sigma_{P_r} \delta_{P_r}))$.

3.5 Training Procedure

Stage 1: Warmup. Our approach includes two training stages: the warmup stage and the temporal dynamic grounding stage. The purpose of warmup is to provide prior 3D geometry and dynamics information to the next stage and thus reduce the difficulty of dynamic grounding. In this stage, we take T consecutive images $\{\mathbf{I}_t^v\}_{t=1}^T$ uniformly collected by a monocular in 1 second. Each image is taken from a random viewpoint. In addition, we also take only a few multi-view images $\{\mathbf{I}_{t_0}^v\}_{v=1}^V$ of the scene at the first timestamp. Generally, we employ a clustering algorithm to initialize $\mathcal{V}_{\text{density}}^{t_0}$ with N channels based on the grid-level appearance (*e.g.*: geometry, color) and dynamics information. The assumption is that voxels belonging to the same object should be clustered together, exhibiting similar motion and appearance features. Conversely, voxels corresponding to objects in different spatial locations should be separated and exhibit diverse features. Specifically, we first select valid voxels $\{X_k\}_{k=1}^K$ by filtering out invalid locations with density values below a predefined threshold. This filtering step is based on the 3D voxel grids learned in this stage. Subsequently, we construct a feature graph G using these selected voxels. This feature graph incorporates information related to the geometry, color, and dynamics of the voxels. To obtain the dynamics information, we additionally train the f'_ξ module to model the forward deformation field $\Delta x'_{0 \rightarrow t}$. f'_ξ is trained by \mathcal{L}_{Cyc} in Eq. 4. Next, we apply the connected components algorithm on the feature graph G to generate clusters. Preliminary experiments have demonstrated the effectiveness of this method in improving the final performance in the second training stage. The loss function in this stage is defined as $\mathcal{L}_{\text{Warm}} = \sum_{t=1}^T (\mathcal{L}_{\text{Render}} + \alpha_p \mathcal{L}_{\text{Point}} + \alpha_e \mathcal{L}_{\text{Ent}} + \alpha_c \mathcal{L}_{\text{Cyc}})$, where we adopt the empirical values of the hyperparameters from previous literature [18].

Table 1: Novel view synthesis results of our approach compared with D-NeRF [24] and DeVRF [18], as well as their variants for dynamic scenes (see text for details). We evaluate the results averaged over 60 novel views per timestamp.

METHOD	3OBJFALL		3OBJRAND		3OBJMETAL		3FALL+3STILL	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
D-NeRF	28.54	0.946	12.62	0.853	27.83	0.945	24.56	0.908
D-NeRF+STC	29.15	0.954	27.44	0.943	28.59	0.953	25.03	0.913
DeVRF	24.92	0.927	22.27	0.912	25.24	0.931	24.80	0.931
DeVRF-DYN	18.81	0.799	18.43	0.799	17.24	0.769	17.78	0.765
DynaVol	32.11	0.969	30.70	0.964	29.31	0.953	28.96	0.945

METHOD	6OBJFALL		8OBJFALL		3OBJREALSIMP		3OBJREALCMPX	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
D-NeRF	28.27	0.940	27.44	0.923	27.04	0.927	20.73	0.864
D-NeRF+STC	27.20	0.928	26.97	0.919	27.49	0.931	22.72	0.874
DeVRF	24.83	0.905	24.87	0.915	24.81	0.922	21.77	0.891
DeVRF-DYN	17.35	0.738	16.19	0.711	18.64	0.717	17.40	0.778
DynaVol	29.98	0.950	29.78	0.945	30.13	0.952	27.25	0.918

Stage 2: Dynamic grounding. In this stage, we load $\mathcal{V}_{\text{density}}^{t_0}$ and ϕ from the warmup stage. These pretrained parameters introduce valuable grid-level geometry and dynamics priors and facilitate the effective separation of spatial and temporal features in the current stage. The end-to-end optimization throughout the sequential data is beneficial as it enables the refinement of the initial voxel grids with the assistance of object-level representations. The loss function in this training stage is defined as $\mathcal{L}_{\text{Dyn}} = \sum_{t=1}^T (\mathcal{L}_{\text{Render}} + \alpha_p \mathcal{L}_{\text{Point}} + \alpha_e \mathcal{L}_{\text{Ent}})$, where we finetune $(\mathcal{V}_{\text{density}}^{t_0}, \psi)$ from the warmup stage, and train (θ, ω, ϕ) from the scratch.

4 Experiments

4.1 Implementation Details

We set the size of the voxel grid to 110^3 , the assumed number of maximum objects to $N = 10$, and the dimension of slot features to $D = 64$. We use 4 hidden layers with 64 channels in the renderer, and use the Adam optimizer with a batch of 1,024 rays in the two training stages. The base learning rates are 0.1 for the voxel grids and $1e^{-3}$ for all model parameters in the warmup stage and then adjusted to 0.08 and $8e^{-4}$ in the second training stage. The two training stages last for 50k and 35k iterations respectively. The hyperparameters in the loss functions are set to $\alpha_p = 0.1$, $\alpha_e = 0.01$, $\alpha_w = 1.0$, $\alpha_c = 1.0$. All experiments run on an NVIDIA RTX3090 GPU and last for about 3 hours.

4.2 Experimental Setup

We evaluate DynaVol on both scene representation (via scene segmentation) and novel view synthesis in the following 8 scenes. We show its ability on the representative downstream task of dynamic scene editing. For each scene, we capture V -view ($V = 5$) static images of the initial scene and a dynamic sequence with $T = 60$ timestamps which is rendered from viewpoints randomly sampled at different moments from the upper hemisphere as training data and randomly select another different view at each timestamp for the test, all in 512×512 pixels.

Dataset. We build 8 synthetic dynamic scenes using Kubric [7] with different numbers of objects (in different colors and shapes), diverse motions with different initial velocities, different materials, and real-world shapes and textures (All dataset names can be referred from Table 1).

Metrics. For quantitative comparison of the novel view synthesis problem, we report Mean Squared Error (MSE), Peak Signal-to-Noise Ratio (PSNR), and Structural Similarity (SSIM) [31] for synthetic scenes. Additionally, to evaluate segmentation quality in a way that is compatible with 2D methods, we use Foreground Adjusted Rand Index (FG-ARI) as our metric, which measures clustering similarity according to the ground truth foreground objects mask where a random segmentation would score 0 and a perfect segmentation would score 1.

Compared methods. We compare DynaVol with various benchmarks, including 3D scene modeling methods D-NeRF [24] and DeVRF [18], 2D image segmentation methods SAVi [15], SAM [16], and

Table 2: Comparisons with existing approaches based on 2D/3D object-centric representations, *i.e.*, SAVi [5], uORF [34], and SAM [16]. In particular, for uORF, we present novel view synthesis and object segmentation results. For SAVi, since it requires videos with fixed viewpoints, we generate an image sequence at a certain fixed camera position and present the result of SAVi. To compare with it, we evaluate DynaVol (Fix) with images that are also collected at this fixed camera view.

METHOD	3OBJRAND		6OBJFALL		8OBJFALL		3OBJMETAL	
	PSNR \uparrow	FG-ARI \uparrow	PSNR \uparrow	FG-ARI \uparrow	PSNR \uparrow	FG-ARI \uparrow	PSNR \uparrow	FG-ARI \uparrow
SAVi	—	4.38	—	6.85	—	7.87	—	3.38
Ours (Fix)	31.69	84.68	33.11	93.42	31.47	91.22	30.30	94.91
uORF	6.65	38.65	6.63	29.23	6.89	31.93	7.95	22.58
SAM	—	55.52	—	62.66	—	71.68	—	46.80
Ours	30.70	96.01	29.98	94.73	29.78	95.10	29.31	96.06

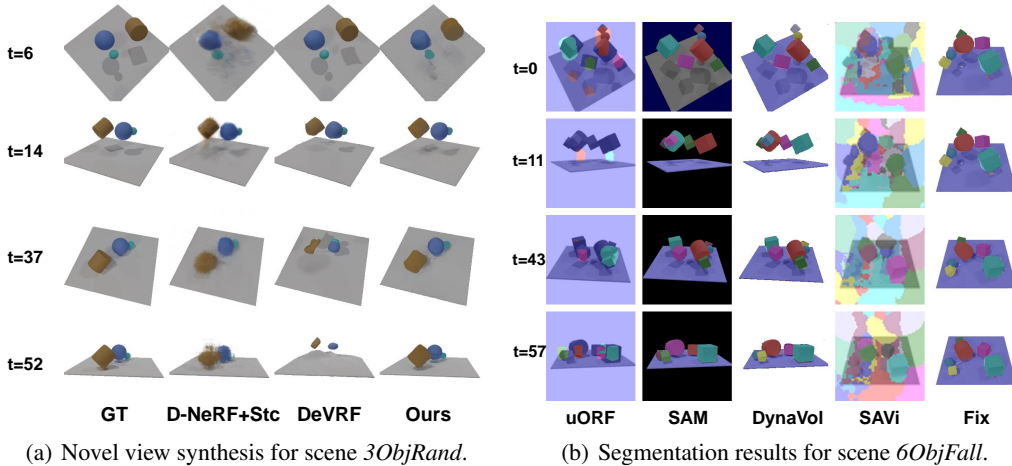


Figure 3: Visualization of novel view synthesis and scene decomposition results.

3D object-centric methods uORF [34]. Since our method also uses static data in synthetic scenes, for a fair comparison, we implement the D-NeRF+Stc which is trained on the static image set and the dynamic sequence. Besides, since DeVRF is trained on the 60-view static image set and 4-view dynamic sequence, we additionally trained DeVRF with the same data setting as ours, *i.e.*, viewpoints are randomly sampled from the upper hemisphere. Both two variants of DeVRF are trained without the optical flow loss proposed in the paper. For SAVi and uORF, we use the models pretrained on MOVi-A [7] and CLEVR-567 [34] respectively, which are similar to our scenes. As for SAM, we employ their pretrained model, which is publicly available and open-sourced. We try to finetune them on our dataset, however, it does not improve the model performance.

4.3 Novel View Synthesis

We evaluate the performance of DynaVol on the novel view synthesis task with other two 3D benchmarks (D-NeRF, DeVRF) and their variants trained on the same data as ours (D-NeRF+Stc and DeVRF-Dyn). As shown in Table 1, DynaVol achieves the best performance in terms of PSNR, SSIM, and MSE in **all** datasets. Notably, our approach outperforms the second-best model by a large margin even in those difficult scenes (*i.e.*, 6ObjFall, 3ObjRealSimp, and 3ObjRealCmpx) with a 15.08% increase in PSNR and a 2.26% increase in SSIM on average. There is no significant difference between D-NeRF and D-NeRF+Stc in most scenes except for 3ObjRand (D-NeRF fails to model 3ObjRand as shown in Figure 4), which illustrates the effectiveness of the use of the static image set. Besides, DeVRF-Dyn has a noticeable decline in performance compared to the standard version due to its heavy dependence on accurate initial scene understanding.

Figure 4 demonstrates qualitative comparisons with other methods and it shows that DynaVol can capture the 3D appearance of different objects and the corresponding motion patterns at an arbitrary timestamp and render a competitive result in a novel view. In contrast, D-NeRF struggles to capture intricate motion patterns, as evidenced by its failure in modeling complex motion in the 3ObjRand

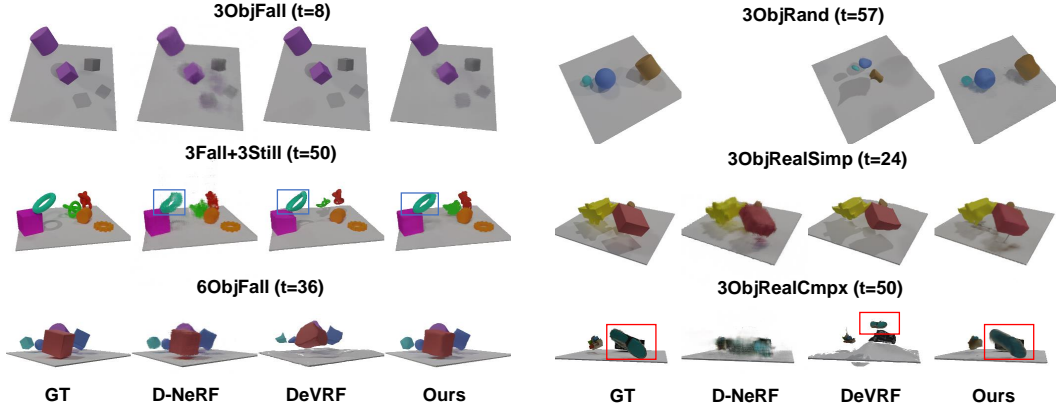


Figure 4: Novel view synthesis for different dynamic scenes. For each scene, we randomly select a novel view at an arbitrary timestamp. Note that the vanilla D-NeRF fails on *3ObjRand*.

scene. Additionally, it performs poorly in accurately representing the appearance of different objects, particularly when dealing with complex textures, as observed in the *3ObjRealCmpx* scene. On the other hand, DeVRF falls short in modeling object dynamics and fails to accurately infer their corresponding spatial locations, as demonstrated in both the *3ObjRand* and *3ObjRealCmpx* scenes.

Figure 3(a) shows a more specific novel view synthetic results on *3ObjRand* scene with timestamps on $t = 6$, $t = 14$, $t = 37$, and $t = 52$ respectively. Considering the synthetic results of D-NeRF on this dataset in Figure 4, we choose its variant as an alternative. It can be found that D-NeRF+Stc produces blurry in all sampled timestamps and DeVRF suffers from severe deformation and position shift in $t = 37$ and $t = 52$, while DynaVol renders relatively clearer and more precise images.

4.4 Scene Decomposition

To get the 2D segmentation results, we assign the rays to different slots according to the contribution of each slot to the final color of the ray. In Table 2 we compare our method with the other three segmentation benchmarks. Since SAVi is a 2D segmentation method that only works on the dynamic scene with a fixed camera position, we evaluate its performance of metrics FG-ARI with DynaVol (labeled as DynaVol(FixCam)) on a fixed one-view dynamic sequence. As for uORF and SAM, they mainly focus on static scenes, so we process the dynamic sequence into T static single scenes as their inputs and evaluate their average performance on the whole sequence. Results show that our method significantly outperforms all approaches, both in reconstruction quality and segmentation results. It is worth mentioning that we take ARI-FG as the segmentation metrics. The higher the ARI-FG, the method has better segmentation results and temporal consistency.

Figure 3(b) shows the results of the qualitative comparison. DynaVol handles well on object occlusion and ensures object-slot correspondence consistency both in 3D and temporal aspects across multiple views (*i.e.* object always has a fixed color). SAVi performs suboptimally in this particular scene, whereas uORF and SAM exhibit inconsistency in both the temporal and spatial dimensions. This inconsistency manifests as the assignment of different slots to the same object at different timestamps or from different viewpoints.

4.5 Dynamic Scene Editing

The object-centric representations learned in DynaVol have practical applications in downstream tasks such as scene editing without requiring additional training. By directly modifying the volumetric representations or altering the deformation function or slot representations, DynaVol enables a range of editing tasks. For instance, in Figure 5(a), DynaVol removes the right hand that is pinching the toys. In Figure 5(b), it modifies the dynamics of the shoes from falling to rotating. Additionally, in Figure 5(c), it replaces the cylinder from the *3ObjFall* scene with the book from the *3ObjRealCmpx* scene. Moreover, object colors can be swapped by exchanging their corresponding slots, as demonstrated in Figure 5(d). This indicates that the model effectively binds the slot features to different objects and learns effective appearance information. For a complete visualization of edited sequences and additional editing tasks, please refer to our supplementary materials.

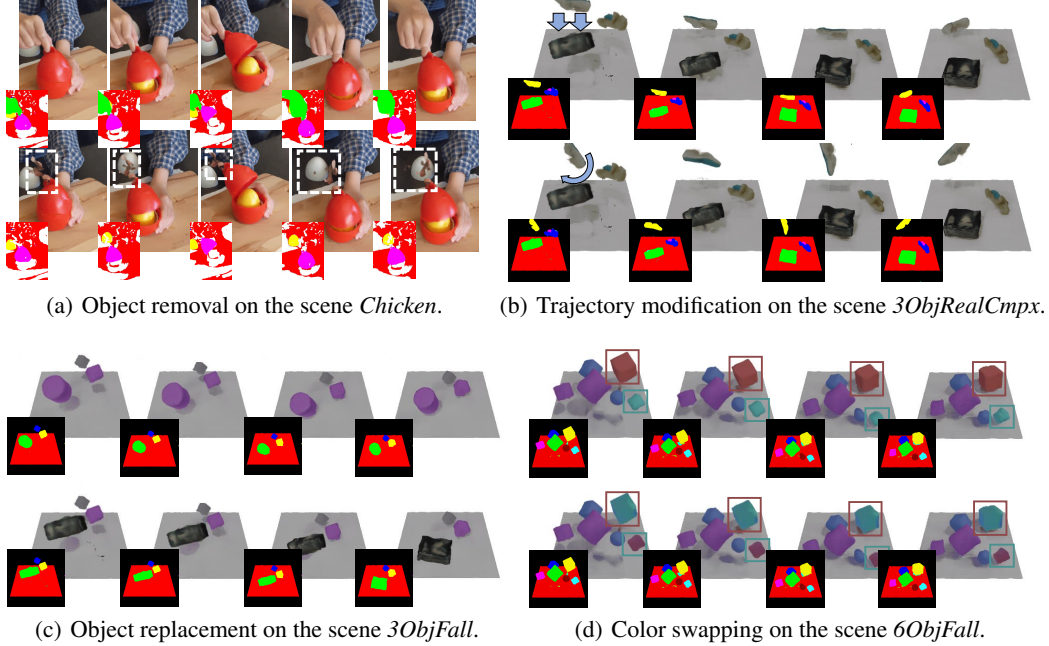


Figure 5: Showcases of dynamic scene editing in both (a) the real-world scene [23], (b-c) the synthetic scenes with real-world object geometry and textures, and (d) the synthetic scene with severe occlusions between objects. The top row in each sub-figure indicates the results of novel view synthesis and scene decomposition. The bottom row indicates the results of scene editing.

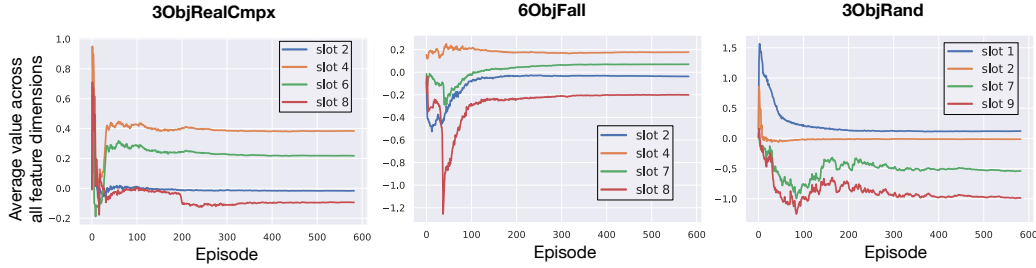


Figure 6: Visualization of the average value across all feature dimensions of \bar{S}_e at different training episodes on *3ObjRealCmpx*, *6ObjFall*, and *3ObjRand*.

4.6 Analysis on Slot Convergence

We explore the convergence of the slot values during the training process. Specifically, we select the four slots (out of ten) that contribute most to image rendering in *3ObjRealCmpx*, *6ObjFall*, and *3ObjRand*. As shown in Figure 6, we present the average value across all dimensions of each slot ($\{\bar{s}_n^e\}$) at different training episodes. The results demonstrate that each slot effectively converges over time to a stable value, indicating that the features become progressively refined and successfully learn time-invariant information about the scene. Furthermore, the noticeable divergence among different slots indicates that DynaVol successfully learned distinct and object-specific features.

5 Conclusion and Limitation

In this paper, we presented DynaVol, an inverse graphics method designed to understand 3D dynamic scenes using object-centric volumetric representations. Our approach demonstrates superior performance over existing techniques in unsupervised scene decomposition in both synthetic and real-world scenarios. Moreover, it goes beyond the 2D counterparts by providing additional capabilities, such as novel view synthesis and dynamic scene editing, which greatly expand its application prospects.

Admittedly, similar to the previous neural rendering technique [34] for static scene decomposition, a notable limitation in DynaVol is its dependence on multi-view images in the warmup stage. We are actively working towards resolving this limitation in our future research.

Acknowledgments

This work was supported by NSFC (62250062, U19B2035, 62106144), Shanghai Municipal Science and Technology Major Project (2021SHZDZX0102), the Fundamental Research Funds for the Central Universities, and Shanghai Sailing Program (21Z510202133) from the Science and Technology Commission of Shanghai Municipality.

References

- [1] Christopher P Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt Botvinick, and Alexander Lerchner. Monet: Unsupervised scene decomposition and representation. In *CVPR*, 2019.
- [2] Chang Chen, Fei Deng, and Sungjin Ahn. ROOTS: Object-centric representation and rendering of 3D scenes. *Journal of Machine Learning Research*, 2021.
- [3] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *EMNLP*, 2014.
- [4] Danny Driess, Zhiao Huang, Yunzhu Li, Russ Tedrake, and Marc Toussaint. Learning multi-object dynamics with compositional neural radiance fields. *arXiv preprint arXiv:2202.11855*, 2022.
- [5] Gamaleldin F. Elsayed, Aravindh Mahendran, Sjoerd van Steenkiste, Klaus Greff, Michael Curtis Mozer, and Thomas Kipf. SAVi++: Towards end-to-end object-centric learning from real-world videos. In *NeurIPS*, 2022.
- [6] Martin Engelcke, Adam R Kosiorek, Oiwi Parker Jones, and Ingmar Posner. Genesis: Generative scene inference and sampling with object-centric latent representations. In *ICLR*, 2020.
- [7] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J. Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, Thomas Kipf, Abhijit Kundu, Dmitry Lagun, Issam H. Laradji, Hsueh-Ti Liu, Henning Meyer, Yishu Miao, Derek Nowrouzezahrai, Cengiz Oztireli, Etienne Pot, Noha Radwan, Daniel Rebain, Sara Sabour, Mehdi S. M. Sajjadi, Matan Sela, Vincent Sitzmann, Austin Stone, Deqing Sun, Suhani Vora, Ziyu Wang, Tianhao Wu, Kwang Moo Yi, Fangcheng Zhong, and Andrea Tagliasacchi. Kubric: A scalable dataset generator. In *CVPR*, 2022.
- [8] Klaus Greff, Raphaël Lopez Kaufman, Rishabh Kabra, Nick Watters, Christopher Burgess, Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-object representation learning with iterative variational inference. In *ICML*, 2019.
- [9] Klaus Greff, Antti Rasmus, Mathias Berglund, Tele Hao, Harri Valpola, and Jürgen Schmidhuber. Tagger: Deep unsupervised perceptual grouping. In *NeurIPS*, 2016.
- [10] Klaus Greff, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. On the binding problem in artificial neural networks. *arXiv preprint arXiv:2012.05208*, 2020.
- [11] Shanyan Guan, Huayu Deng, Yunbo Wang, and Xiaokang Yang. NeuroFluid: Fluid dynamics grounding with particle-driven neural radiance fields. In *ICML*, 2022.
- [12] Xiang Guo, Guanying Chen, Yuchao Dai, Xiaoqing Ye, Jiadai Sun, Xiao Tan, and Errui Ding. Neural deformable voxel grid for fast optimization of dynamic view synthesis. In *ACCV*, 2022.
- [13] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross B. Girshick. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, 2016.
- [14] Rishabh Kabra, Daniel Zoran, Goker Erdogan, Loic Matthey, Antonia Creswell, Matt Botvinick, Alexander Lerchner, and Chris Burgess. SIMONE: View-invariant, temporally-abstracted object representations via unsupervised video decomposition. In *NeurIPS*, 2021.
- [15] Thomas Kipf, Gamaleldin F Elsayed, Aravindh Mahendran, Austin Stone, Sara Sabour, Georg Heigold, Rico Jonschkowski, Alexey Dosovitskiy, and Klaus Greff. Conditional object-centric learning from video. In *ICLR*, 2022.

- [16] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023.
- [17] Yunzhu Li, Shuang Li, Vincent Sitzmann, Pulkit Agrawal, and Antonio Torralba. 3D neural scene representations for visuomotor control. In *CoRL*, 2022.
- [18] Jia-Wei Liu, Yan-Pei Cao, Weijia Mao, Wenqiao Zhang, David Junhao Zhang, Jussi Keppo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. DeVRF: Fast deformable voxel radiance fields for dynamic scenes. In *NeurIPS*, 2022.
- [19] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. In *NeurIPS*, 2020.
- [20] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015.
- [21] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 1995.
- [22] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- [23] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *ACM Trans. Graph.*, 40(6), dec 2021.
- [24] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural radiance fields for dynamic scenes. In *CVPR*, 2020.
- [25] Mehdi SM Sajjadi, Daniel Duckworth, Aravindh Mahendran, Sjoerd van Steenkiste, Filip Pavetić, Mario Lučić, Leonidas J Guibas, Klaus Greff, and Thomas Kipf. Object scene representation transformer. In *NeurIPS*, 2022.
- [26] Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning. In *NeurIPS*, volume 30, 2017.
- [27] Gautam Singh, Yi-Fu Wu, and Sungjin Ahn. Simple unsupervised object-centric learning for complex and naturalistic videos. In *NeurIPS*, 2022.
- [28] Karl Stelzner, Kristian Kersting, and Adam R Kosiorek. Decomposing 3D scenes into objects via unsupervised volume segmentation. *arXiv preprint arXiv:2104.01148*, 2021.
- [29] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022.
- [30] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [31] Zhou Wang, Alan Conrad Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 2004.
- [32] Jiajun Wu, Ilker Yildirim, Joseph J Lim, Bill Freeman, and Josh Tenenbaum. Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. In *NeurIPS*, volume 28, 2015.
- [33] Tianhao Wu, Fangcheng Zhong, Andrea Tagliasacchi, Forrester Cole, and Cengiz Öztireli. D2NeRF: Self-supervised decoupling of dynamic and static objects from a monocular video. In *NeurIPS*, 2022.
- [34] Hong-Xing Yu, Leonidas J Guibas, and Jiajun Wu. Unsupervised discovery of object radiance fields. In *ICLR*, 2022.

Appendix

A Data Description

- **3ObjFall.** The scene consists of two cubes and a cylinder. Initially, these objects are positioned randomly within the scene, and then undergo a free-fall motion along the Z-axis.
- **3ObjRand.** We use random initial velocities along the X and Y axes for each object in *3ObjFall*.
- **3ObjMetal.** We change the material of each object in *3ObjFall* from “Rubber” to “Metal”.
- **3Fall+3Still.** We add another three static objects with complex geometry to *3ObjFall*.
- **6ObjFall & 8ObjFall.** We increase the number of objects in *3ObjFall* to 6 and 8.
- **3ObjRealSimp.** We modify *3ObjFall* with real-world objects that have simple textures.
- **3ObjRealCmpx.** We modify *3ObjFall* with real-world objects that have complex textures.
- **Real-world data.** We take *Chicken* from [23].

B Experimental Details of Scene Decomposition and Editing

Scene decomposition. To get the 2D segmentation results, we assign the rays to different slots according to the contribution of each slot to the final color of the ray. Specifically, suppose σ_{in} is the density of slot n at point i , we have

$$w_{in} = \frac{\sigma_{in}}{\sum_{n=0}^N \sigma_{in}}, \quad \beta_n = \sum_{i=1}^P T_i(1 - \exp(-\sigma_i \delta_i)) w_{in}, \quad (5)$$

where w_{in} is the corresponding density probability and β_n is the color contribution to the final color of slot n . We can then predict the label of 2D segmentation by $\hat{y}(\mathbf{r}) = \operatorname{argmax}_n(\beta_n)$.

Scene editing. The object-centric representations acquired through DynaVol demonstrate the capability for seamless integration into scene editing workflows, eliminating the need for additional training. By manipulating the 4D voxel grid $\mathcal{V}_{\text{density}}$ learned in DynaVol, objects can be replaced, removed, duplicated, or added according to specific requirements. For example, we can switch the color between objects by switching the corresponding slots assigned to each object. Moreover, the deformation field of individual objects can be replaced with user-defined trajectories (*e.g.*, rotations and translations), enabling precise animation of the objects in the scene.

C Further Experimental Results

Ablation study of the loss function. In Table 3, we explore the effectiveness of $\mathcal{L}_{\text{point}}$ in both of the training phases by setting the weight (α_p) of per-point RGB loss as $\alpha_p = 0.0$, $\alpha_p = 0.1$ (DynaVol’s setting), and $\alpha_p = 1.0$. It can be found that $\alpha_p = 0.1$ performs better than $\alpha_p = 0.0$, indicating that a conservative use of $\mathcal{L}_{\text{point}}$ can improve the performance. A possible reason is that it eases the training process by moderately penalizing the discrepancy of nearby sampling points on the same ray. When the weight of $\mathcal{L}_{\text{point}}$ increases to $\alpha_p = 1.0$, the performance of the method decreases significantly. Such a substantial emphasis on $\mathcal{L}_{\text{point}}$ is intuitively unreasonable and can potentially introduce bias to the neural rendering process, thereby negatively impacting the final results.

Table 3: Ablation study on the impact of the $\mathcal{L}_{\text{Point}}$ loss. We present novel view synthesis and decomposition results on three datasets.

α_p	3OBJFALL			3OBJRAND			3OBJREALCMPX		
	PSNR \uparrow	SSIM \uparrow	ARI-FG \uparrow	PSNR \uparrow	SSIM \uparrow	ARI-FG \uparrow	PSNR \uparrow	SSIM \uparrow	ARI-FG \uparrow
0.0	31.49	0.965	95.30	30.46	0.963	95.06	26.84	0.914	95.18
0.1(OURS)	32.11	0.969	96.95	30.70	0.964	96.01	27.25	0.918	95.26
1.0	29.90	0.953	89.33	29.25	0.956	93.65	25.49	0.906	55.52

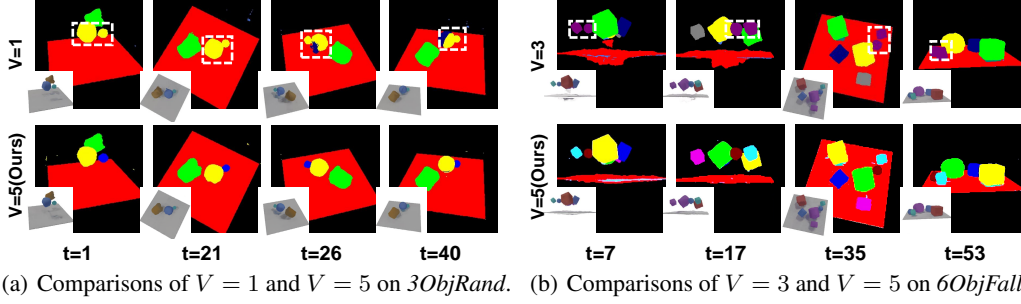


Figure 7: Scene decomposition results of using different numbers ($V = \{1, 3, 5\}$) of images captured from different viewpoints at the initial timestamp.

Table 4: Quantitative results with error bars (mean \pm std) of DynaVol, corresponding to Table 1 and Table 2 in the main manuscript.

3OBJFALL		3OBJRAND		3OBJMETAL		3FALL+3STILL	
PSNR	FG-ARI	PSNR	FG-ARI	PSNR	FG-ARI	PSNR	FG-ARI
32.05 \pm 0.06	97.03 \pm 0.08	30.79 \pm 0.11	96.04 \pm 0.04	29.37 \pm 0.07	96.02 \pm 0.11	28.97 \pm 0.02	94.36 \pm 0.06
6OBJFALL		8OBJFALL		3OBJREALSIMP		3OBJREALCMPX	
PSNR	FG-ARI	PSNR	FG-ARI	PSNR	FG-ARI	PSNR	FG-ARI
29.99 \pm 0.07	94.75 \pm 0.17	29.79 \pm 0.06	95.12 \pm 0.03	30.17 \pm 0.07	94.10 \pm 0.13	27.14 \pm 0.10	95.23 \pm 0.05

Analysis of the initial image set. In general, the performance of the model improves as we incorporate a more diverse set of initial images captured from different viewpoints. In practice, however, obtaining a large number of initial images from various views in a dynamic scene can be challenging. To strike a balance between the difficulty of capturing initial images and the performance of our DynaVol model, we choose the image set size to be $V = 5$. In Figure 7, we validate the effectiveness of different sizes of $\{\mathbf{I}_{t_0}^v\}_{v=1}^V$. We observe that DynaVol achieves better segmentation results at $V = 5$. By contrast, using a smaller number of initial images ($V = 1, 3$) leads to under-segmentation results.

Error bars. To assess the performance stability of DynaVol, we perform separate training processes using three distinct seeds. The results, presented in Table 4, showcase the mean and standard deviation of the PSNR (Peak Signal-to-Noise Ratio) and ARI-FG (Adjusted Rand Index for Foreground) values for novel view synthesis and scene decomposition tasks. These additional results serve as a valuable supplement to those presented in Table 1 and Table 2 in the main manuscript, demonstrating the consistent and reliable performance of DynaVol across multiple training trials.

D Enlarged Visualization of Edited Dynamic Scenes

In Figures 8–11, we provide enlarged visualizations of the edited dynamic scenes for better clarity and observation. These figures provide a closer look at the specific changes made to the dynamic scenes, enabling a better understanding of the editing process.

For more examples, including object adding and duplication, please refer our project page: <https://sites.google.com/view/dynavol/>, which demonstration showcases further instances of scene editing using DynaVol, providing a comprehensive overview of its capabilities.

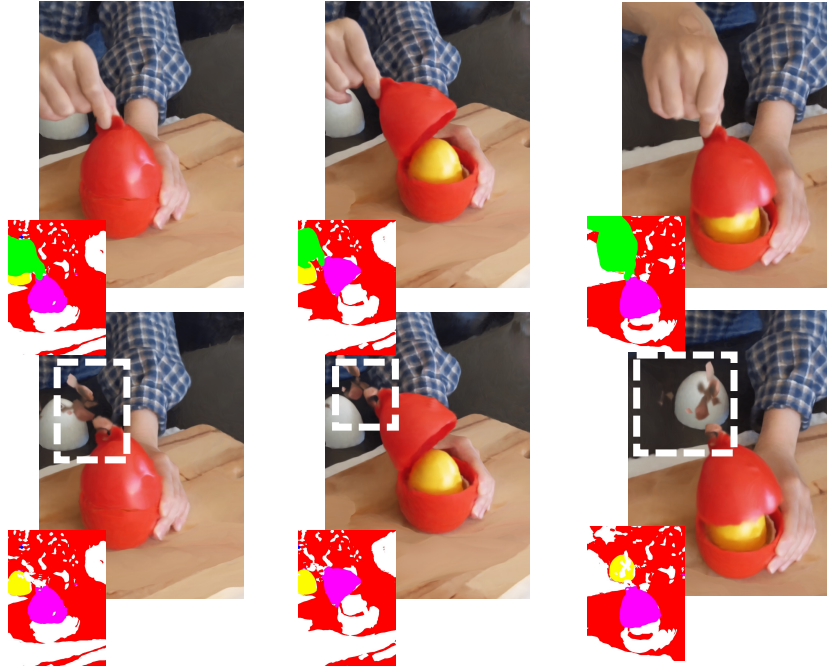


Figure 8: Real-world decomposition and scene editing. The top images illustrate the original scene prior to any editing, while the bottom images present the scene after object removal. It is an enlarged version of Figure 5(a) in the main manuscript.

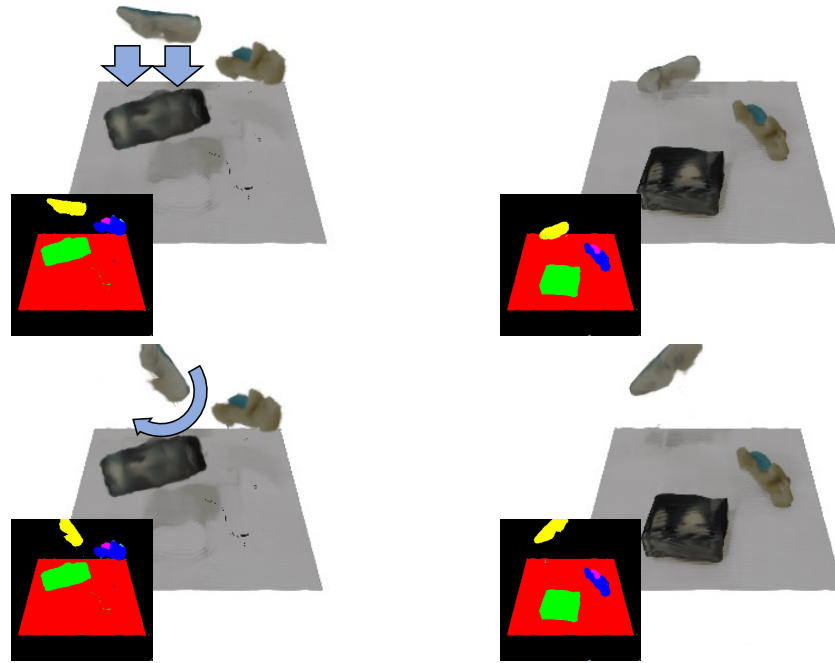


Figure 9: Trajectory modification based on *3ObjRealCmpx* (Top: before editing; Bottom: after editing). It is an enlarged version of Figure 5(b) in the main manuscript.

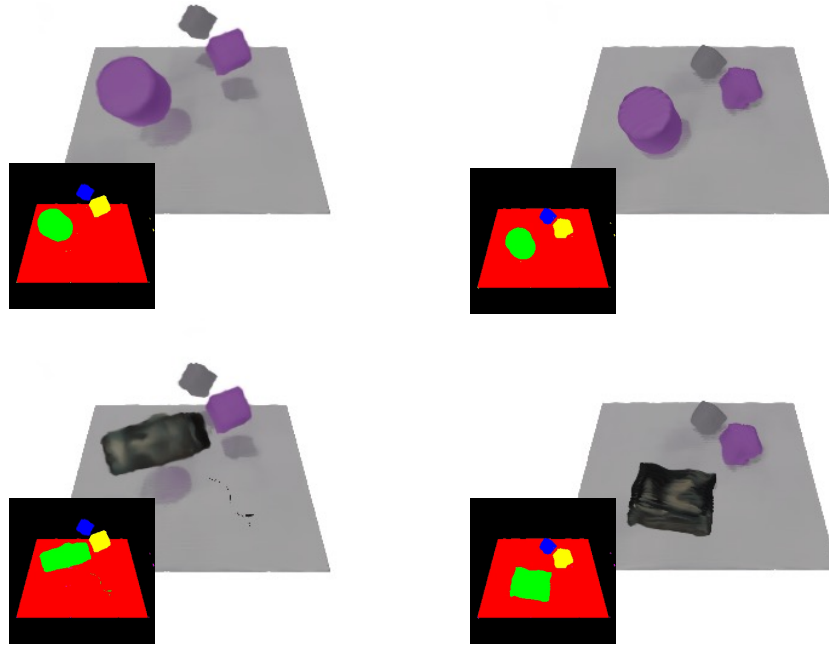


Figure 10: Object replacement based on *3ObjFall* and *3ObjRealCmpx* (Top: before editing; Bottom: after editing). It is an enlarged version of Figure 5(c) in the main manuscript.

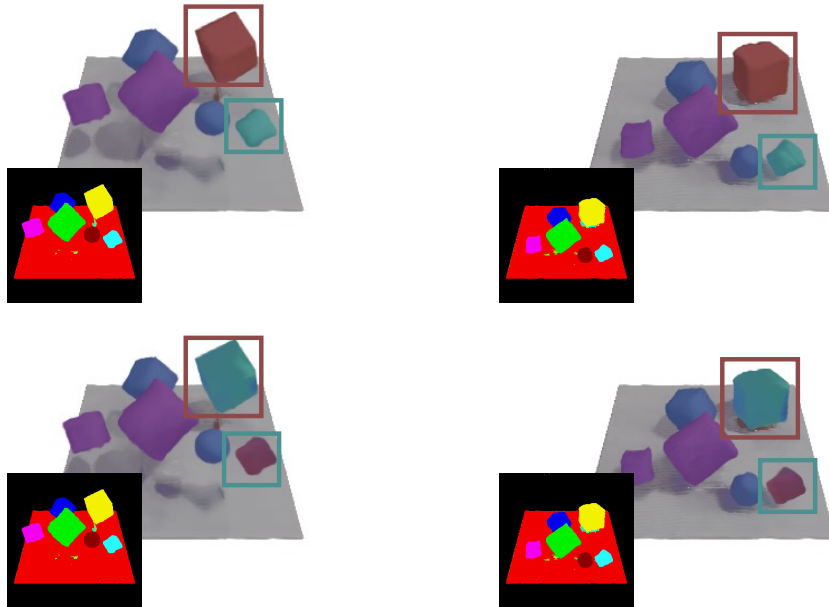


Figure 11: Color swapping based on *6ObjFall* (Top: before editing; Bottom: after editing). It is an enlarged version of Figure 5(d) in the main manuscript.