

DynVideo-E: Harnessing Dynamic NeRF for Large-Scale Motion- and View-Change Human-Centric Video Editing

Jia-Wei Liu^{1*}, Yan-Pei Cao^{3†}, Jay Zhangjie Wu¹, Weijia Mao¹, Yuchao Gu¹, Rui Zhao¹,
Jussi Keppo², Ying Shan³, Mike Zheng Shou^{1†}

¹Show Lab, ²National University of Singapore ³ARC Lab, Tencent PCG



Figure 1. Given a reference subject image and a background style image, our DynVideo-E enables highly consistent editing of large-scale motion- and view-change human-centric videos (a-c).

Abstract

Despite recent progress in diffusion-based video editing, existing methods are limited to short-length videos due to the contradiction between long-range consistency and frame-wise editing. Prior attempts to address this challenge by introducing video-2D representations encounter signifi-

cant difficulties with large-scale motion- and view-change videos, especially in human-centric scenarios. To overcome this, we propose to introduce the dynamic Neural Radiance Fields (NeRF) as the innovative video representation, where the editing can be performed in the 3D spaces and propagated to the entire video via the deformation field. To provide consistent and controllable editing, we propose the image-based video-NeRF editing pipeline with a set of innovative designs, including multi-view multi-pose Score Dis-

* Work is partially done during internship at ARC Lab, Tencent PCG.

† Corresponding Authors.

tillation Sampling (SDS) from both the 2D personalized diffusion prior and 3D diffusion prior, reconstruction losses, text-guided local parts super-resolution, and style transfer. Extensive experiments demonstrate that our method, dubbed as DynVideo-E, significantly outperforms SOTA approaches on two challenging datasets by a large margin of 50% ~ 95% for human preference. Code will be released at <https://showlab.github.io/DynVideo-E/>.

1. Introduction

The remarkable success of powerful image diffusion models [44] has sparked considerable interests in extending them to support video editing [55]. Despite promising, it presents significant challenges in maintaining high temporal consistency. To tackle this problem, existing diffusion-based video editing approaches have evolved to extract and incorporate various correspondences from source videos into the frame-wise editing process, including attention maps [30, 40], spatial maps [57, 64], optical flows and nn-fields [12]. While these works have demonstrated enhanced temporal consistency of editing results, the inherent contradiction between long-range consistency and frame-wise editing limits these methods to short-length videos with small motions and viewpoint changes.

Another line of research seeks to introduce intermediate video-2D representations to degrade video editing to image editing, such as decomposing videos using the layered neural atlas [19] and mapping spatial-temporal contents to 2D UV maps. As such, editing can be performed on a single frame [16, 24] or on the atlas itself [2, 6, 7], with the edited results consistently propagating to other frames. More recently, CoDeF [33] proposes the 2D hash-based canonical image coupled with a 3D deformation field to further improve the video representative capability. However, these approaches are 2D representations of video contents, and thus they encounter significant difficulties in representing and editing videos with large-scale motion and viewpoint changes, especially in human-centric scenarios.

This motivates us to introduce the video-3D representation for large-scale motion- and view-change human-centric video editing. Recent advances in dynamic NeRF [18, 28, 54] show that the 3D dynamic human space coupled with the human pose guided deformation field can effectively reconstruct single human-centric videos with large motions and viewpoints changes. Therefore, in this paper, we propose DynVideo-E that for the first time introduces the dynamic NeRF as the innovative video representation for challenging human-centric video editing. Such a video-NeRF representation effectively aggregates the large-scale motion- and view-change video information into a 3D background space and a 3D dynamic human space through the human pose guided deformation field, and thus the editing

can be performed in the 3D spaces and propagated to the entire video via the deformation field.

To provide consistent and controllable editing, we propose the image-based video-NeRF editing pipeline with a set of effective designs. These include 1) reconstruction losses on the reference image under reference human pose and camera viewpoint to inject subject contents from the reference image to the 3D dynamic human space. 2) To improve the 3D consistency and animatability of the edited 3D dynamic human space, we design a multi-view multi-pose Score Distillation Sampling (SDS) from both the 2D personalized diffusion prior and 3D diffusion prior, as well as a set of training strategies under various human pose and camera pose configurations. 3) To improve the resolution and geometric details of 3D dynamic human space, we utilize the text-guided local parts zoom-in super-resolution with 7 semantic body regions augmented with view conditions. 4) We employ a style transfer module to transfer the reference style to our 3D background model. After training, our video-NeRF model can render highly consistent videos along source video viewpoints by propagating the edited contents through the deformation field, and it can achieve 360° free-viewpoint high-fidelity novel view synthesis for edited dynamic scenes.

We extensively evaluate our DynVideo-E on HOS-NeRF [28] and NeuMan [18] dataset with 24 editing prompts on 11 challenging dynamic human-centric videos. As shown in Fig. 1, our DynVideo-E generates photorealistic video editing results with very high temporal consistency, and significantly outperforms SOTA approaches by a large margin of 50% ~ 95% in terms of human preference.

To summarize, the major contributions of our paper are:

- We present a novel framework of DynVideo-E that for the first time introduces the dynamic NeRF as the innovative video representation for large-scale motion- and view-change human-centric video editing.
- We propose a set of effective designs and training strategies for the image-based 3D dynamic human and static background space editing in our video-NeRF model.
- DynVideo-E significantly outperforms SOTA approaches on two challenging datasets by a large margin of 50% ~ 95% for human preference and achieves high-fidelity free-viewpoint novel view synthesis for edited scenes.

2. Related Work

2.1. Diffusion-based Video Editing

Thanks to the power of diffusion models, prior works have extended their support to video editing [40, 55] and generation [4, 58]. Pioneer Tune-A-Video [55] inflates the image diffusion with cross-frame attention and fine-tunes the source video, aiming to implicitly learn the source motion and transfer it to the target video. Although Tune-A-

Video [55] demonstrates versatility across different video editing applications, it exhibits inferior temporal consistency. Subsequent works extract various correspondences from the source video and employ them to improve temporal consistency. FateZero [40] and Video-P2P [30] extract the cross- and self-attention from the source video to control the spatial layout. Rerender-A-Video [57], ControlVideo [64], and TokenFlow [12] extract and align optical flows, spatial maps, and nn-fields from the source video, resulting in improved consistency of editing results. Although these works have shown promising results, they are typically used in short-form video editing scenarios with small-scale motions and view changes.

Another line of video editing work relies on a powerful video representation, namely, the layered neural atlas [19], as an intermediate editing representation. The layered neural atlas factorizes the input video using a layered presentation and maps the subject and background of all frames to 2D UV maps. Once the layered neural atlas is learned, editing can occur either on keyframes [16, 24] or on the atlas itself [2, 6, 7], and the editing results consistently propagate to other frames. CoDeF [33] incorporates the 3D deformation field with the 2D hash-based canonical image to further improve the video representative capability. However, both the layered neural atlas [19] and canonical image [33] are pseudo-3D representations of video contents, and they encounter difficulties in reconstructing videos with large-scale motion and viewpoint changes.

2.2. Dynamic NeRFs

Remarkable progress has been made in the field of novel view synthesis since the introduction of Neural Radiance Fields (NeRF) [32]. Subsequent studies have extended it to reconstruct dynamic NeRFs from monocular videos by either learning a deformation field that maps sampled points from the deformed space to the canonical space [34, 35, 39, 52] or building 4D spatio-temporal radiance fields [11, 26, 56]. Other studies have introduced voxel grids [9, 27, 51] or planar representations [5, 10] to improve the training efficiency of dynamic NeRFs. While these approaches have shown promising results, they are limited to short videos with simple deformations. Another series of work focus on human modelling and leverage estimated human pose priors [37, 54] to reconstruct dynamic humans with complex motions. Recently, NeuMan [18] reconstructs the dynamic human NeRF together with static scene NeRF to model human-centric scenes. HOSNeRF [28] further proposes to represent the complex human-object-scene with the state-conditional dynamic human model and unbounded background model, achieving 360° free-viewpoint renderings from single videos. In contrast, we aim to introduce the dynamic NeRF as the innovative video-NeRF representation for human-centric video editing.

2.3. NeRF-based Editing and Generation

Since the introduction of diffusion models, text-guided 3D NeRF editing and generation has evolved from CLIP-based [14, 17, 53] to 2D diffusion-based [23, 25, 31, 48, 65] methods. SINE [1] supports editing a local region of static NeRF from a single view by delivering edited contents to multi-views through pretrained NeRF priors. ST-NeRF [59] presents a spatiotemporal neural layered radiance representation to represent dynamic scenes with layered NeRFs, and it can achieve simple editing such as affine transform or duplication by manipulating the NeRF layer. However, it requires 16 cameras to capture a dynamic scene and cannot edit the contents of layered NeRFs. Subsequent works such as Control4D [49] and Dyn-E [63] propose to edit the contents of dynamic NeRFs. However, Control4D [49] is limited to human-only scenes with small motions and short video length, while Dyn-E [63] only supports editing the local appearance with explicit user manipulation.

3. Method

3.1. Video-NeRF Model

Motivation. Given single videos with large viewpoint changes, intricate scene contents, and complex human motions, we seek to represent such videos using dynamic NeRFs for video editing. HOSNeRF [28] has been recently proposed to reconstruct the dynamic neural radiance fields for dynamic human-object-scene interactions from a single monocular in-the-wild video and achieves new SOTA performances. It proposes the state-conditional 3D dynamic human-object model and 3D background model to separately represent the dynamic human-object and static background. Therefore, we harness HOSNeRF [28] as our video-NeRF model to represent large-scale motion- and view-change human-centric videos that consists of dynamic humans, dynamic objects, and static backgrounds. Since our goal is to edit the dynamic human and unbounded background while keep the interacted objects unchanged, we utilize the original reconstructed HOSNeRF model to keep interacted objects, and simplify HOSNeRF [28] to HSNerF by removing the object state designs for video editing. Therefore, our video-NeRF model consists of a dynamic human model Ψ^H and a static scene model Ψ^S .

3D Dynamic Human Model Ψ^H aggregates the dynamic information across all video frames into a 3D canonical human space Ψ_c^H that maps 3D points to color \mathbf{c} and density d , and a human pose-guided deformation field Ψ_d^H that maps deformed points \mathbf{x}_d^i from the deformed space at frame i to canonical points \mathbf{x}_c^i in the canonical space (i omitted for simplicity).

$$\Psi_c^H(\gamma(\mathbf{x}_c)) \mapsto (\mathbf{c}, d), \Psi_d^H(\mathbf{x}_d, \mathcal{J}, \mathcal{R}) \mapsto (\mathbf{x}_c), \quad (1)$$

where $\gamma(\mathbf{x})$ is the standard positional encoding function,

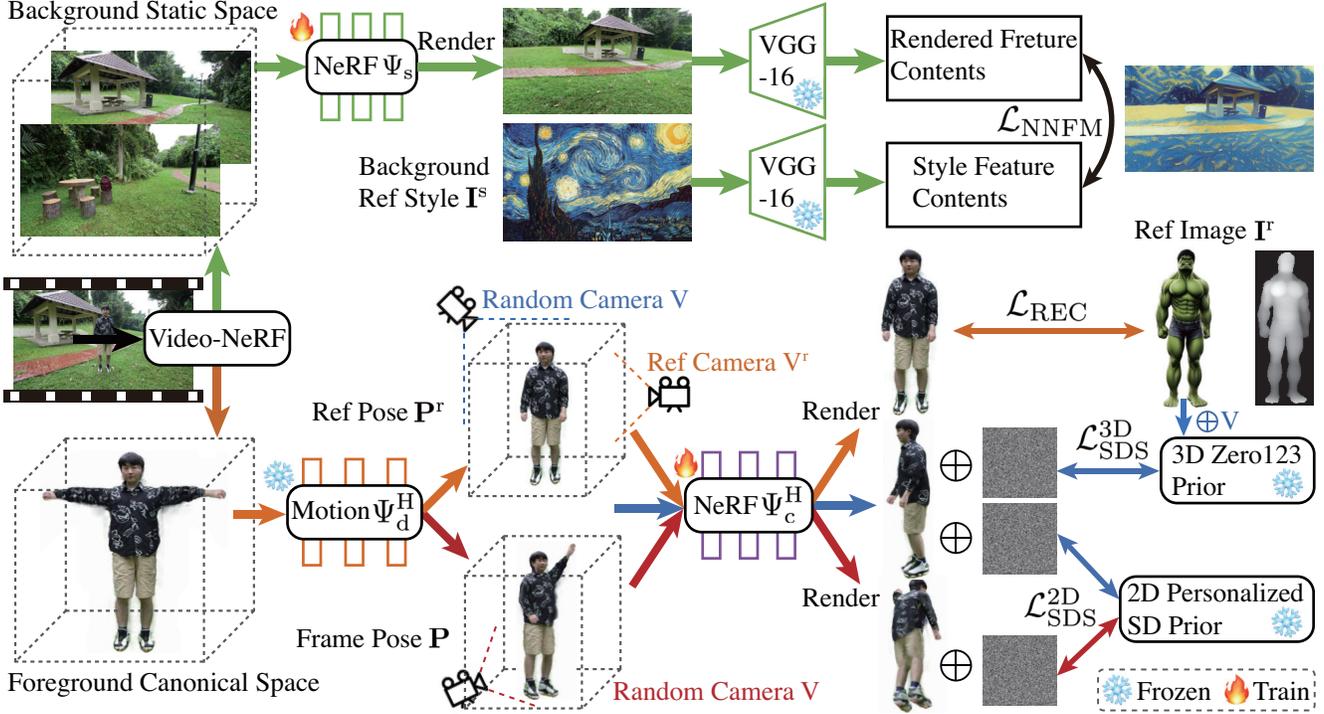


Figure 2. **Overview of DynVideo-E.** (1) Our **video-NeRF model** represents the input video as a 3D dynamic human space coupled with the deformation field and a 3D static background space. (2) **Orange flowchart:** Given the reference subject image, we edit the animatable 3D dynamic human space under multi-view multi-pose configurations by leveraging reconstruction losses, 2D personalized diffusion priors, 3D diffusion priors, and local parts super-resolution. (3) **Green flowchart:** A style transfer loss in feature spaces is utilized to transfer the reference style to our 3D background model. (4) **Edited videos** can be accordingly rendered by volume rendering in the edited video-NeRF model under source video camera poses, and we can also achieve high-fidelity free-viewpoint renderings of edited dynamic scenes.

and $\mathcal{J} = \{\mathbf{J}_i\}$ and $\mathcal{R} = \{\omega_i\}$ are 3D human joints and local joint axis-angle rotations, respectively.

Following HOSNeRF [28] and HumanNeRF [54], we decompose the deformation field Ψ_d^H into a coarse human skeleton-driven deformation $\Psi_d^{H,coarse}$ and a fine non-rigid deformation conditioned on human poses $\Psi_d^{H,fine}$:

$$\mathbf{x}'_c = \Psi_d^{H,coarse}(\mathbf{x}_d, \mathcal{J}, \mathcal{R}), \quad \mathbf{x}_c = \mathbf{x}'_c + \Psi_d^{H,fine}(\mathbf{x}'_c, \mathcal{R}). \quad (2)$$

3D Static Scene Model Ψ^S aggregates intricate static scene contents into a Mip-NeRF 360 [3] space that maps contracted Gaussian parameters $(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$ to color \mathbf{c} and density d .

$$\Psi_s(\hat{\boldsymbol{\gamma}}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})) \mapsto (\mathbf{c}, \sigma), \quad (3)$$

where $\hat{\boldsymbol{\gamma}}$ is the integrated positional encoding (IPE) [3]:

$$\hat{\boldsymbol{\gamma}}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) = \left\{ \begin{array}{l} \sin(2^\ell \hat{\boldsymbol{\mu}}) \exp\left(-2^{2\ell-1} \text{diag}(\hat{\boldsymbol{\Sigma}})\right) \\ \cos(2^\ell \hat{\boldsymbol{\mu}}) \exp\left(-2^{2\ell-1} \text{diag}(\hat{\boldsymbol{\Sigma}})\right) \end{array} \right\}_{\ell=0}^{L-1}. \quad (4)$$

To obtain the contracted Gaussian parameters, we first split the casted rays into a set of intervals $T_i = [t_i, t_{i+1})$ and compute their corresponding conical frustums' mean and covariance as $(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathbf{r}(T_i)$ [3]. Then we adopt the contraction function $f(\mathbf{x})$ proposed in Mip-NeRF 360 [3] to

distribute distant points proportionally to disparity, and parameterize the Gaussian parameters for unbounded scenes as follows,

$$f(\mathbf{x}) = \begin{cases} \mathbf{x} & \|\mathbf{x}\| \leq 1 \\ \left(2 - \frac{1}{\|\mathbf{x}\|}\right) \left(\frac{\mathbf{x}}{\|\mathbf{x}\|}\right) & \|\mathbf{x}\| > 1 \end{cases}, \quad (5)$$

and $f(\mathbf{x})$ is applied to $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ to obtain contracted Gaussian parameters:

$$(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) = \left(f(\boldsymbol{\mu}), \mathbf{J}_f(\boldsymbol{\mu}) \boldsymbol{\Sigma} \mathbf{J}_f(\boldsymbol{\mu})^T\right), \quad (6)$$

where $\mathbf{J}_f(\boldsymbol{\mu})$ is the Jacobian of f at $\boldsymbol{\mu}$.

Video-NeRF Optimization. Given single videos with camera poses calibrated using COLMAP [46, 47], our video-NeRF model is trained by minimizing the difference between the rendered pixel colors and ground-truth pixel colors. To render pixel colors, we shoot rays and query the scene properties in the 3D dynamic human model and scene model, and re-order all sampled properties based on their distances from the camera center. Then, the pixel color can be calculated through the volume rendering [32]:

$$\hat{\mathbf{C}}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - e^{-\sigma_i \delta_i}) \mathbf{c}_i, \quad T_i = e^{-\sum_{j=1}^{i-1} \sigma_j \delta_j}. \quad (7)$$

Following HOSNeRF [28], we optimize our video-NeRF representation by minimizing the photometric MSE loss, patched-based perceptual LPIPS [62] loss, and the regularization losses proposed by Mip-NeRF 360 [3] to avoid background collapse, deformation cycle consistency, and indirect optical flow supervisions. Please refer to HOSNeRF [28] for more details.

3.2. Image-based Video-NeRF Editing

Motivation. Previous video editing works [2, 6, 20, 33, 57] primarily describe intended editing through text prompts. However, finer-grained details and the concept’s identity are better conveyed through reference images. To this end, we focus on the image based editing for finer and direct controllability. As shown in Fig. 2, our video-NeRF model represents the large-scale motion- and view-change human-centric video with a 3D dynamic human space and a 3D background space. Therefore, to better disentangle the foreground and background editing, we propose to edit the 3D dynamic human space with both the reference subject image and its text description, and edit the background static space with the reference style image.

3.2.1 Image-based 3D Dynamic Human Editing

Challenges. Consistent and high-quality image-based video editing requires the edited 3D dynamic human space to 1) keep the subject contents of the reference image; 2) animatable by the human poses from the source video; 3) consistent along large-scale motion and viewpoint changes; and 4) high-resolution with fine details. To address these challenges, we design a set of strategies below.

Reference Image Reconstruction Loss. We utilize a reference subject image \mathbf{I}^r to provide finer identity controls and allow for personalized human editing. To ensure that the reference image has a similar human pose with respect to the source human, we leverage ControlNet [61] to generate the reference subject image conditioned on a source human pose \mathbf{P}^r , as exemplified in Fig. 2. Then, we use a pretrained monocular depth estimator [43] to estimate the pseudo depth \mathbf{D}^r of reference subject and use SAM [22] to obtain its mask \mathbf{M}^r . During training, we assume the reference image viewpoint to be the front view (Ref Camera \mathbf{V}^r in Fig. 2) and render the subject image $\hat{\mathbf{I}}^r$ driven by the source human pose \mathbf{P}^r at \mathbf{V}^r under our video-NeRF representation. We additionally compute the rendered mask $\hat{\mathbf{M}}^r$ and depth $\hat{\mathbf{D}}^r$ at \mathbf{V}^r by integrating the volume density and sampled distances along the ray of each pixel. Following Magic123 [41], we supervise our framework at \mathbf{V}^r viewpoint using the mean squared error (MSE) loss on the reference image and mask, as well as the normalized negative

Pearson correlation on the pseudo depth map.

$$\mathcal{L}_{\text{REC}} = \lambda_{\text{rgb}} \left\| \mathbf{M} \odot \left(\hat{\mathbf{I}}^r - \mathbf{I}^r \right) \right\|_2^2 + \lambda_{\text{mask}} \left\| \hat{\mathbf{M}}^r - \mathbf{M}^r \right\|_2^2 + \frac{1}{2} \lambda_{\text{depth}} \left(1 - \frac{\text{cov} \left(\mathbf{M}^r \odot \mathbf{D}^r, \mathbf{M}^r \odot \hat{\mathbf{D}}^r \right)}{\sigma \left(\mathbf{M}^r \odot \mathbf{D}^r \right) \sigma \left(\mathbf{M}^r \odot \hat{\mathbf{D}}^r \right)} \right) \quad (8)$$

where λ_{rgb} , λ_{mask} , λ_{depth} are the loss weights, \odot is the Hadamard product, $\text{cov}(\cdot)$ is the covariance, and $\sigma(\cdot)$ is the standard deviation.

Score Distillation Sampling (SDS) from 3D Diffusion Prior. Although \mathcal{L}_{REC} can provide supervision on the reference image contents, it only works on the source human pose \mathbf{P}^r at the reference view \mathbf{V}^r . To provide more 3D supervision from the reference image, we utilize the Zero-1-to-3 [29] pretrained on the Objaverse-XL [8] as the 3D diffusion prior to distill the inherent 3D geometric and texture information from the reference image using the SDS loss [38]. Given the 3D diffusion model ϕ with the noise prediction network $\epsilon_\phi(\cdot)$, the SDS loss works by directly minimizing the injected noise ϵ added to the encoded rendered images \mathbf{I} and the predicted noise. Therefore, we render images \mathbf{I} from the 3D dynamic human space driven by the source human pose \mathbf{P}^r at random camera viewpoints $\mathbf{V} = [\mathbf{R}, \mathbf{T}]$, and the SDS loss of Zero-1-to-3 [29] can be computed with the reference image \mathbf{I}^r and the camera pose $[\mathbf{R}, \mathbf{T}]$ as conditions:

$$\nabla_\theta \mathcal{L}_{\text{SDS}}^{\text{3D}}(\phi, F_\theta) = \lambda_{\text{3D}} \cdot \mathbb{E}_{t, \epsilon} \left[w(t) \left(\epsilon_\phi(\mathbf{z}_t; \mathbf{I}^r, t, \mathbf{R}, \mathbf{T}) - \epsilon \right) \frac{\partial \mathbf{I}}{\partial \theta} \right] \quad (9)$$

where \mathbf{z}_t is the noised latent image by injecting a random Gaussian noise of level t to the encoded rendered images \mathbf{I} . $w(t)$ is a weighting function that depends on the noise level t . θ is the optimizable parameters of our DynVideo-E.

SDS from 2D Personalized Diffusion Prior. The reference image guided supervisions above are limited to edit the 3D human space driven only by the source human pose \mathbf{P}^r , and thus are not sufficient to produce a satisfactory 3D dynamic human space that can be animated by the frame human poses from source videos. To this end, we further animate the 3D dynamic human space with the frame human poses \mathbf{P} from the source video and render images \mathbf{I} at random camera poses \mathbf{V} , and we further propose to use the 2D text-based diffusion prior [44] to guide these rendered views. However, naively using the 2D diffusion prior hinders the personalization contents learned from the reference image because the 2D diffusion prior tends to imagine the subject contents purely from text descriptions, as validated in Fig. 4. To solve this problem, we further propose to use 2D personalized diffusion prior that is first finetuned

on the reference image using Dreambooth-LoRA [15, 45]. To generate more inputs for Dreambooth-LoRA, we augment the reference image with random backgrounds and use Magic123 [41] to augment reference image with multiple views. With such Dreambooth-LoRA finetuned 2D personalized diffusion prior ϕ' with the noise prediction network $\epsilon_{\phi'}(\cdot)$, we further employ the 2D SDS loss to supervise the rendered images \mathbf{I} of the 3D dynamic human space animated by the human poses from the source video and rendered at random camera poses.

$$\nabla_{\theta} \mathcal{L}_{\text{SDS}}^{2\text{D}}(\phi', F_{\theta}) = \lambda_{2\text{D}} \mathbb{E}_{t, \epsilon} \left[w(t) (\epsilon_{\phi'}(\mathbf{z}_t; y, t) - \epsilon) \frac{\partial \mathbf{I}}{\partial \theta} \right] \quad (10)$$

where y is the text embedding.

Text-guided Local Parts Super-Resolution. Due to the GPU memory limitation, our DynVideo-E is trained with (128×128) resolutions, resulting in coarse geometry and blurry textures. To solve this problem, inspired by DreamHuman [23], we utilize the text-guided local parts super-resolution to render and supervise the local parts of zoom-in humans, which improves the effective resolution. Because our dynamic human model is a human pose-driven 3D canonical space under the ‘‘T-pose’’ configuration, we can accurately render the zoom-in local human body parts by directly locating the camera close to the corresponding parts. Specifically, we utilize 7 semantic regions: full body, head, upper body, midsection, lower body, left arm, and right arm, and we accordingly modify the input text prompt with body parts and additionally augment these prompts with view-conditional prompts: front view, side view, and back view. Since it is difficult to track the arm’s position under all human poses due to occlusions, we only zoom in on the arms under the ‘‘T-pose’’. We provide 8 visualization examples of text-guided local parts super-resolution sampled during training in supplementary materials.

Dynamic Objects. For human-centric videos with dynamic interacted objects, we utilize the original reconstructed HOSNeRF model to render the interacted objects. During inference, we query the original HOSNeRF model for the rays within object masks, and query the edited video-NeRF model for the rays outside the object masks. As such, we can maintain the dynamic objects in our edited videos.

3.2.2 Image-based 3D Background Editing

We aim at transferring the artistic features of an arbitrary 2D reference style image to our 3D unbounded scene model. As shown in the green flowchart of Fig. 2, we take inspiration from ARF [60] and adopt its nearest neighbor feature matching (NNFM) style loss to transfer the semantic visual details from the 2D reference image \mathbf{I}^s to our 3D background model Ψ^s . We additionally utilize the deferred back-propagation [60] to directly optimize our model on

full-resolution renderings. Specifically, we render the background images \mathbf{I} and extract the VGG [50] feature maps \mathbf{F} and \mathbf{F}^s for \mathbf{I} and \mathbf{I}^s , respectively, and $\mathcal{L}_{\text{NNFM}}$ minimizes the cosine distance between the rendered feature map and its nearest neighbor in the reference feature map.

$$\mathcal{L}_{\text{NNFM}} = \lambda_{\text{NNFM}} \cdot \frac{1}{N} \sum_{i,j} \min_{i',j'} D(\mathbf{F}(i,j), \mathbf{F}^s(i',j')) \quad (11)$$

$$D(\mathbf{v}_1, \mathbf{v}_2) = 1 - \frac{\mathbf{v}_1^T \mathbf{v}_2}{\sqrt{\mathbf{v}_1^T \mathbf{v}_1 \mathbf{v}_2^T \mathbf{v}_2}} \quad (12)$$

To prevent the 3D scene model from deviating much from the source contents, we also add an additional L2 loss penalizing the difference between \mathbf{F} and \mathbf{F}^s [60].

3.3. Training Objectives

The training of DynVideo-E consists of 2 stages. Firstly, we reconstruct our video-NeRF model on the source video. Secondly, we edit the 3D dynamic human space and 3D unbounded scene space given reference images and text prompts. After training, we render the edited videos using our edited video-NeRF model along source video camera viewpoints, and we also achieve free-viewpoint renderings. **Multi-view Multi-pose Training for 3D Dynamic Human Spaces.** As shown in Fig. 2, we design a multi-view multi-pose training process with three conditions during training.

- Orange flowchart in Fig. 2: only \mathcal{L}_{REC} is used to supervise the rendered images under source human pose \mathbf{P}^r at the reference camera view \mathbf{V}^r .
- Blue flowchart in Fig. 2: $\mathcal{L}_{\text{SDS}}^{3\text{D}}$ and $\mathcal{L}_{\text{SDS}}^{2\text{D}}$ are jointly used to supervise the rendered images under source human pose \mathbf{P}^r at random camera view \mathbf{V} .
- Red flowchart in Fig. 2: only $\mathcal{L}_{\text{SDS}}^{2\text{D}}$ is used to supervise the rendered images under the frame human pose \mathbf{P} from the source video at random camera view \mathbf{V} .

4. Experiments

Dataset. To evaluate our DynVideo-E on both long and short videos, we utilize HOSNeRF [28] dataset with [300, 400] frames per video and NeuMan [28] dataset with [30, 90] frames per video, all at a resolution of (1280×720) . In total, we design 24 editing prompts on 11 challenging dynamic human-centric videos to evaluate our DynVideo-E and all SOTA Approaches.

4.1. Comparisons with SOTA Approaches

Baselines. We compare our method against five SOTA approaches, including Text2Video-Zero [20], Rerender-A-Video [57], Text2LIVE [2], StableVideo [6], and CoDeF [33]. We utilize Midjourney* to generate the text descriptions of the reference images to train these baselines.

* <https://www.midjourney.com/>



Figure 3. Qualitative comparisons of DynVideo-E against SOTA approaches on the Backpack scene (a) and Jogging scene (b).

Qualitative Results. We present a visual comparison of our approach against all baselines in Fig. 3 for a long video (a) and a short video (b). Since both videos contain large motions and viewpoint changes, all baselines fail to edit the foreground or background and their results cannot preserve consistent structures. In contrast, our DynVideo-E produces high-quality edited videos that accurately edit both the foreground subject and background style and maintains high temporal consistency, which largely outperforms SOTA approaches. We provide more visual comparisons of all methods, the editing time comparison of all methods, and video comparisons of all methods on 24 editing prompts in the supplementary material.

It is worth noting that for challenging videos with large-scale motions and viewpoint changes, CoDeF [33], Text2LIVE [2], and StableVideo [6] largely overfit to input video frames and learn meaningless canonical images or neural atlas, and thus cannot generate meaningful editing results. We show examples of their learned canonical images and neural atlas in the supplementary material.

Quantitative Results. We quantify our method against baselines through standard metrics and human preferences. We measure the textual faithfulness by computing the average CLIPScore [13] between all frames of output edited videos and corresponding text descriptions. As shown in Tab. 1, our DynVideo-E achieves the highest textual faithfulness score among all approaches.

Human Preference. We show the pairwise comparing videos and textual descriptions to raters, and ask them to select their preference videos in terms of textual faithfulness, temporal consistency, and overall quality. We utilize Amazon MTurk[†] to recruit 10 participants for each comparison (Each comparison may recruit different raters), and compute their preferences over all comparisons on 24 editing prompts. For each comparison, we show our result and one baseline result (shuffled order in questionnaires), together with textual descriptions to raters and ask their preferences. In total, we collected 1140 comparisons over all pairwise results from 32 different raters. As shown in Tab. 1, we report the comparison “ $p_1\%$ v.s. $p_2\%$ ” where p_1 represents the percentage of a baseline is preferred and p_2 denotes our method is preferred. As evident in Tab. 1, our method achieves the highest human preference in all aspects and outperforms all baselines by a large margin of 50% ~ 95%.

4.2. Ablation Study

We conduct ablation studies on 2 videos from HOSNeRF dataset [28] and NeuMan dataset [18]. To evaluate the effectiveness of each proposed component in DynVideo-E, we progressively ablate each component from local parts super-resolution, reconstruction loss, 2D personalized SDS, 3D SDS, and 2D personalization LoRA. To provide the

[†] <https://requester.mturk.com/>

	METRICS	HUMAN PREFERENCE		
	CLIPScore (↑)	Textual Faithfulness (↑)	Temporal Consistency (↑)	Overall Quality (↑)
Text2Video-Zero [20]	26.70	9.17 v.s. 90.83 (Ours)	21.25 v.s. 78.75 (Ours)	12.08 v.s. 87.92 (Ours)
Rerender-A-Video [57]	26.11	6.67 v.s. 93.33 (Ours)	25.00 v.s. 75.00 (Ours)	9.58 v.s. 90.42 (Ours)
Text2LIVE [2]	22.77	3.81 v.s. 96.19 (Ours)	26.67 v.s. 73.33 (Ours)	9.05 v.s. 90.95 (Ours)
StableVideo [6]	22.02	4.29 v.s. 95.71 (Ours)	24.29 v.s. 75.71 (Ours)	6.19 v.s. 93.81 (Ours)
CoDeF [33]	16.77	1.25 v.s. 98.75 (Ours)	3.75 v.s. 96.25 (Ours)	1.25 v.s. 98.75 (Ours)
DynVideo-E (Ours)	31.31	–	–	–

Table 1. Quantitative comparisons of our DynVideo-E against SOTA approaches on HOSNeRF dataset [28] and NeuMan dataset [18].



Figure 4. Qualitative ablation results of our method on each proposed component for (a) Backpack scene and (b) Lab scene.

Ablation components	Backpack Lab	
Full model	0.756	0.647
w/o Super-solution	0.736	0.645
w/o Super-solution, Rec	0.728	0.617
w/o Super-solution, Rec, 2D SDS	0.679	0.517
w/o Super-solution, Rec, 3D SDS	0.711	0.613
w/o Super-solution, Rec, 3D SDS, 2D LoRA	0.698	0.539

Table 2. Quantitative ablation results of our method for the Backpack and Lab scene (higher score means better performance).

quantitative results of our ablation study, we compute the average cosine similarity between the CLIP [42] image embeddings of all frames of output edited videos and the corresponding reference subject image. As evident in Tab. 2, the CLIP score progressively drops with the disabling of each component, with the full model achieving the best performances, which clearly demonstrates the effectiveness of our designs. In addition, we provide the qualitative results of our ablations in Fig. 4, which further demonstrates the effectiveness of our designs. More ablation results on more videos are provided in the supplementary material.

5. Conclusion

We introduced a novel framework of DynVideo-E to consistently edit large-scale motion- and view-change human-centric videos. We first proposed to harness dynamic NeRF as our innovative video representation where the editing can be performed in dynamic 3D spaces and accurately propagated to the entire video via deformation fields. Then, we proposed a set of effective image-based video-NeRF editing designs, including multi-view multi-pose Score Distillation Sampling (SDS) from both the 2D personalized diffusion prior and 3D diffusion prior, reconstruction losses on the reference image, text-guided local parts super-resolution, and style transfer for 3D background spaces. Finally, extensive experiments demonstrated DynVideo-E produced significant improvements over SOTA approaches.

Limitations and Future Work. Although DynVideo-E achieves remarkable progress in video editing, its NeRF-based representation is time-consuming. Using voxel or hash grid in the video-NeRF model can largely reduce the training time and we leave it as a faithful future direction.

References

- [1] Chong Bao, Yinda Zhang, Bangbang Yang, Tianxing Fan, Zesong Yang, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. Sine: Semantic-driven nerf editing with prior-guided editing field. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20919–20929, 2023. 3
- [2] Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kasten, and Tali Dekel. Text2live: Text-driven layered image and video editing. In *European conference on computer vision*, pages 707–723. Springer, 2022. 2, 3, 5, 6, 7, 8, 12, 13, 15
- [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. 4, 5
- [4] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023. 2
- [5] Ang Cao and Justin Johnson. Hexplane: a fast representation for dynamic scenes. *arXiv preprint arXiv:2301.09632*, 2023. 3
- [6] Wenhao Chai, Xun Guo, Gaoang Wang, and Yan Lu. Stable-video: Text-driven consistency-aware diffusion video editing. *arXiv preprint arXiv:2308.09592*, 2023. 2, 3, 5, 6, 7, 8, 12, 13, 15
- [7] Paul Couairon, Clément Rambour, Jean-Emmanuel Haugeard, and Nicolas Thome. Videdit: Zero-shot and spatially aware text-driven video editing. *arXiv preprint arXiv:2306.08707*, 2023. 2, 3
- [8] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. Objaverse-xl: A universe of 10m+ 3d objects. *arXiv preprint arXiv:2307.05663*, 2023. 5
- [9] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. 3
- [10] Sara Fridovich-Keil, Giacomo Meanti, Frederik Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. *arXiv preprint arXiv:2301.10241*, 2023. 3
- [11] 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. 3
- [12] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arXiv:2307.10373*, 2023. 2, 3
- [13] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. 7
- [14] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. Avatarclip: Zero-shot text-driven generation and animation of 3d avatars. *arXiv preprint arXiv:2205.08535*, 2022. 3
- [15] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 6
- [16] Jiahui Huang, Leonid Sigal, Kwang Moo Yi, Oliver Wang, and Joon-Young Lee. Inve: Interactive neural video editing. *arXiv preprint arXiv:2307.07663*, 2023. 2, 3
- [17] Ajay Jain, Ben Mildenhall, Jonathan T Barron, Pieter Abbeel, and Ben Poole. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 867–876, 2022. 3
- [18] Wei Jiang, Kwang Moo Yi, Golnoosh Samei, Oncel Tuzel, and Anurag Ranjan. Neuman: Neural human radiance field from a single video. In *European Conference on Computer Vision*, pages 402–418. Springer, 2022. 2, 3, 7, 8, 12, 13, 15
- [19] 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. 2, 3
- [20] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023. 5, 6, 8, 13, 15
- [21] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 12
- [22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 5
- [23] Nikos Kolotouros, Thiemo Alldieck, Andrei Zanfir, Eduard Gabriel Bazavan, Mihai Fieraru, and Cristian Sminchisescu. Dreamhuman: Animatable 3d avatars from text. *arXiv preprint arXiv:2306.09329*, 2023. 3, 6
- [24] Yao-Chih Lee, Ji-Ze Genevieve Jang, Yi-Ting Chen, Elizabeth Qiu, and Jia-Bin Huang. Shape-aware text-driven layered video editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14317–14326, 2023. 2, 3
- [25] Yuhan Li, Yishun Dou, Yue Shi, Yu Lei, Xuanhong Chen, Yi Zhang, Peng Zhou, and Bingbing Ni. Focaldreamer: Text-driven 3d editing via focal-fusion assembly. *arXiv preprint arXiv:2308.10608*, 2023. 3
- [26] 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*, pages 6498–6508, 2021. 3

- [27] 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. *arXiv preprint arXiv:2205.15723*, 2022. 3
- [28] Jia-Wei Liu, Yan-Pei Cao, Tianyuan Yang, Eric Zhongcong Xu, Jussi Keppo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Hosnerf: Dynamic human-object-scene neural radiance fields from a single video. *arXiv preprint arXiv:2304.12281*, 2023. 2, 3, 4, 5, 6, 7, 8, 12, 15
- [29] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. *arXiv preprint arXiv:2303.11328*, 2023. 5
- [30] Shaoteng Liu, Yuechen Zhang, Wenbo Li, Zhe Lin, and Jiaya Jia. Video-p2p: Video editing with cross-attention control. *arXiv preprint arXiv:2303.04761*, 2023. 2, 3
- [31] Aryan Mikaeili, Or Perel, Daniel Cohen-Or, and Ali Mahdavi-Amiri. Sked: Sketch-guided text-based 3d editing. *arXiv preprint arXiv:2303.10735*, 2023. 3
- [32] 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. 3, 4
- [33] Hao Ouyang, Qiuyu Wang, Yuxi Xiao, Qingyan Bai, Juntao Zhang, Kecheng Zheng, Xiaowei Zhou, Qifeng Chen, and Yujun Shen. Codef: Content deformation fields for temporally consistent video processing. *arXiv preprint arXiv:2308.07926*, 2023. 2, 3, 5, 6, 7, 8, 12, 13, 15
- [34] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021. 3
- [35] 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. 3
- [36] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 12
- [37] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9054–9063, 2021. 3
- [38] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 5
- [39] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. 3
- [40] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing. *arXiv preprint arXiv:2303.09535*, 2023. 2, 3
- [41] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skokhodov, Peter Wonka, Sergey Tulyakov, et al. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. *arXiv preprint arXiv:2306.17843*, 2023. 5, 6
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 8
- [43] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020. 5
- [44] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 5
- [45] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 6
- [46] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 4
- [47] 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. 4
- [48] Etai Sella, Gal Fiebelman, Peter Hedman, and Hadar Averbuch-Elor. Vox-e: Text-guided voxel editing of 3d objects. *arXiv preprint arXiv:2303.12048*, 2023. 3
- [49] Ruizhi Shao, Jingxiang Sun, Cheng Peng, Zerong Zheng, Boyao Zhou, Hongwen Zhang, and Yebin Liu. Control4d: Dynamic portrait editing by learning 4d gan from 2d diffusion-based editor. *arXiv preprint arXiv:2305.20082*, 2023. 3
- [50] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 6
- [51] Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. Nerfplayer: A streamable dynamic scene representation with decomposed neural radiance fields. *arXiv preprint arXiv:2210.15947*, 2022. 3

- [52] 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. 3
- [53] Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Nerf-art: Text-driven neural radiance fields stylization. *IEEE Transactions on Visualization and Computer Graphics*, 2023. 3
- [54] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16210–16220, 2022. 2, 3, 4
- [55] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei, Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. *arXiv preprint arXiv:2212.11565*, 2022. 2, 3
- [56] 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. 3
- [57] Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change Loy. Rerender a video: Zero-shot text-guided video-to-video translation. *arXiv preprint arXiv:2306.07954*, 2023. 2, 3, 5, 6, 8, 13, 15
- [58] David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation, 2023. 2
- [59] Jiakai Zhang, Xinhang Liu, Xinyi Ye, Fuqiang Zhao, Yan-shun Zhang, Minye Wu, Yingliang Zhang, Lan Xu, and Jingyi Yu. Editable free-viewpoint video using a layered neural representation. *ACM Transactions on Graphics (TOG)*, 40(4):1–18, 2021. 3
- [60] Kai Zhang, Nick Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snavely. Arf: Artistic radiance fields. In *European Conference on Computer Vision*, pages 717–733. Springer, 2022. 6
- [61] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 5
- [62] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 5
- [63] Shangzhan Zhang, Sida Peng, Yinji ShenTu, Qing Shuai, Tianrun Chen, Kaicheng Yu, Hujun Bao, and Xiaowei Zhou. Dyn-e: Local appearance editing of dynamic neural radiance fields. *arXiv preprint arXiv:2307.12909*, 2023. 3
- [64] Min Zhao, Rongzhen Wang, Fan Bao, Chongxuan Li, and Jun Zhu. Controlvideo: Adding conditional control for one shot text-to-video editing. *arXiv preprint arXiv:2305.17098*, 2023. 2, 3
- [65] Jingyu Zhuang, Chen Wang, Lingjie Liu, Liang Lin, and Guanbin Li. Dreameditor: Text-driven 3d scene editing with neural fields. *arXiv preprint arXiv:2306.13455*, 2023. 3

Appendix

The supplementary material is structured as follows:

- Sec. A presents implementation details on the network designs and optimization parameters of DynVideo-E.
- Sec. B summarizes additional comparisons and ablations of our DynVideo-E against SOTA approaches.

Furthermore, we provide a **supplementary video** showcasing all 24 edited video comparisons of our method against baselines, as well as 360° free-viewpoint renderings of edited dynamic scenes from our DynVideo-E.

A. Implementation Details

DynVideo-E Network Details. As shown in Fig. 5, we employ a 10-layer multilayer perceptron (MLP) as our state-conditional background network (a) and a 8-layer MLP as our state-conditional canonical human-object network (b). To edit the dynamic human, we establish a 9-layer canonical human network (c) where the parameters of its first 8 layers are initialized from the reconstructed human-object model (b). During optimization, we train the 3D background model (a) and the 3D dynamic human model (c) while freeze the reconstructed dynamic human-object model (b). During inference, for the source video that contains dynamic objects, we query the original dynamic human-object model (b) for the pixels within the object masks to keep the dynamic objects, while we query the edited dynamic human model (c) and edited background model for other pixels to obtain the colors and densities of edited contents. For the human-background videos, we only need to query the edited dynamic human model and edited background model to obtain the edited contents.

Optimization Parameters. We optimize our DynVideo-E using Adam optimizer [21]. We set the learning rate for our training process as 0.0005 with 20000 training iterations. We balance the loss terms using the following weighting factors: $\lambda_{\text{rgb}} = 5$, $\lambda_{\text{mask}} = 0.5$, $\lambda_{\text{depth}} = 0.01$, $\lambda_{3\text{D}} = 40$, $\lambda_{2\text{D}} = 1.0$, $\lambda_{\text{NNFM}} = 1.0$. The guidance scale of the 3D diffusion prior and 2D personalized diffusion prior are set to 5 and 20, respectively. We conducted all our experiments on 1 NVIDIA A100 GPU, using the PyTorch [36] deep learning framework.

Visualization of Text-guided Local Parts Super-Resolution. To improve the effective resolution during training, we utilize the text-guided local parts super-resolution to render and supervise the local parts of zoom-in humans and augment with view-conditional prompts. We provide 8 visualization examples of text-guided local parts super-resolution sampled during training in Fig. 6. As shown in Fig. 6, even though all figures are rendered in (128×128) resolutions, rendering local parts can largely improve the effective resolution and thus we can supervise the detailed geometry and textures of edited

Ablation components	Average CLIP Score
Full model	0.674
w/o Super-solution	0.659
w/o Super-solution, Rec	0.650
w/o Super-solution, Rec, 2D SDS	0.572
w/o Super-solution, Rec, 3D SDS	0.641
w/o Super-solution, Rec, 3D SDS, 2D LoRA	0.593

Table 3. Averaged quantitative ablation results of our method. human body with diffusion priors.

B. Additional Results

More Qualitative Results. We present two more visual comparisons of our approach against all baselines in Fig. 7 and Fig. 8. As shown in the figures, our DynVideo-E achieves the best performances with photo-realistic edited videos, which clearly demonstrates the superiority of our model against other approaches on editing large-scale motion- and view-change human-centric videos. Comparing the long (a) and short (b) video editing results of Fig. 7, we find that baseline approaches perform better on short videos than long videos, but still none of them can edit the correct subject “Thanos” due to the large subject motions and viewpoint changes in videos. In contrast, our DynVideo-E produces high-quality editing results on both short and long videos. Please refer to our supplementary video for all 24 edited video comparisons of our method against baselines.

Additional Ablation Results. We conduct ablation studies on more videos from HOSNeRF dataset [28] and NeuMan dataset [18]. To evaluate the effectiveness of each proposed component in DynVideo-E, we progressively ablate each component from local parts super-resolution, reconstruction loss, 2D personalized SDS, 3D SDS, and 2D personalization LoRA. We observe that the model even fails to converge on some videos when we disable several components of our model. We compute the average CLIP score of all successfully edited videos in Tab. 3, where the CLIP score progressively drops with the disabling of each component, with the full model achieving the best performances, which clearly demonstrates the effectiveness of our designs.

Visualization of Canonical Images from CoDeF [33] and Atlas from Text2LIVE [2] and StableVideo [6]. For challenging videos with large-scale motions and viewpoint changes, CoDeF [33], Text2LIVE [2], and StableVideo [6] largely overfit to input video frames and learn meaningless canonical images or neural atlas, and thus cannot generate meaningful editing results. We show several examples of their learned canonical images [33] and neural atlas [2, 6] in Fig. 9, where Text2LIVE [2] and StableVideo [6] utilizes the same foreground and background atlas during editing. As shown in Fig. 9, canonical images and atlas all fail to represent the challenging large-scale motion- and view-change videos, and thus they cannot gen-

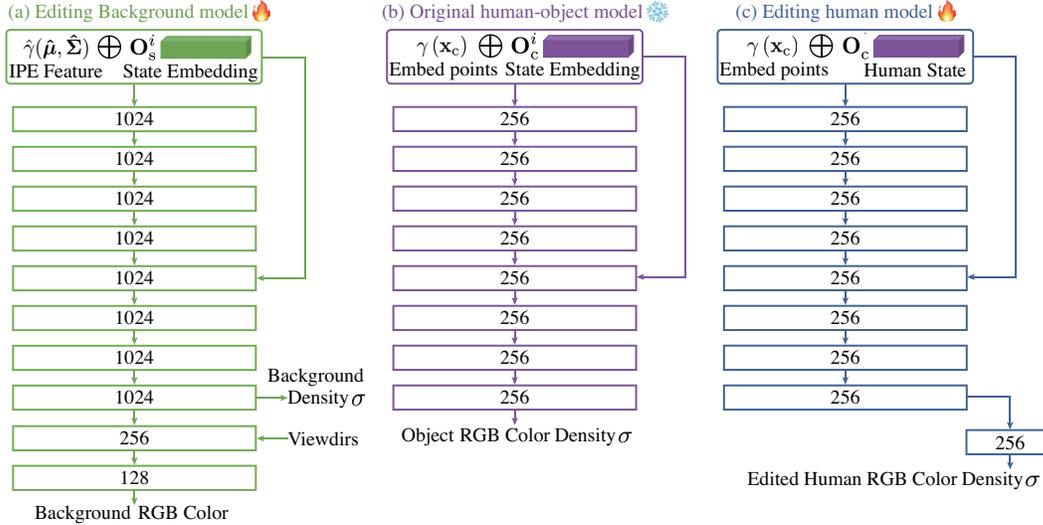


Figure 5. DynVideo-E network designs: (a) Editing Background model, (b) Original human-object model, (c) Editing human model.

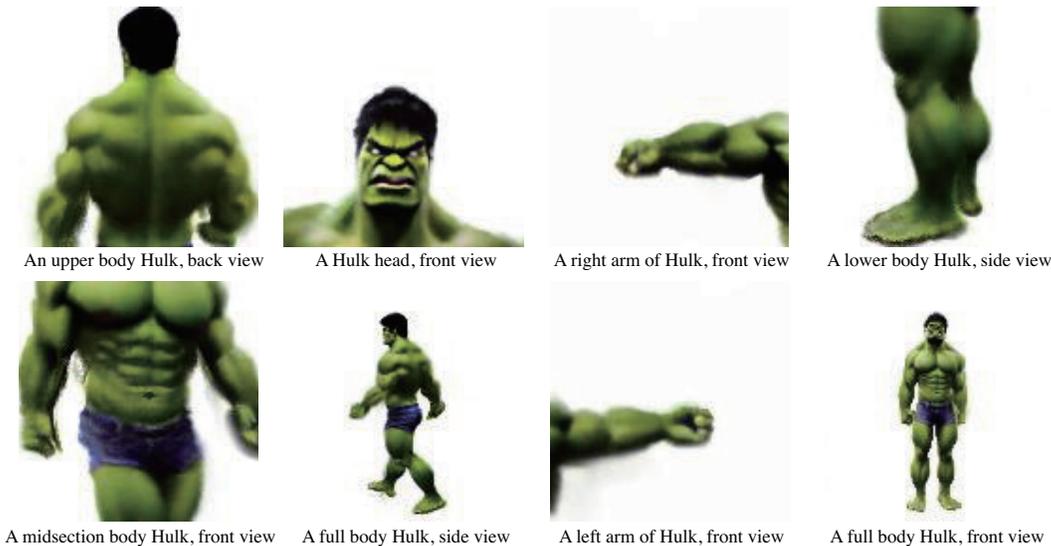


Figure 6. Visualization examples of text-guided local parts super-resolution sampled during training.

erate satisfactory editing results. In addition, the atlas performs better for short videos in NeuMan dataset [18] than long videos with a better background atlas, but the foreground atlas still cannot represent the humans with large motions. In contrast, our DynVideo-E represents videos with the dynamic NeRFs to effectively aggregate the large-scale motion- and view-change video information into a 3D dynamic human space and a 3D background space, and achieves high-quality video editing results by editing the 3D dynamic spaces.

Editing Operation Time Comparison. We compare the editing operation time of our DynVideo-E against other approaches on a long video of the HOSNeRF dataset ([300, 400] frames) using a single A100 GPU in Tab. 4. Although other approaches are faster than ours, 2D-video representation-based methods such as CoDeF [33], Stable-

Video [6], and Text2LIVE [2] cannot accurately reconstruct large-scale motion- and view-change videos and thus fail to generate meaningful editing results, as validated in Fig. 9. Text2Video-Zero [20] and Rerender-A-Video [57] fail to edit the challenging human-centric videos with large-scale motion and viewpoint changes and their editing results are highly inconsistent. Therefore, previous approaches cannot handle the challenging human-centric videos no matter how many computation resources are provided. In contrast, our method is the first work to achieve highly consistent long-term video editing that outperforms previous approaches by a large margin of 50% ~ 95% in terms of human preference, and we leave accelerating our model with voxel or hash grid representation as a faithful future direction.

Example of Human Preference Questionnaire. We uti-

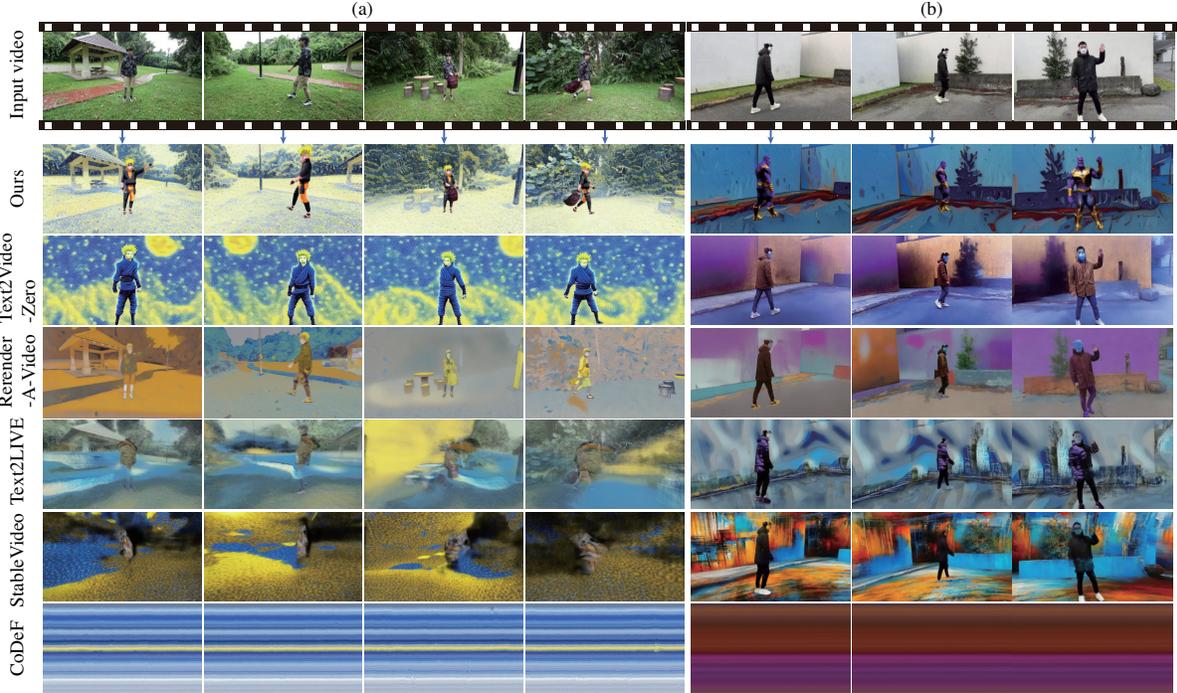


Figure 7. More qualitative comparisons of DynVideo-E against SOTA approaches on the Backpack scene (a) and Parkinglot scene (b).

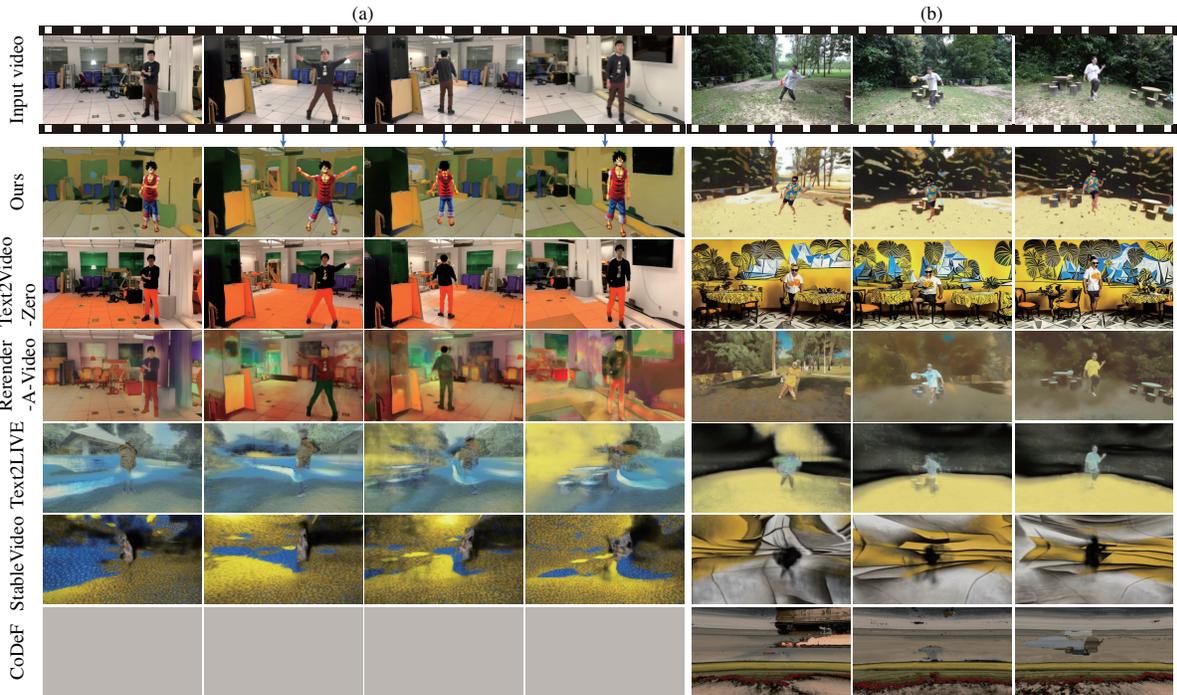


Figure 8. More qualitative comparisons of DynVideo-E against SOTA approaches on the Lab scene (a) and Dance scene (b).

lize Amazon MTurk [‡] to recruit raters to rate our pairwise comparing videos. For each comparison, we show our result and one baseline result (shuffled order in questionnaires), together with textual descriptions to raters and ask their preferences. In total, we collected 1140 comparisons

[‡] <https://requester.mturk.com/>

over all pairwise results from 32 different raters. Fig. 10 illustrate one comparison example in our questionnaires.

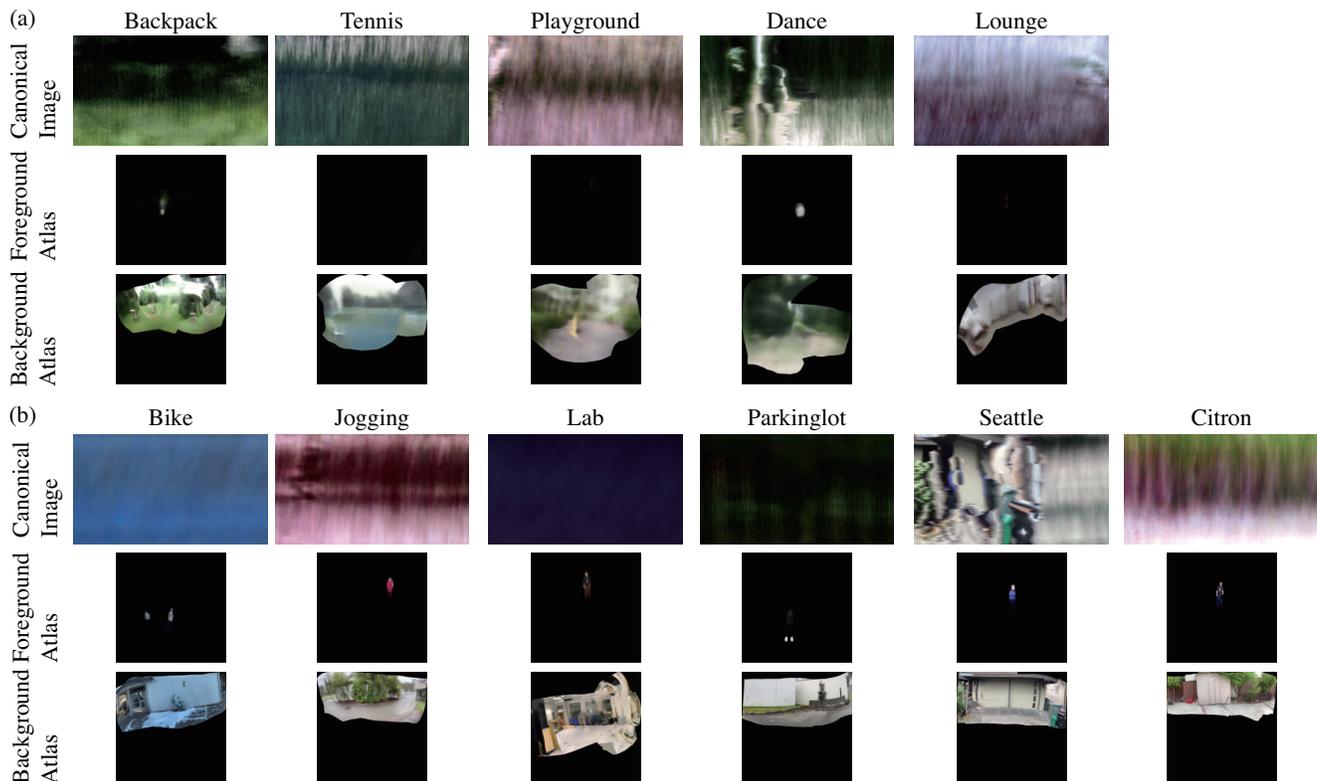


Figure 9. Visualization of canonical images from CoDeF [33], and foreground and background atlases from Text2LIVE [2] and StableVideo [6] on (a) HOSNeRF dataset [28] and (b) NeuMan dataset [18].

Method	CoDeF [33]	Text2Video-Zero [20]	Rerender-A-Video [57]	StableVideo [6]	Text2LIVE [2]	DynVideo-E (Ours)
Time	~ 1 mins	15 mins	1.2 hrs	~ 1 mins	~ 2 hrs	7.3 hrs

Table 4. Editing operation time comparison of our method against other approaches.

Instructions

Please watch two videos (best viewed in **full-screens**), and answer the following questions:

- **Text alignment:** Which video better matches the caption?
- **Temporal consistency:** Which video looks more natural in terms of human motion?
- **Overall quality:** Aesthetically, which video is better?

Option 1



Option 2



Question

1. Which video better matches the description **"Luffy"**?

Option 1 Option 2

2. Which video looks more natural in terms of human motion?

Option 1 Option 2

3. Aesthetically, which video is better?

Option 1 Option 2

Figure 10. One comparison example from our questionnaires.