

# SyncMV4D: Synchronized Multi-view Joint Diffusion of Appearance and Motion for Hand-Object Interaction Synthesis

Lingwei Dang<sup>1</sup> Zonghan Li<sup>1</sup> Juntong Li<sup>1</sup>  
 Hongwen Zhang<sup>2</sup> Liang An<sup>3</sup> Yebin Liu<sup>3</sup> Qingyao Wu<sup>1†</sup>

<sup>1</sup> South China University of Technology <sup>2</sup> Beijing Normal University <sup>3</sup> Tsinghua University

<https://droliven.github.io/SyncMV4D>

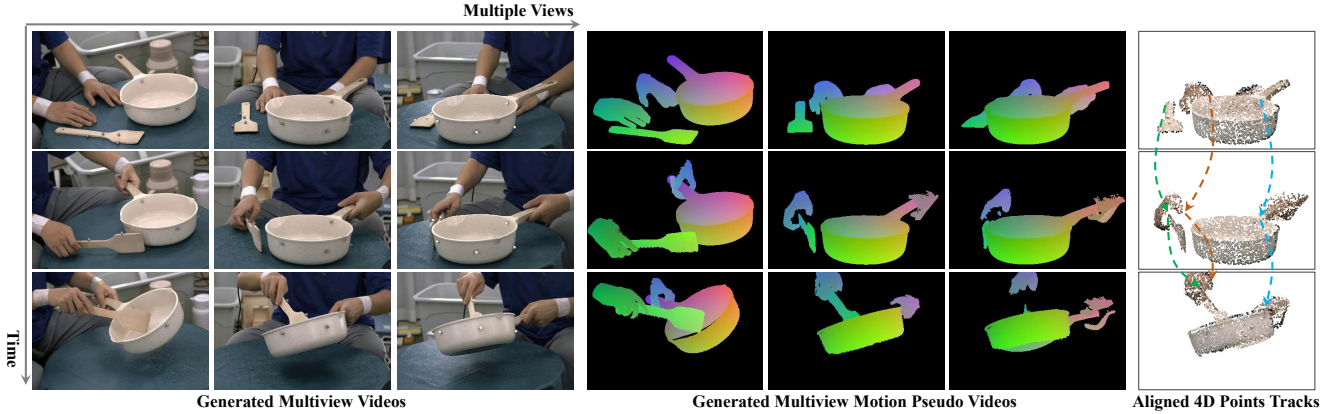


Figure 1. Our synchronized multi-view joint diffusion (SyncMV4D) simultaneously models multi-view geometry, visual appearance, and motion dynamics. It is capable of generating both multi-view hand-object interaction videos (left) and 4D motion sequences, comprising intermediate coarse pseudo videos (middle) and refined point tracks (right), with results achieving visual realism, dynamic plausibility, and geometric consistency.

## Abstract

*Hand-Object Interaction (HOI) generation plays a critical role in advancing applications across animation and robotics. Current video-based methods are predominantly single-view, which impedes comprehensive 3D geometry perception and often results in geometric distortions or unrealistic motion patterns. While 3D HOI approaches can generate dynamically plausible motions, their dependence on high-quality 3D data captured in controlled laboratory settings severely limits their generalization to real-world scenarios. To overcome these limitations, we introduce SyncMV4D, the first model that jointly generates synchronized multi-view HOI videos and 4D motions by unifying visual prior, motion dynamics, and multi-view geometry. Our framework features two core innovations: (1) a Multi-view Joint Diffusion (MJD) model that co-generates HOI videos and intermediate motions, and (2) a Diffusion Points Aligner (DPA) that refines the coarse intermediate motion into globally aligned 4D metric point tracks. To tightly*

*couple 2D appearance with 4D dynamics, we establish a closed-loop, mutually enhancing cycle. During the diffusion denoising process, the generated video conditions the refinement of the 4D motion, while the aligned 4D point tracks are reprojected to guide next-step joint generation. Experimentally, our method demonstrates superior performance to state-of-the-art alternatives in visual realism, motion plausibility, and multi-view consistency.*

## 1. Introduction

Realistic human-/hand-object interaction (HOI) generation plays a vital role in animation production [63] and dexterous robotic manipulation [4, 40, 47]. Current 3D HOI methods [6, 8, 15, 19, 27, 29, 30, 35, 41, 46, 54, 58, 64, 73] predominantly generate 6D object poses and parameterized human/hand poses (e.g., SMPL [39] or MANO [48]). Although these approaches yield kinematically plausible results, their reliance on high-quality motion capture data [9, 17, 18, 32, 36–38, 50, 60, 65, 71, 74] limits scalability and diversity, exposing a fundamental generalization bot-

<sup>†</sup> Corresponding Author. Email: qyw@scut.edu.cn.

tleneck.

Recent advances in generative video foundation models [7, 26, 42, 45, 52, 66] have shown great potential for high-fidelity generation, even serving as “world models” that implicitly capture spatial geometry, temporal dynamics, and physical rules [3, 51, 68]. These developments open new opportunities for HOI generation, yet effectively leveraging such visual priors remains challenging. Some methods enable controllable HOI video generation by injecting pose and appearance conditions [13, 21, 43, 63, 77]. However, they typically require pre-defined or pre-generated pose sequences, which limits their practical applicability. Others integrate video and motion generation within a joint diffusion process [10, 11, 14, 76] to improve physical plausibility, yet suffer from inefficient motion representations (*e.g.*, 2D optical flow, sparse keypoints, or depth maps). More critically, most operate under a single-view generation setting, hindering full 3D geometric perception.

Synchronized multi-view HOI generation provides a comprehensive understanding of object geometry, making it particularly valuable in heavily occluded scenarios such as hand-object interaction and embodied dexterous manipulation. Existing multi-view video generation methods, however, face certain limitations. Some methods generate novel views one at a time, conditioned on a source video [23, 34, 69, 70]. Despite the use of 3D intermediate representations, the separately generated results struggle to maintain multi-view consistency. Others can generate multiple viewpoint videos simultaneously [2, 53, 56, 67]. However, these methods are either confined to generating simple, background-free 4D assets or can only produce videos from a constrained set of viewpoints.

Our core insight is that synchronized multi-view HOI generation enables a comprehensive understanding of object geometry, thereby improving generation quality and applicability. Meanwhile, joint 2D video and 4D motion diffusion further enhances spatiotemporal coherence and physical plausibility. Based on this, we propose **SyncMV4D**, the first model that generates multi-view HOI videos and 4D motions using only text prompts and a reference image, as shown in Fig. 1. Using a pre-trained Diffusion Transformer (DiT) [44] backbone, SyncMV4D introduces inter-view geometry attention and motion modulation modules. This multi-view joint diffusion (**MJD**) framework unifies modeling of visual appearance, motion dynamics, and multi-view geometry via sequential inter-view attention, intra-view spatiotemporal attention, and multi-modal modulation. It produces view-consistent 2D videos and 4D motions without 3D object models or predefined poses.

Inspired by DaS [20], we represent motion with enhanced 4D point tracks to ensure temporal smoothness and 3D geometric awareness. Specifically, each point track is represented by three channels: the first two store the 2D

pixel coordinates of its anchor point in the first frame, while the last channel stores its metric depth at each frame. After normalization, they can be treated as “pseudo videos” and embedded into a shared VAE latent space, simplifying the adaptation of the DiT model. Furthermore, to align 4D motion across views, we design a Diffusion Points Aligner (**DPA**) module. It first converts the generated coarse motion pseudo-videos into world coordinates and uses them as conditions to synthesize globally aligned 4D point tracks. Owing to the similar diffusion pipelines of MJD and DPA, we establish a closed-loop feedback cycle: at each denoising step, the MJD output is refined by the DPA, reprojected into pseudo videos, and fed back to the MJD for iterative refinement, enabling mutual enhancement.

Experimental results show that our method achieves state-of-the-art performance in three aspects: multi-view video quality, motion plausibility, and cross-view consistency. Comprehensive ablation studies validate the effectiveness of each component in our framework.

Our contributions are threefold:

- The first synchronized multi-view HOI generation method that can synthesize results with high visual quality, motion plausibility, and view consistency, requires only reference images and text.
- A multi-view joint diffusion (MJD) framework of video and motion that unifies visual prior, motion dynamics, and multi-view geometry modeling.
- A Diffusion Points Aligner (DPA) module with closed-loop feedback refines the per-view misaligned coarse motions from MJD into globally aligned 4D point tracks and is co-optimized with the multi-view joint diffusion.

## 2. Related Work

**Multi-view video synthesis** follows two main paradigms: the first treats the task as a “translation” from a source video to a novel view, while the second generates all views jointly in one step. In the first paradigm, ReCamMaster [1] directly synthesizes the target view from the input video and a novel camera pose. Other methods [23, 34, 69, 70] rely on explicit 3D representations via monocular reconstruction, novel-view rendering, and inpainting or warping. DaS [20] and GS-DiT [5] further integrate tracking signals for improved quality. However, their separate, view-by-view generation often undermines multi-view consistency. The second paradigm enables simultaneous multi-view synthesis from text or a reference video [2, 53, 56, 67], but typically focuses on simple background-free 4D assets or supports only a few fixed viewpoints. In contrast, our approach generates synchronized multi-view HOIs in a single step and enhances physical plausibility through joint diffusion of 2D video and 4D motion signals.

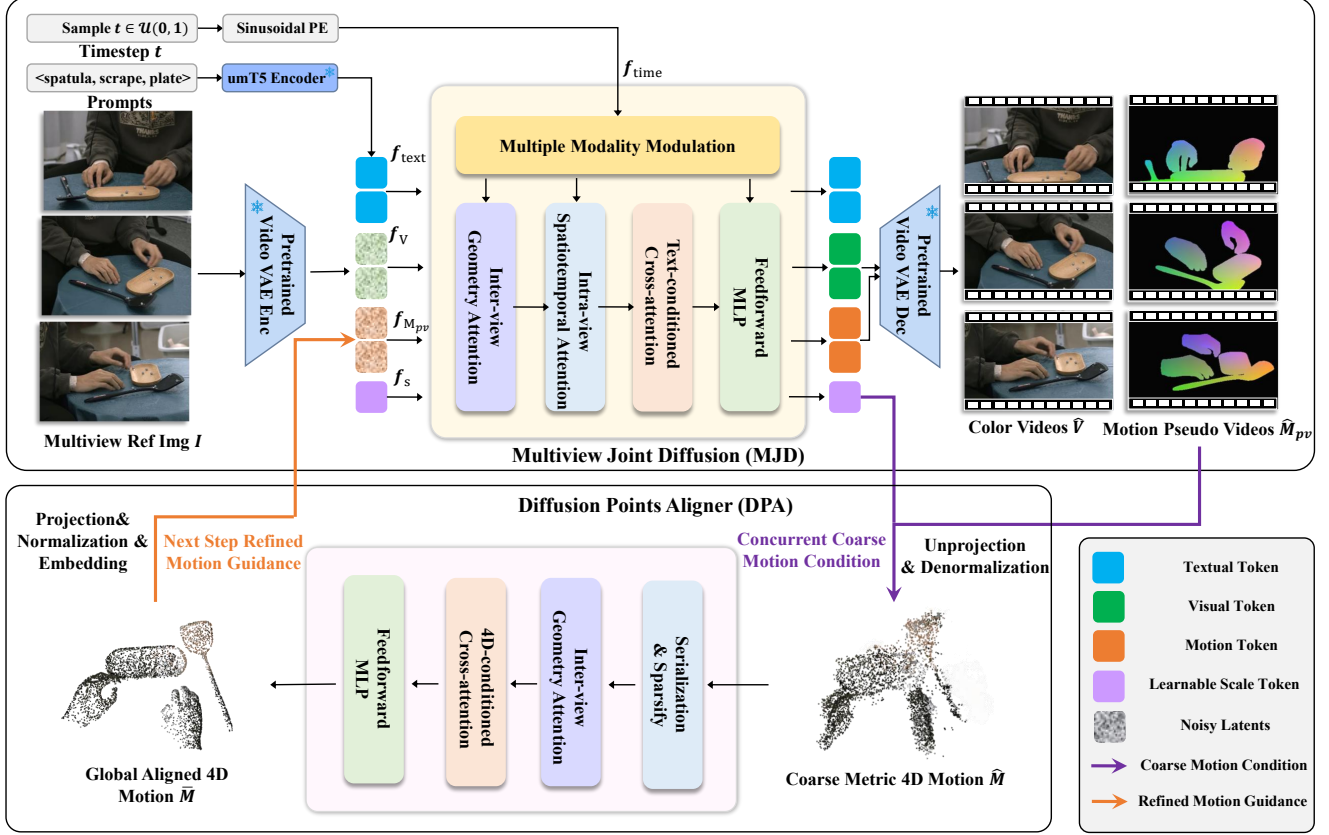


Figure 2. Our SyncMV4D consists of two key components: First, the Multi-view Joint Diffusion (MJD) module generates synchronized multi-view color videos, intermediate motion pseudo videos, and metric depth scales (Sec. 3.3). Second, the Diffusion Points Aligner (DPA) module takes the resulting coarse 4D motions as a conditioning signal to reconstruct globally aligned 4D point tracks (Sec. 3.4). Furthermore, since both MJD and DPA are iterative denoisers, the refined 4D point tracks from DPA are fed back to guide MJD in subsequent denoising steps, forming a closed-loop mutual enhancement cycle (Sec. 3.5).

**HOI Video models.** Recent video foundation models [26, 52, 66] have significantly advanced HOI video generation. Some approaches [21, 63, 77] extend UNet-based architectures with pose guides and appearance reference networks to enable pose-controlled synthesis. However, the incremental temporal modeling of UNets often causes flickering artifacts, and these methods require predefined or pre-generated pose sequences as input. More recent works [10, 14, 76] leverage the native spatio-temporal attention of Diffusion Transformers (DiT) [16] and adopt a video-motion co-generation paradigm to improve physical plausibility. Yet motion representation remains an open challenge. VideoJam [10] uses 2D optical flow, which lacks explicit 3D awareness. SViMo [14] relies on sparse key-points and suffers from limited precision. Tesseract [76] employs pixel-aligned depth but lacks inter-frame smoothness. In contrast, our multi-view joint diffusion simultaneously generates 2D videos and metric 4D point tracks, achieving both 3D awareness and temporal stability.

**3D HOI generation** primarily relies on high-precision

3D motion capture data [9, 18, 37, 71]. Some works [8, 15, 27, 28, 30, 31, 33, 35, 73] enhance kinematic plausibility by predicting intermediate contact maps or affordances. Others [6, 41, 55, 59, 61] integrate complex physics simulators to improve dynamic realism. However, the limited dataset scale and diversity constrain their generalization. A few approaches [72, 75] leverage semantic knowledge from multimodal vision-language models (VLMs) to boost HOI generalization, but their multi-stage pipelines are prone to error accumulation.

### 3. Method

Given multi-view reference images  $I \in \mathbb{R}^{V \times H \times W \times 3}$  and a textual prompt  $P$ , we aim to synthesize synchronized multi-view hand-object interaction (HOI) videos  $V \in \mathbb{R}^{V \times N \times H \times W \times 3}$  along with a metric 4D motion sequence represented as point tracks  $M \in \mathbb{R}^{V \times N \times K \times 3}$ , where  $V$ ,  $N$ ,  $H$ ,  $W$ , and  $K$  denote the number of viewpoints, temporal frames, height, width, and tracked points, respectively.

### 3.1. Preliminary: Basic Video Foundation Model

Our framework is built upon a pre-trained foundation model for text-and-image-to-video generation. It comprises two key components: a spatio-temporal variational autoencoder (VAE) [25] that compresses the original video  $V$  into a more compact latent space  $z$ , and a Diffusion Transformer (DiT) [44] based video generator to synthesize video latents  $\hat{z}$ . The model employs the Rectified Flow framework [16] for noise scheduling and denoising operations. During training, given clean video latents  $z_0$ , Gaussian noise  $z_1 \in \mathcal{N}(\mathbf{0}, I)$ , and a timestep  $t \in [0, 1]$ , intermediate latents  $z_t$  are constructed through linear interpolation  $z_t = (1-t) \cdot z_0 + t \cdot z_1$ . The corresponding ground-truth velocity  $v_t$  is defined as:  $v_t = dz_t/dt = z_1 - z_0$ . The model parameterized with  $\Theta$  is trained to predict the velocity field using mean squared error loss:

$$\mathcal{L} = \mathbb{E}_{z_0, z_1, c, t} \|v_t - \hat{v}_\Theta(z_t, c, t)\|_2^2. \quad (1)$$

In the inference phase, the framework is processed with iterative sampling:

$$z_{t-1} = z_t + \Delta t \cdot \hat{v}_\Theta(z_t, c, t). \quad (2)$$

### 3.2. Framework Overview

Synchronized multi-view HOI generation enables comprehensive perception of full object geometry, making it especially valuable in heavily occluded hand-object interaction scenarios. However, existing multi-view video generation methods face two key limitations. Approaches based on a “video translation” paradigm generate only one novel view at a time, leading to inconsistency across separately synthesized views. Others can produce multi-view videos simultaneously but typically focus on simple, background-free 4D assets or support only a few fixed viewpoints. To address these issues, we propose a novel end-to-end framework for synchronized multi-view video and motion generation. As shown in Fig. 2, our framework comprises two core components. The first is a Multi-view Joint Diffusion (MJD) model, built upon a pre-trained single-view video foundation model, which jointly generates color videos, pseudo videos encoding 4D motion, and their metric scale (Sec. 3.3). The second component is the Diffusion Points Aligner (DPA), which refines MJD’s coarse multi-view motions into globally consistent 4D point tracks (Sec. 3.4). Moreover, leveraging the shared diffusion structure of MJD and DPA, we establish a closed-loop cycle that enables their mutual refinement (Sec. 3.5).

### 3.3. Synchronized Multi-view Joint Diffusion

MJD learns to generate multi-view color videos  $\hat{V}$ , pseudo videos of 4D motions  $\hat{M}_{pv}$ , alongside with the metric scale  $s$ . A detailed description is provided below.

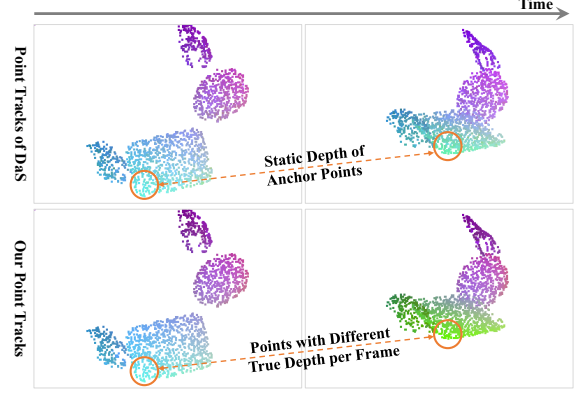


Figure 3. Comparison of motion representations between that of DaS [20] and our 4D point tracks. For each point, the first two dimensions represent the pixel coordinates of the tracked point in the first frame. The difference lies in the third dimension: DaS uses the static depth from the first frame, whereas we use the actual per-frame depth to enhance 3D perceptual capability.

**Data Representation and Embedding.** In DaS’s “tracking video” representation [20], the depth of points in each frame is fixed to that of the first-frame anchor points, which limits the model’s perception of actual depth, as shown in Fig. 3. To ensure temporal stability and 3D awareness, we enhance it by using the actual depth of points in each frame. Each frame contains  $K$  points, each defined by three values: the first two are the pixel coordinates of the tracked anchor point in the first frame, and the third is the metric depth in the current frame. The first two dimensions ensure temporal smoothness, while the depth introduces 3D geometric awareness, contributing to physically plausible results. To convert  $M$  into a pseudo video  $M_{pv}$  compatible with DiT, we normalize the first two dimensions by the video resolution and apply min-max normalization to the third dimension, mapping it to  $[0, 1]$ . The normalized values are then scaled by 255 and rendered to obtain the motion pseudo videos  $M_{pv}$ . A visualization of this representation is provided in Fig. 3. Note that the normalized pseudo video loses metric scale, so we introduce an additional mechanism to regress the min-max scale  $s \in \mathbb{R}^{V \times 2}$  of the multi-view depth.

The text instruction  $P$  is encoded by the frozen Google umT5 model [12] and projected to  $f_{\text{text}} \in \mathbb{R}^{L \times d}$ . Multi-view reference frames  $I$  are encoded by a video VAE into  $z^I$ , while the color video  $V$  and motion pseudo video  $M_{pv}$  are encoded into  $z^V$  and  $z^{M_{pv}} \in \mathbb{R}^{V \times \frac{N-1}{rn} \times \frac{H}{rh} \times \frac{W}{rw} \times d_{\text{VAE}}}$ , where  $(rn, rh, rw)$  denote the video VAE’s compression ratios for temporal, height, and width dimensions, respectively. Using the forward diffusion process (Sec. 3.1), we obtain noisy latents  $z_t^V$  and  $z_t^{M_{pv}}$ . Following the foundation model’s conditioning mechanism, the first frame



of the noisy  $z_t^V$  is replaced with  $z^I$ , while  $z_t^{M_{pv}}$  remains unchanged. Both are then tokenized and flattened to produce visual tokens  $\mathbf{f}_V$  and motion tokens  $\mathbf{f}_{M_{pv}} \in \mathbb{R}^{V \times \frac{N-1}{rn} \times \frac{H}{rh*2} \times \frac{W}{rw*2} \times d}$ .

**Multi-view Multimodality Diffusion Blocks.** In addition to processing text tokens  $\mathbf{f}_{\text{text}}$ , multi-view visual tokens  $\mathbf{f}_V$ , and motion tokens  $\mathbf{f}_{M_{pv}}$ , the Diffusion Blocks also process  $t_s$  learnable scale tokens  $\mathbf{f}_s \in \mathbb{R}^{t_s \times d}$  and  $t_r$  registration tokens  $\mathbf{f}_r \in \mathbb{R}^{t_r \times d}$  to regress the min-max scale of the normalized depth in the motion pseudo video.

We enhance the original single-view video DiT Blocks by incorporating additional inter-view attention modules to learn multi-view spatial geometry, and by introducing extra modulation modules to handle motion modality features. Each DiT block consists of alternating multimodality adaptive modulation, inter-view geometry attention, intra-view spatiotemporal attention, text-conditioned cross-attention, and feedforward MLP layers. First, due to the semantic and distributional divergence between video and motion features, directly sharing the visual token modulator from the DiT backbone for their co-generation is suboptimal. Therefore, we introduce a dedicated motion modulation module, analogous to the video modality and conditioned on the denoising timestep. For inter-view geometry attention, the features are permuted and reshaped to  $[(B \cdot \frac{N-1}{rn}), (V \cdot \frac{H}{rh*2} \cdot \frac{W}{rw*2}), d]$  to model relationships between tokens from different views at the same timestep. For intra-view spatiotemporal attention, features are reshaped to  $[(B \cdot V), (\frac{N-1}{rn} \cdot \frac{H}{rh*2} \cdot \frac{W}{rw*2}), d]$  to capture dependencies among tokens across different frames within the same view.

**Training Objectives.** The output video, motion, and scale features from the final DiT block are projected through linear layers and depatchified to yield the multi-view video latent velocity  $\hat{\mathbf{v}}^V$ , the motion pseudo video latent velocity  $\hat{\mathbf{v}}^{M_{pv}}$ , and the predicted metric scale  $\hat{s}$ . Thus, the MJD is optimized with the following objective:

$$\begin{cases} (\hat{\mathbf{v}}^V, \hat{\mathbf{v}}^{M_{pv}}, \hat{s}) = \mathcal{G}_{\text{MJD}} \left[ \mathbf{z}_t^V, \mathbf{z}_t^{M_{pv}}, \mathbf{z}^I, \mathbf{P}, t \right], \\ \mathcal{L}_{\text{MJD}} = \mathbb{E} \left[ \|\hat{\mathbf{v}}^V - \mathbf{v}^V\|_2^2 + \|\hat{\mathbf{v}}^{M_{pv}} - \mathbf{v}^{M_{pv}}\|_2^2 + \|\hat{s} - s\|_2^2 \right], \end{cases} \quad (3)$$

where  $\mathbf{v}^V$  and  $\mathbf{v}^{M_{pv}}$  are the ground-truth velocities, as derived in Sec. 3.1.

### 3.4. Diffusion Points Aligner

Given the intermediate motion pseudo-video  $\hat{\mathbf{M}}_{pv}$  and the predicted depth scale  $\hat{s}$ , we obtain the 4D motion tracks  $\bar{\mathbf{M}}$  through denormalization and unprojection. However, in practice, since the video VAE and DiT are not designed or trained for 4D motion data, the resulting 4D motion tracks may contain inaccuracies. To address this issue, DPA formulates the task as a conditional generation problem: it

**Algorithm 1** Joint training and mutual enhancement of MJD and DPA.

---

**Input:** Multi-view reference image  $\mathbf{I}$ , text prompt  $\mathbf{P}$ , target multi-view video  $\mathbf{V}$ , target 4D motion  $\mathbf{M}$ , video VAE encoder  $\mathcal{E}$  and decoder  $\mathcal{D}$ , MJD network  $\mathcal{G}_{\text{MJD}}^\theta$ , DPA network  $\mathcal{G}_{\text{DPA}}^\phi$ .

**Output:** Optimized parameters  $\theta^*$  and  $\phi^*$ .

- 1:  $\mathbf{s} = \text{Min-Max}(\mathbf{M}[\dots, -1])$   $\triangleright$  GT metric depth scale
- 2:  $\mathbf{M}_{pv} = \text{Proj \& Norm}(\mathbf{M})$   $\triangleright$  motion pseudo video
- 3:  $\mathbf{z}_0^V = \mathcal{E}(\mathbf{V})$ ,  $\mathbf{z}_0^{M_{pv}} = \mathcal{E}(\mathbf{M}_{pv})$ ,  $\mathbf{z}_I = \mathcal{E}(\mathbf{I})$   $\triangleright$  latents
- 4: **while** not converged **do**
- 5:    $t \sim \mathcal{U}\{0, 1\}$   $\triangleright$  sample time step  $t$
- 6:   Diffuse to get  $(\mathbf{z}_t^V, \mathbf{z}_t^{M_{pv}}, \mathbf{M}_t)$   $\triangleright$  Sec. 3.1
- 7:    $\tilde{\mathbf{M}} = \mathcal{G}_{\text{DPA}}^{\text{no-grad}} \left( \mathbf{M}_t, \text{Denorm \& Dec}(\mathbf{z}_t^{M_{pv}}, \mathbf{s}), t \right)$   $\triangleright$  Eq. 4
- 8:    $\tilde{\mathbf{M}}_{pv} = \text{Proj \& Norm}(\tilde{\mathbf{M}})$ ,  $\tilde{\mathbf{z}}^{M_{pv}} = \mathcal{E}(\tilde{\mathbf{M}})$
- 9:    $(\hat{\mathbf{v}}^V, \hat{\mathbf{v}}^{M_{pv}}, \hat{s}) = \mathcal{G}_{\text{MJD}}(\mathbf{z}_t^V, \mathbf{z}_t^{M_{pv}} + \tilde{\mathbf{z}}^{M_{pv}}, \mathbf{z}_I, \mathbf{P}, t)$   $\triangleright$  Eq. 3: co-generation with refined motion guidance
- 10:    $\tilde{\mathbf{z}}^{M_{pv}} = \mathbf{z}_t^{M_{pv}} + \Delta t \cdot \hat{\mathbf{v}}^{M_{pv}}$   $\triangleright$  Sec. 3.1
- 11:    $\bar{\mathbf{M}} = \mathcal{G}_{\text{DPA}} \left[ \mathbf{M}_t, \text{Denorm \& Dec}(\mathbf{z}_t^{M_{pv}}, \hat{s}), t \right]$   $\triangleright$  conditional points refine Eq. 4
- 12:    $\mathcal{L} = \mathcal{L}_{\text{MJD}} + \mathcal{L}_{\text{DPA}}$   $\triangleright$  Eq. 3, 4
- 13:   update parameters  $\theta$  and  $\phi$  by gradient descent
- 14: **end while**
- 15: **return**  $\theta^* = \theta$ ,  $\phi^* = \phi$

---

takes coarse point tracks as input conditions and generates globally aligned 4D point sequences  $\bar{\mathbf{M}}$ . Specifically, DPA consists of multiple sequentially connected inter-view geometry attention and coarse motion conditioned cross-attention modules, each built upon the sparse convolutions of Point Transformer V3 [57]. The DPA is also formalized with the flow matching operation in Sec. 3.1 and optimized with the following loss function:

$$\begin{cases} \bar{\mathbf{v}}^M = \mathcal{G}_{\text{DPA}} \left[ \mathbf{M}_t, \hat{\mathbf{M}}, t \right], & \bar{\mathbf{M}} = \mathbf{M}_t - t \cdot \bar{\mathbf{v}}^M, \\ \mathcal{L}_{\text{DPA}} = \mathbb{E} \left[ \text{MSE}(\bar{\mathbf{v}}^M, \mathbf{v}^M) + \text{D}_{\text{chamfer}}(\bar{\mathbf{M}}, \mathbf{M}) \right]. \end{cases} \quad (4)$$

### 3.5. Close-loop Mutual Enhancement Cycle

Owing to the similar diffusion pipeline between MJD and DPA, we establish a closed-loop feedback: MJD’s outputs serves as concurrent coarse motion condition for DPA, while DPA’s globally aligned points are projected and normalized into a motion pseudo video to guide the next-step denoising of MJD. This mutually enhancing process during training is summarized in Alg. 1. The inference pseudocode is included in the Supp.

## 4. Experiments

### 4.1. Settings

**TACO dataset** [37] is a large-scale hand-object interaction benchmark that provides high-precision 3D object models, pose sequences, and synchronized high-resolution videos from 12 viewpoints. It captures diverse interactions, each structured as a “tool–action–target” triplet, denoting the use of a tool to perform an action on a target object. Tools and targets span 20 physical categories (196 distinct instances), with 15 action types performed by 14 participants. For video preprocessing, we crop hand-object regions from both allocentric-view videos (original resolution  $4096 \times 3000$ ) and egocentric-view videos (1080p), then resize them to 480p to meet our video foundation model’s input requirements. The frame rate is uniformly downsampled from 30 fps to 8 fps. Our 4D point tracks are computed from ground-truth vertex coordinates of the MANO hand model and object meshes. We group the 12 camera views into three position-based clusters: left, right, and center. In each training iteration, we randomly sample one view from each group to enable synchronous three-view generation. Additionally, we adopt a two-stage split: first, we hold out samples involving specific objects (*e.g.*, hammer) and actions (*e.g.*, measure) as a separate test set for generalization evaluation. The remaining data is then divided into training and validation sets in a 9:1 ratio.

**Metrics.** Our method is capable of simultaneously generating synchronized multi-view videos and 4D motions. For video quality evaluation, we assess both single-view quality and multi-view consistency. For the former, we adopt two key metrics from VBench [22]: **Subject Consistency** and **Dynamic Degree**, which measure the temporal consistency of the subject across frames and the overall dynamic quality of the video, respectively. For multi-view consistency, we follow SynCamMaster [2] and employ two metrics: **Matching Pixels** and **CLIP-Views**. The former uses the GIM [49] image matching method to quantify the number of aligned pixels between different views, while the latter evaluates semantic consistency by measuring CLIP feature similarity across viewpoints of the same frame. For motion quality evaluation, we similarly examine both single-view and multi-view aspects. Single-view metrics include **Chamfer Distance** and **Motion Smoothness**, which assess the accuracy and temporal coherence of the motion. In addition, following GeometryCrafter [62], we introduce **Relative Point Error (RPE)** and **Percentage of Inliers (PT, threshold 0.25)** to evaluate the accuracy of the point tracks. Multi-view motion consistency is evaluated using the same two metrics after multi-view reprojection. Due to page limitations, detailed formulations of the evaluation metrics will be provided in the supplementary material.

### 4.2. Implementation Details.

**Model Architecture.** Our multi-view joint diffusion model extends the pre-trained text-and-image-to-video foundation model WAN2.2-5B-480P [52] to simultaneously generate hand-object interaction (HOI) videos from  $V = 3$  viewpoints. The VAE in WAN2.2 employs compression ratios  $(r_n, r_h, r_w) = (4, 16, 16)$ , and its visual tokenizer applies an additional  $2\times$  downsampling, requiring input resolutions divisible by 32. We thus adopt a video resolution of  $(N \times H \times W) = (49 \times 480 \times 704)$ . To recover metric scale, we introduce  $t_s = 2$  learnable scale tokens and  $t_r = 6$  register tokens. With a maximum text token length of  $L = 512$ , the total token count is  $512 + 2 + 6 + 2 \times [13 \times (480/32) \times (704/32)] = 9100$ , with hidden dimension  $d = 3072$ . The model comprises 30 DiT blocks. Our Diffusion Points Aligner is built upon Point Transformer V3 [57], where tracked points in each frame are voxelized and downsampled to  $K = 1600$  points.

**Training Details.** Our model was trained on 8 NVIDIA A800 GPUs. To boost training efficiency, we first ran a 5K-step warm-up phase using 10% of the data, separately training the Diffusion Points Aligner and the base single-view WAN-I2V model. This was followed by full-parameter training of the entire SyncMV4D model for 30K steps on the full dataset. We employed several memory optimization techniques, including DeepSpeed ZeRO-3, gradient checkpointing, and mixed-precision training with BF16. Additionally, we adopted a progressive training strategy: starting at a lower resolution ( $256 \times 384$ ) and later fine-tuning at the target resolution ( $480 \times 704$ ) to further improve efficiency.

### 4.3. Comparison with Baselines

**Baselines.** For video generation quality evaluation, we compare against single-view image-to-video (I2V) models, WAN2.2 [52] and SViMo [14], by generating multi-view videos one-by-one using reference frames from different viewpoints. We also compare with the single-view video-to-video (V2V) recamera method DaS [20]: a randomly selected video from the nine non-ground-truth views is used as input to render three novel-view videos along distinct camera trajectories. Additionally, we evaluate against the synchronized multi-view V2V method SV4D 2.0 [67], which conditions on one randomly chosen non-ground-truth view and simultaneously generates four-view videos. We select three of these for comparison. Notably, SV4D 2.0 removes the background in the reference video and focuses on foreground content, so our evaluation is restricted to foreground regions.

For 4D motion evaluation, as no existing method directly generates 4D sequences from reference images, we compare our Image-to-4D generation (Ito4G) with single-view Video-to-4D reconstruction (Vto4R) approaches, Geo4D [24] and GeometryCrafter [62]. These methods pro-

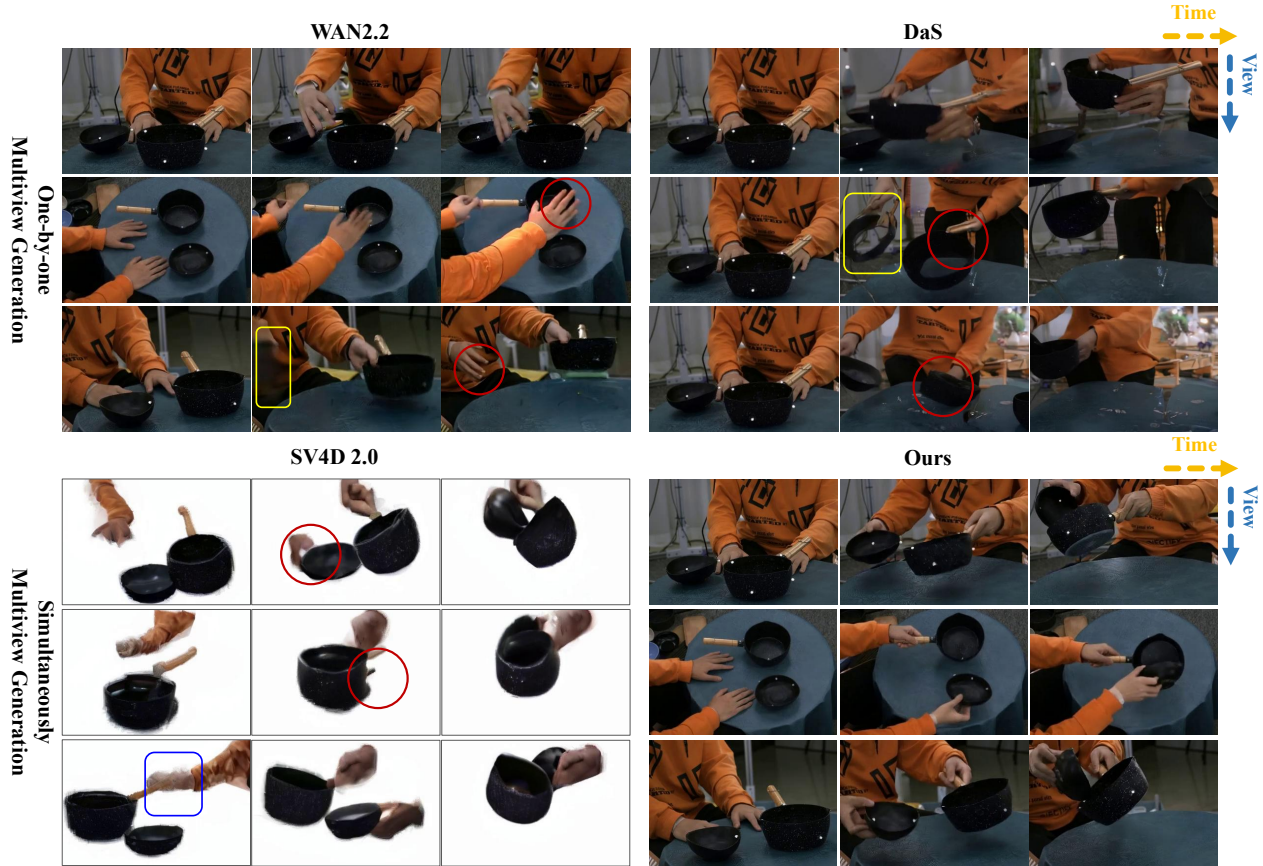


Figure 4. Visualization of the generated multi-view videos from different methods. Red circles indicate multi-view inconsistencies, yellow boxes highlight video distortions, and blue boxes denote blurring artifacts.

Table 1. Comparison of video quality. The best and second best results are highlighted with **bold** and underlined fonts. Note that “SV” means single view, and “MV” means multi-view.

Method	View	Type	Single-view		Multi-view	
			Subj. Cons.	Dyn. Deg.	Mat. Pix.	CLIP-V
WAN2.2 [52]	SV	I2V	<b>0.9562</b>	0.4521	33.8	76.78
SViMo [14]		I2VM	0.9285	0.9601	137.5	82.88
DaS [20]		V2V	0.8828	<b>0.9962</b>	<u>182.3</u>	<u>83.34</u>
SV4D 2.0 [67]	MV	V2V	0.7988	0.8628	108.7	79.31
Ours		I2VM	<u>0.9351</u>	<u>0.9877</u>	<b>529.4</b>	<b>83.67</b>

duce up-to-scale pointmaps without metric scale. Following GeometryCrafter, we apply optimization-based alignment to their outputs before metric computation. In contrast, our method directly yields metric point tracks and requires no such transformation.

**Qualitative and Quantitative Comparison.** The comparison results on video generation quality are shown in Tab. 1 and Fig. 4. Our method achieves the best multi-view consistency and the second-best single-view quality. Specifically, although WAN2.2 [52] obtains the high-

Table 2. Quantitative Comparison of Motions. Best in **Bold**.

Method	Type	Single-view				Multi-view	
		RPE	PI	Cham. Dis.	Smoo.	RPE	PI
Geo4D [24]	Vto4R	21.7	71.5	0.0134	0.0228	81.1	1.56
GeoCrafter [62]	Vto4R	17.0	87.3	0.0115	0.0222	67.3	4.16
Ours	Ito4G	<b>15.2</b>	<b>98.2</b>	<b>0.0103</b>	<b>0.0152</b>	<b>32.7</b>	<b>39.1</b>

est subject consistency score, this is largely due to its extremely low dynamic degree (Fig. 4, left top). This indicates that both metrics are partial and should not be interpreted in isolation. Among the three single-view methods, both WAN2.2 [52] and SViMo [14] suffer from poor multi-view consistency. DaS [20] achieves relatively higher multi-view consistency but still lags behind our method. Moreover, because DaS employs continuously varying camera trajectories, it attains the highest dynamic degree score, unlike our static-camera setting. However, as shown in the top-right of Fig. 4, its outputs exhibit noticeable distortions. The multi-view method SV4D 2.0 [67] focuses on generating simple, background-free 4D asset videos. It underperforms on both



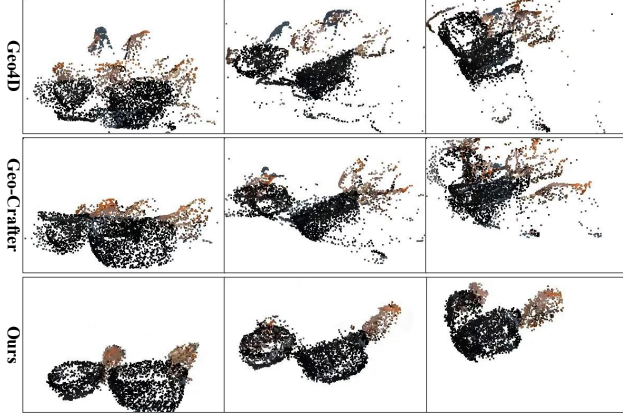


Figure 5. Visualization of multi-view points reprojected onto the same coordinate system.

Table 3. Ablation studies on the key components of SyncMV4D, including simultaneous multi-view generation (Multi-view), joint video-motion diffusion (MJD), diffusion points aligner (DPA), and the mutual enhancement cycle (Cycle).

Multi-view	MJD	DPA	Cycle	Mat. Pix. (MV)	RPE (MV)
	✓	✓	✓	122.7	75.3
✓				462.6	-
✓	✓			501.3	46.3
✓	✓	✓		498.2	33.5
Ours				<b>529.4</b>	<b>32.7</b>

single-view and multi-view metrics and suffers from severe temporal flickering and blurriness.

For motion generation, as shown in Table 2 and Fig. 5, our method demonstrates superior performance on both single-view and multi-view motion metrics. It outpaces the second-best method on RPE by 11% and 51% in the respective categories. These gains stem from our method’s joint generation of multi-view consistent points, in contrast to the baseline methods that suffer from inherent inconsistencies due to their per-view video reconstruction.

#### 4.4. Ablation Study

**Effectiveness of the Multi-view Generation.** We argue that simultaneously generating HOI from multiple viewpoints is crucial for improving viewpoint consistency. To verify this, we adapt SyncMV4D into a single-view image-to-video and motion generation framework, as shown in the first row of Tab. 3. Both the multi-view consistency of the generated videos and the 4D motion significantly degrade. This is because, influenced by the diverse generation capability of video foundation models, one-by-one sequential viewpoint generation cannot guarantee consistent motion patterns of HOI across different views.

**Ablation on Video-Motion Joint Diffusion.** We retain

only the multi-view synchronized generation architecture of our method, removing all other components to form a pure image-to-multi-view-video framework. As shown in the second row of Tab. 3, while its multi-view video consistency significantly surpasses the aforementioned single-view setting, it still lags behind our full SyncMV4D. This is because pure video generation lacks physical constraints, resulting in artifacts such as blurriness and deformation.

**Impact of the Diffusion Points Aligner.** We remove the DPA module and retain only the multi-view video motion joint generation module (MJD), directly combining the motion pseudo-video generated by MJD with the estimated metric scale for denormalization and unprojection to obtain 4D point tracks. As shown in the third row of Tab. 3, the motion quality drops significantly. This is because the video foundation model is not trained specifically for 4D motion, and the video VAE incurs information loss when encoding and decoding such pseudo-videos, thereby limiting the accuracy of the generated 4D motion.

**Influence of the Mutual Enhancement Cycle.** We simply input the detached coarse motion into the DPA, removing the “next-step refined motion as guidance” pathway. As shown in the second-to-last row of Tab. 3, this leads to a slight degradation in both video and motion consistency. In contrast, our final method achieves superior results through the closed-loop mutual enhancement cycle between MJD and DPA.

## 5. Conclusion

In this work, we propose a method for synchronized multi-view co-generation of videos and motions in hand-object interactions. Our approach uses a joint appearance-motion diffusion model to ensure visual realism, plausible motion, and geometric consistency across views. Motion is represented by stable, 3D-aware point tracks, which are aligned across viewpoints using a novel Diffusion Point Aligner. A dual-branch diffusion architecture with closed-loop feedback further enables mutual enhancement between the two components. The method requires only a reference image and a text instruction, making it highly accessible, especially effective in occluded scenarios such as hand-object interaction. We believe our framework also offers valuable insights for building physics-aware video world models.

**Limitation.** Our method currently requires multi-view reference images to generate videos and motions. A more appealing and user-friendly alternative would be to produce synchronized multi-view outputs from a single reference image, which can be achieved by integrating advanced novel-view image synthesis techniques. Additionally, the system could support user-specified camera viewpoints for controllable multi-view generation. This can be realized by incorporating camera conditional guidance modules, trained on densely sampled multi-view data.



## References

- [1] Jianhong Bai, Menghan Xia, Xiao Fu, Xintao Wang, Lianrui Mu, Jinwen Cao, Zuozhu Liu, Haoji Hu, Xiang Bai, Pengfei Wan, and Di Zhang. Recammaster: Camera-controlled generative rendering from a single video. In *ICCV*, pages 14834–14844, 2025. 2
- [2] Jianhong Bai, Menghan Xia, Xintao WANG, Ziyang Yuan, Zuozhu Liu, Haoji Hu, Pengfei Wan, and Di ZHANG. Syn-cammaster: Synchronizing multi-camera video generation from diverse viewpoints. In *ICLR*, pages 58038–58060, 2025. 2, 6
- [3] Philip J. Ball, Jakob Bauer, Frank Belletti, Bethanie Brownfield, Ariel Ephrat, Shlomi Fruchter, Agrim Gupta, Kristian Holsheimer, Aleksander Holynski, Jiri Hron, et al. Genie 3: A new frontier for world models. URL <https://deepmind.google/blog/genie-3-a-new-frontier-for-world-models/>, 2025. 2
- [4] Homanga Bharadhwaj, Debidatta Dwibedi, Abhinav Gupta, Shubham Tulsiani, Carl Doersch, Ted Xiao, Dhruv Shah, Fei Xia, Dorsa Sadigh, and Sean Kirmani. Gen2act: Human video generation in novel scenarios enables generalizable robot manipulation. In *CoRL*, 2025. 1
- [5] Weikang Bian, Zhaoyang Huang, Xiaoyu Shi, Yijin Li, Fu-Yun Wang, and Hongsheng Li. Gs-dit: Advancing video generation with dynamic 3d gaussian fields through efficient dense 3d point tracking. In *CVPR*, pages 21717–21727, 2025. 2
- [6] Jona Braun, Sammy Christen, Muhammed Kocabas, Emre Aksan, and Otmar Hilliges. Physically plausible full-body hand-object interaction synthesis. In *2024 International Conference on 3D Vision (3DV)*, pages 464–473. IEEE, 2024. 1, 3
- [7] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. 2024. URL <https://openai.com/research/video-generation-models-as-world-simulators>, 3:1, 2024. 2
- [8] Junuk Cha, Jihyeon Kim, Jae Shin Yoon, and Seungryul Baek. Text2hoi: Text-guided 3d motion generation for hand-object interaction. In *CVPR*, pages 1577–1585, 2024. 1, 3
- [9] Yu-Wei Chao, Wei Yang, Yu Xiang, Pavlo Molchanov, Ankur Handa, Jonathan Tremblay, Yashraj S Narang, Karl Van Wyk, Umar Iqbal, Stan Birchfield, et al. Dexycb: A benchmark for capturing hand grasping of objects. In *CVPR*, pages 9044–9053, 2021. 1, 3
- [10] Hila Chefer, Uriel Singer, Amit Zohar, Yuval Kirstain, Adam Polyak, Yaniv Taigman, Lior Wolf, and Shelly Sheynin. VideoJAM: Joint appearance-motion representations for enhanced motion generation in video models. In *ICML*, 2025. 2, 3
- [11] Junyi Chen, Haoyi Zhu, Xianglong He, Yifan Wang, Jianjun Zhou, Wenzheng Chang, Yang Zhou, Zizun Li, Zhoujie Fu, Jiangmiao Pang, et al. Deepverse: 4d autoregressive video generation as a world model. *arXiv preprint arXiv:2506.01103*, 2025. 2
- [12] Hyung Won Chung, Xavier Garcia, Adam Roberts, Yi Tay, Orhan Firat, Sharan Narang, and Noah Constant. Unimax: Fairer and more effective language sampling for large-scale multilingual pretraining. In *ICLR*, 2023. 4
- [13] Enric Corona, Andrei Zanfir, Eduard Gabriel Bazavan, Nikos Kolotouros, Thiemo Alldieck, and Cristian Sminchisescu. Vlogger: Multimodal diffusion for embodied avatar synthesis. In *CVPR*, pages 15896–15908, 2025. 2
- [14] Lingwei Dang, Ruizhi Shao, Hongwen Zhang, Wei Min, Yebin Liu, and Qingyao Wu. Svimo: Synchronized diffusion for video and motion generation in hand-object interaction scenarios. *NeurIPS*, 2025. 2, 3, 6, 7
- [15] Christian Diller and Angela Dai. Cg-hoi: Contact-guided 3d human-object interaction generation. In *CVPR*, pages 19888–19901, 2024. 1, 3
- [16] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *ICML*, 2024. 3, 4
- [17] Zicong Fan, Omid Taheri, Dimitrios Tzionas, Muhammed Kocabas, Manuel Kaufmann, Michael J Black, and Otmar Hilliges. Arctic: A dataset for dexterous bimanual hand-object manipulation. In *CVPR*, pages 12943–12954, 2023. 1
- [18] Rao Fu, Dingxi Zhang, Alex Jiang, Wanjia Fu, Austin Funk, Daniel Ritchie, and Srinath Sridhar. Gigahands: A massive annotated dataset of bimanual hand activities. In *CVPR*, pages 17461–17474, 2025. 1, 3
- [19] Jiawei Gao, Ziqin Wang, Zeqi Xiao, Jingbo Wang, Tai Wang, Jinkun Cao, Xiaolin Hu, Si Liu, Jifeng Dai, and Jiangmiao Pang. Coohei: Learning cooperative human-object interaction with manipulated object dynamics. In *NeurIPS*, pages 79741–79763, 2024. 1
- [20] Zekai Gu, Rui Yan, Jiahao Lu, Peng Li, Zhiyang Dou, Chenyang Si, Zhen Dong, Qifeng Liu, Cheng Lin, Ziwei Liu, et al. Diffusion as shader: 3d-aware video diffusion for versatile video generation control. In *SIGGRAPH*, pages 1–12, 2025. 2, 4, 6, 7, 1
- [21] Li Hu. Animate anyone: Consistent and controllable image-to-video synthesis for character animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8153–8163, 2024. 2, 3
- [22] Ziqi Huang, Yanan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *CVPR*, pages 21807–21818, 2024. 6
- [23] Hyeonho Jeong, Suhyeon Lee, and Jong Chul Ye. Reangle-a-video: 4d video generation as video-to-video translation. In *ICCV*, 2025. 2
- [24] Zeren Jiang, Chuanxia Zheng, Iro Laina, Diane Larlus, and Andrea Vedaldi. Geo4d: Leveraging video generators for geometric 4d scene reconstruction. In *ICCV*, 2025. 6, 7
- [25] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 4
- [26] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video

- generative models. *arXiv preprint arXiv:2412.03603*, 2024. 2, 3
- [27] Nilesh Kulkarni, Davis Rempe, Kyle Genova, Abhijit Kundu, Justin Johnson, David Fouhey, and Leonidas Guibas. Nifty: Neural object interaction fields for guided human motion synthesis. In *CVPR*, pages 947–957, 2024. 1, 3
- [28] Jihyun Lee, Shunsuke Saito, Giljoo Nam, Minhyuk Sung, and Tae-Kyun Kim. Interhandgen: Two-hand interaction generation via cascaded reverse diffusion. In *CVPR*, pages 527–537, 2024. 3
- [29] Jiaman Li, Jiajun Wu, and C Karen Liu. Object motion guided human motion synthesis. *ACM Transactions on Graphics (TOG)*, 42(6):1–11, 2023. 1
- [30] Jiaman Li, Alexander Clegg, Roozbeh Mottaghi, Jiajun Wu, Xavier Puig, and C Karen Liu. Controllable human-object interaction synthesis. In *ECCV*, pages 54–72. Springer, 2024. 1, 3
- [31] Quanzhou Li, Jingbo Wang, Chen Change Loy, and Bo Dai. Task-oriented human-object interactions generation with implicit neural representations. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3035–3044, 2024. 3
- [32] Kun Liu, Qi Liu, Xinchun Liu, Jie Li, Yongdong Zhang, Jiebo Luo, Xiaodong He, and Wu Liu. Hoigen-1m: A large-scale dataset for human-object interaction video generation. In *CVPR*, pages 24001–24010, 2025. 1
- [33] Siqi Liu, Yong-Lu Li, Zhou Fang, Xinpeng Liu, Yang You, and Cewu Lu. Primitive-based 3d human-object interaction modelling and programming. In *AAAI*, pages 3711–3719, 2024. 3
- [34] Tianqi Liu, Zihao Huang, Zhaoxi Chen, Guangcong Wang, Shoukang Hu, Liao Shen, Huiqiang Sun, Zhiguo Cao, Wei Li, and Ziwei Liu. Free4d: Tuning-free 4d scene generation with spatial-temporal consistency. 2025. 2
- [35] Xueyi Liu and Li Yi. Geneoh diffusion: Towards generalizable hand-object interaction denoising via denoising diffusion. In *ICLR*, 2024. 1, 3
- [36] Yunze Liu, Yun Liu, Che Jiang, Kangbo Lyu, Weikang Wan, Hao Shen, Boqiang Liang, Zhoujie Fu, He Wang, and Li Yi. Hoi4d: A 4d egocentric dataset for category-level human-object interaction. In *CVPR*, pages 21013–21022, 2022. 1
- [37] Yun Liu, Haolin Yang, Xu Si, Ling Liu, Zipeng Li, Yuxiang Zhang, Yebin Liu, and Li Yi. Taco: Benchmarking generalizable bimanual tool-action-object understanding. In *CVPR*, pages 21740–21751, 2024. 3, 6
- [38] Yun Liu, Chengwen Zhang, Ruofan Xing, Bingda Tang, Bowen Yang, and Li Yi. Core4d: A 4d human-object-human interaction dataset for collaborative object rearrangement. In *CVPR*, pages 1769–1782, 2025. 1
- [39] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 851–866. 2023. 1
- [40] Hao Luo, Yicheng Feng, Wanpeng Zhang, Sipeng Zheng, Ye Wang, Haoqi Yuan, Jiazheng Liu, Chaoyi Xu, Qin Jin, and Zongqing Lu. Being-h0: vision-language-action pre-training from large-scale human videos. *arXiv preprint arXiv:2507.15597*, 2025. 1
- [41] Zhengyi Luo, Jinkun Cao, Sammy Christen, Alexander Winkler, Kris Kitani, and Weipeng Xu. Omnigrasp: Grasping diverse objects with simulated humanoids. In *NeurIPS*, pages 2161–2184, 2024. 1, 3
- [42] Guoqing Ma, Haoyang Huang, Kun Yan, Liangyu Chen, Nan Duan, Shengming Yin, Changyi Wan, Ranchen Ming, Xiaoniu Song, Xing Chen, et al. Step-video-t2v technical report: The practice, challenges, and future of video foundation model. *arXiv preprint arXiv:2502.10248*, 2025. 2
- [43] Youxin Pang, Ruizhi Shao, Jiajun Zhang, Hanzhang Tu, Yun Liu, Boyao Zhou, Hongwen Zhang, and Yebin Liu. Manivideo: Generating hand-object manipulation video with dexterous and generalizable grasping. In *CVPR*, pages 12209–12219, 2025. 2
- [44] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 2, 4
- [45] Xiangyu Peng, Zangwei Zheng, Chenhui Shen, Tom Young, Xinying Guo, Binluo Wang, Hang Xu, Hongxin Liu, Mingyan Jiang, Wenjun Li, et al. Open-sora 2.0: Training a commercial-level video generation model in 200 k. *arXiv preprint arXiv:2503.09642*, 2025. 2
- [46] Huaijin Pi, Sida Peng, Minghui Yang, Xiaowei Zhou, and Hujun Bao. Hierarchical generation of human-object interactions with diffusion probabilistic models. In *ICCV*, pages 15061–15073, 2023. 1
- [47] Yuzhe Qin, Yueh-Hua Wu, Shaowei Liu, Hanwen Jiang, Ruihan Yang, Yang Fu, and Xiaolong Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. In *ECCV*, pages 570–587. Springer, 2022. 1
- [48] Javier Romero, Dimitrios Tzionas, and Michael J. Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6), 2017. 1
- [49] Xuelun Shen, Zhipeng Cai, Wei Yin, Matthias Müller, Zijun Li, Kaixuan Wang, Xiaozhi Chen, and Cheng Wang. Gim: Learning generalizable image matcher from internet videos. In *ICLR*, 2024. 6
- [50] Omid Taheri, Nima Ghorbani, Michael J Black, and Dimitrios Tzionas. Grab: A dataset of whole-body human grasping of objects. In *ECCV*, pages 581–600. Springer, 2020. 1
- [51] HunyuanWorld Team, Zhenwei Wang, Yuhao Liu, Junta Wu, Zixiao Gu, Haoyuan Wang, Xuhui Zuo, Tianyu Huang, Wenhuan Li, Sheng Zhang, et al. Hunyuanworld 1.0: Generating immersive, explorable, and interactive 3d worlds from words or pixels. *arXiv preprint arXiv:2507.21809*, 2025. 2
- [52] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Fei Wu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 2, 3, 6, 7
- [53] Chaoyang Wang, Peiye Zhuang, Tuan Duc Ngo, Willi Menapace, Aliaksandr Siarohin, Michael Vasilkovsky, Ivan Skorokhodov, Sergey Tulyakov, Peter Wonka, and Hsin-Ying Lee. 4real-video: Learning generalizable photo-realistic 4d video diffusion. In *CVPR*, pages 17723–17732, 2025. 2

- [54] Rong Wang, Wei Mao, and Hongdong Li. Deepsimho: Stable pose estimation for hand-object interaction via physics simulation. In *NeurIPS*, pages 79685–79697, 2023. 1
- [55] Yinhuai Wang, Jing Lin, Ailing Zeng, Zhengyi Luo, Jian Zhang, and Lei Zhang. Physshoi: Physics-based imitation of dynamic human-object interaction. *arXiv preprint arXiv:2312.04393*, 2023. 3
- [56] Rundi Wu, Ruiqi Gao, Ben Poole, Alex Trevithick, Changxi Zheng, Jonathan T Barron, and Aleksander Holynski. Cat4d: Create anything in 4d with multi-view video diffusion models. In *CVPR*, pages 26057–26068, 2025. 2
- [57] Xiaoyang Wu, Li Jiang, Peng-Shuai Wang, Zhijian Liu, Xihui Liu, Yu Qiao, Wanli Ouyang, Tong He, and Hengshuang Zhao. Point transformer v3: Simpler faster stronger. In *CVPR*, pages 4840–4851, 2024. 5, 6
- [58] Sirui Xu, Zhengyuan Li, Yu-Xiong Wang, and Liang-Yan Gui. Interdiff: Generating 3d human-object interactions with physics-informed diffusion. In *ICCV*, pages 14928–14940, 2023. 1
- [59] Sirui Xu, Yu-Xiong Wang, Liangyan Gui, et al. Interdreamer: Zero-shot text to 3d dynamic human-object interaction. In *NeurIPS*, pages 52858–52890, 2024. 3
- [60] Sirui Xu, Dongting Li, Yucheng Zhang, Xiyan Xu, Qi Long, Ziyin Wang, Yunzhi Lu, Shuchang Dong, Hezi Jiang, Akshat Gupta, et al. Interact: Advancing large-scale versatile 3d human-object interaction generation. In *CVPR*, pages 7048–7060, 2025. 1
- [61] Sirui Xu, Hung Yu Ling, Yu-Xiong Wang, and Liang-Yan Gui. Intermimic: Towards universal whole-body control for physics-based human-object interactions. In *CVPR*, 2025. 3
- [62] Tian-Xing Xu, Xiangjun Gao, Wenbo Hu, Xiaoyu Li, Song-Hai Zhang, and Ying Shan. Geometrycrafter: Consistent geometry estimation for open-world videos with diffusion priors. In *ICCV*, 2025. 6, 7
- [63] Ziyi Xu, Ziyao Huang, Juan Cao, Yong Zhang, Xiaodong Cun, Qing Shuai, Yuchen Wang, Linchao Bao, Jintao Li, and Fan Tang. Anchorcrafter: Animate cyberanchors saling your products via human-object interacting video generation. *arXiv preprint arXiv:2411.17383*, 2024. 1, 2, 3
- [64] Mengqing Xue, Yifei Liu, Ling Guo, Shaoli Huang, and Changxing Ding. Guiding human-object interactions with rich geometry and relations. In *CVPR*, pages 22714–22723, 2025. 1
- [65] Lixin Yang, Kailin Li, Xinyu Zhan, Fei Wu, Anran Xu, Liu Liu, and Cewu Lu. Oakink: A large-scale knowledge repository for understanding hand-object interaction. In *CVPR*, pages 20953–20962, 2022. 1
- [66] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Yuxuan Zhang, Weihang Wang, Yean Cheng, Bin Xu, Xiaotao Gu, Yuxiao Dong, and Jie Tang. Cogvideox: Text-to-video diffusion models with an expert transformer. In *ICLR*, 2025. 2, 3
- [67] Chun-Han Yao, Yiming Xie, Vikram Voleti, Huaizu Jiang, and Varun Jampani. Sv4d 2.0: Enhancing spatio-temporal consistency in multi-view video diffusion for high-quality 4d generation. In *ICCV*, pages 13248–13258, 2025. 2, 6, 7
- [68] Hong-Xing Yu, Haoyi Duan, Charles Herrmann, William T Freeman, and Jiajun Wu. Wonderworld: Interactive 3d scene generation from a single image. In *CVPR*, pages 5916–5926, 2025. 2
- [69] Mark YU, Wenbo Hu, Jinbo Xing, and Ying Shan. Trajectorycrafter: Redirecting camera trajectory for monocular videos via diffusion models. In *ICCV*, 2025. 2
- [70] Wangbo Yu, Jinbo Xing, Li Yuan, Wenbo Hu, Xiaoyu Li, Zhipeng Huang, Xiangjun Gao, Tien-Tsin Wong, Ying Shan, and Yonghong Tian. Viewcrafter: Taming video diffusion models for high-fidelity novel view synthesis. *IEEE TPAMI*, 2025. 2
- [71] Xinyu Zhan, Lixin Yang, Yifei Zhao, Kangrui Mao, Hanlin Xu, Zenan Lin, Kailin Li, and Cewu Lu. Oakink2: A dataset of bimanual hands-object manipulation in complex task completion. In *CVPR*, pages 445–456, 2024. 1, 3
- [72] Jinlu Zhang, Yixin Chen, Zan Wang, Jie Yang, Yizhou Wang, and Siyuan Huang. Interactanything: Zero-shot human object interaction synthesis via llm feedback and object affordance parsing. In *CVPR*, pages 7015–7025, 2025. 3
- [73] Jiajun Zhang, Yuxiang Zhang, Liang An, Mengcheng Li, Hongwen Zhang, Zonghai Hu, and Yebin Liu. Manidext: Hand-object manipulation synthesis via continuous correspondence embeddings and residual-guided diffusion. *IEEE TPAMI*, 2025. 1, 3
- [74] Xiaohan Zhang, Bharat Lal Bhatnagar, Sebastian Starke, Vladimir Guzov, and Gerard Pons-Moll. Couch: Towards controllable human-chair interactions. In *ECCV*, pages 518–535. Springer, 2022. 1
- [75] Zhenhao Zhang, Ye Shi, Lingxiao Yang, Suting Ni, Qi Ye, and Jingya Wang. Openhoi: Open-world hand-object interaction synthesis with multimodal large language model. In *NeurIPS*, 2025. 3
- [76] Haoyu Zhen, Qiao Sun, Hongxin Zhang, Junyan Li, Siyuan Zhou, Yilun Du, and Chuang Gan. Tesseract: learning 4d embodied world models. In *ICCV*, 2025. 2, 3
- [77] Shenhao Zhu, Junming Leo Chen, Zuozhuo Dai, Zilong Dong, Yinghui Xu, Xun Cao, Yao Yao, Hao Zhu, and Siyu Zhu. Champ: Controllable and consistent human image animation with 3d parametric guidance. In *ECCV*, pages 145–162. Springer, 2024. 2, 3

# SyncMV4D: Synchronized Multi-view Joint Diffusion of Appearance and Motion for Hand-Object Interaction Synthesis

## Supplementary Material

### 6. More Ablation Studies

In Sec. 3.3 and Fig. 3, we detail our motion representation. The “tracking video” in DaS [20] lacks per-frame object depth, making it a pseudo-3D or “2.5D” representation. In contrast, our representation includes explicit depth for each frame, offering stronger 3D awareness. To evaluate how these representations affect final performance, we conduct an ablation study. We replace our 4D point tracks with the DaS’s tracking video. Since this representation cannot recover per-frame depth or point trajectories, the DPA module becomes irrelevant. We therefore remove both DPA and the mutual enhancement circle, keeping only the MJD module. This variant is called “Ours-MJD-TrackingVideo”. Results in Tab. 4 show that this representation causes a clear decline in multi-view consistency. The main reason is that the 2.5D representation cannot adequately recover multi-view 3D geometry.

Table 4. Ablation studies on the motion representation.

Setting	Mat. Pix.	RPE
Ours-MJD-TrackingVideo	483.8	-
Ours-MJD-PointTracks	503.1	46.3
Ours	<b>529.4</b>	<b>32.7</b>

### 7. Additional Demonstrations

In this work, we focus on the task of image-to-video and action generation. To provide a more vivid and comprehensive presentation of our results, we have included a **supplementary video** in the appendix. This video showcases a wide range of generated examples by our method, along with systematic comparisons against several baseline approaches. These visual comparisons clearly demonstrate the superior performance of our method in terms of visual fidelity, dynamic plausibility, and multi-view consistency. Furthermore, in the final segment of the video, we present generation results on unseen hand-object interaction (HOI) tasks as well as on real-world captured data, highlighting the generalization capability of our proposed approach.