

OmniDrag: Enabling Motion Control for Omnidirectional Image-to-Video Generation

WeiQi Li^{*1,2}, Shijie Zhao^{†2}, Chong Mou¹, Xuhan Sheng¹, Zhenyu Zhang¹,
Qian Wang¹, Junlin Li², Li Zhang², Jian Zhang^{†1}

¹ School of Electronic and Computer Engineering, Peking University, ² ByteDance Inc
<https://lwq20020127.github.io/OmniDrag>

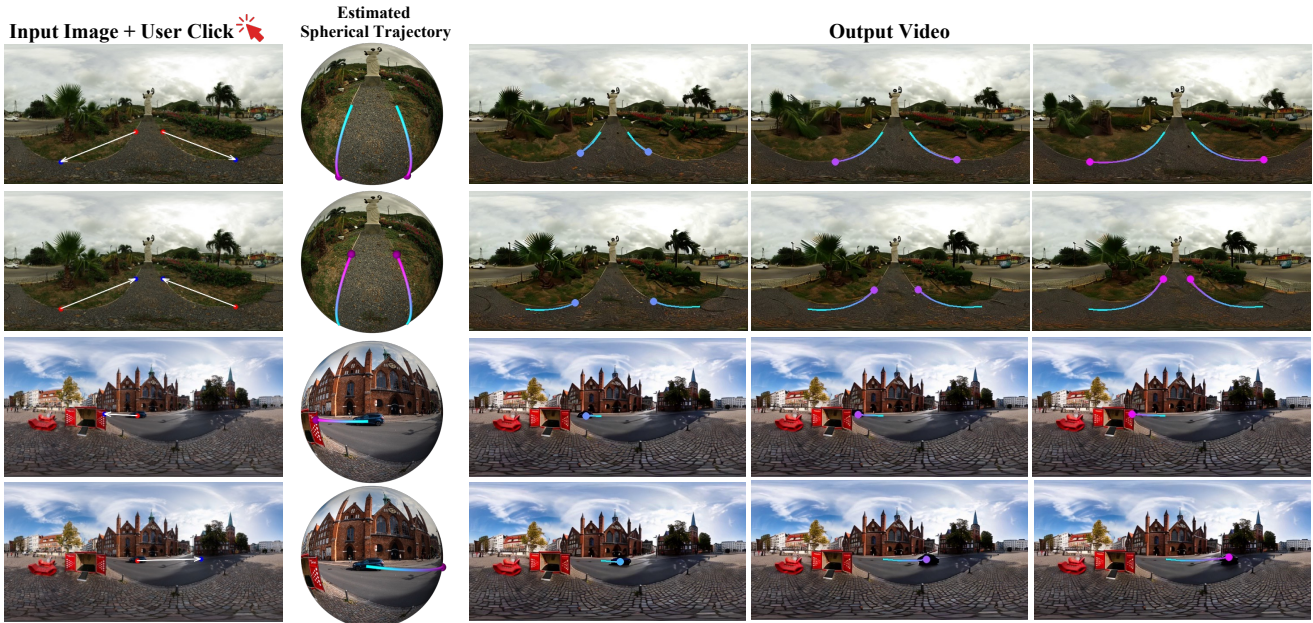


Figure 1. **Omnidirectional videos generated by proposed OmniDrag.** It enables drag-style synthesis from a reference omnidirectional image and user-specified points, providing both scene-level (top) and object-level (bottom) accurate, high-quality controllable generation.

Abstract

As virtual reality gains popularity, the demand for controllable creation of immersive and dynamic omnidirectional videos (ODVs) is increasing. While previous text-to-ODV generation methods achieve impressive results, they struggle with content inaccuracies and inconsistencies due to reliance solely on textual inputs. Although recent motion control techniques provide fine-grained control for video generation, directly applying these methods to ODVs often results in spatial distortion and unsatisfactory performance, especially with complex spherical motions. To tackle these challenges, we propose **OmniDrag**, the first approach enabling both scene- and object-level motion control for accurate, high-quality omnidirectional image-to-video generation. Building on pretrained video diffusion models, we in-

roduce an omnidirectional control module, which is jointly fine-tuned with temporal attention layers to effectively handle complex spherical motion. In addition, we develop a novel spherical motion estimator that accurately extracts motion-control signals and allows users to perform drag-style ODV generation by simply drawing handle and target points. We also present a new dataset, named **Move360**, addressing the scarcity of ODV data with large scene and object motions. Experiments demonstrate the significant superiority of OmniDrag in achieving holistic scene-level and fine-grained object-level control for ODV generation. The project page is available at <https://lwq20020127.github.io/OmniDrag>.

1. Introduction

Omnidirectional video (ODV) [66, 80], also known as 360° or panoramic video, has gained increasing attention due to its immersive and interactive capabilities, as well as its wide

* This work was done during the internship at ByteDance.

† Corresponding author.

applications in virtual and augmented reality. It provides a full $360^\circ \times 180^\circ$ field of view and is typically captured using an array of high-resolution fisheye cameras. Such a process is expensive in terms of both time and hardware resources in real-world scenarios [1]. Therefore, there is an urgent need for developing ODV generation methods.

In the field of 2D video generation, numerous diffusion-based models such as Gen-2 [18], Stable Video Diffusion (SVD) [5], and Sora [7] have achieved great success by leveraging powerful generative priors learned using large-scale training data and substantial computation resources. For ODV generation, 360DVD [61] introduces a plug-and-play 360-Adapter to enable text-to-ODV synthesis. However, this paradigm relies solely on text input, which often provides overly broad generation freedom and fails to precisely determine video frames, leading to inaccurate and inconsistent content control. While 360DVD offers optical flow-based control, obtaining ODV optical flow is challenging for users [52], thus limiting its practical utility.

Recently, trajectory-based motion control has emerged as a more user-friendly and effective solution for controllable video generation. Drawing trajectories offers a simple yet flexible approach, compared to other control signals like optical flow or depth maps [21]. Based on this approach, efforts such as DragNUWA [71], MotionCtrl [62], and DragAnything [65] encode sparse trajectories or camera motions into latent space to effectively guide object movements. Despite these advanced methods for 2D video synthesis, directly applying them to ODV generation presents three significant challenges: *Firstly*, unlike controlling traditional 2D videos, which generally involve simple motions, the motion patterns in ODVs are often spherical. Previous approaches applied in this task can lead to spatial distortions in generated results due to their inability to model complicated spherical motions. *Secondly*, since ODVs are generally stored in equirectangular projection (ERP) format, controlling them is more difficult than controlling 2D videos, as drawing reasonable and precise spherical motion trajectories on ERP images is challenging for human users. *Thirdly*, existing ODV datasets contain samples with limited motion magnitudes, constraining the effectiveness of deep controllable ODV generation models when faced with users’ requirements for larger motion ranges.

To address these problems, in this paper, we propose **OmniDrag**, the first method to enable motion control for omnidirectional image-to-video generation based on powerful pretrained video diffusion models. As demonstrated in Fig. 1, OmniDrag achieves high-quality, controllable ODV generation with simple user input, enabling both scene-level and object-level drag-style control using a unified model. In OmniDrag, we introduce an omnidirectional controller that takes trajectory as input to provide fine-grained motion controllability. To effectively learn complicated spherical mo-

tions in ODVs, we propose jointly fine-tuning the temporal attention components with our controller. For accurate and easy motion control, we develop a novel spherical motion estimator (SME). During training, SME tracks object motion using an equal-area iso-latitude spherical point initialization [19] and samples through a filter based on spherical distance to capture important movements uniformly and accurately. During inference, SME estimates motion trajectories via spherical interpolation, allowing users to provide only the handle and target points. Furthermore, we introduce a new high-quality ODV dataset named **Move360**, featuring significant scene-level and object-level motions. Move360 comprises more than 1,500 video clips across diverse scenes, captured by an Insta360 Titan mounted on a filming car. Experiments show that training on Move360 enhances OmniDrag’s ability for scene-level movement.

In summary, our contributions are:

- We propose **OmniDrag**, a novel method enabling motion control for ODV generation. It learns spherical motion patterns by jointly fine-tuning an omnidirectional controller and temporal attention layers in the UNet denoiser.
- We develop a novel spherical motion estimator (SME) that accurately captures control signals during training and allows users to simply draw handle and target points during inference, providing user-friendly controllability.
- We introduce **Move360**, a new high-quality ODV dataset, featuring large camera and object movements in samples captured by an Insta360 Titan mounted on a filming car, which enhances OmniDrag’s scene-level controllability.
- Extensive experiments demonstrate OmniDrag’s effectiveness and superior performance in generating smooth and visually appealing ODVs under interactive motion control, including both scene- and object-level control.

2. Related Work

2.1. Controllable Image and Video Generation

Recent developments in diffusion models [17, 23] have significantly enhanced image and video generation capabilities. Leading image generation frameworks, such as Stable Diffusion [49], Imagen [50], and DALL-E2 [48], utilize textual inputs to guide the generation process. Approaches like ControlNet [73] and T2I-Adapter [43] incorporate additional control modules into these pre-trained diffusion models to achieve finer controllability. In video generation, early methods similarly depend on text-based conditions, as demonstrated by Video LDM [6], Imagen Video [24], and AnimateDiff [20]. However, text prompts often fall short in handling complex scenarios, prompting recent research [5, 18, 67, 75] to adopt image-based conditions for more precise and effective control. For example, Video ControlNet [11, 78] extends the ControlNet architecture to video generation by conditioning on sequences of control signals such as depth and edge maps. ControlNeXt [45] further im-

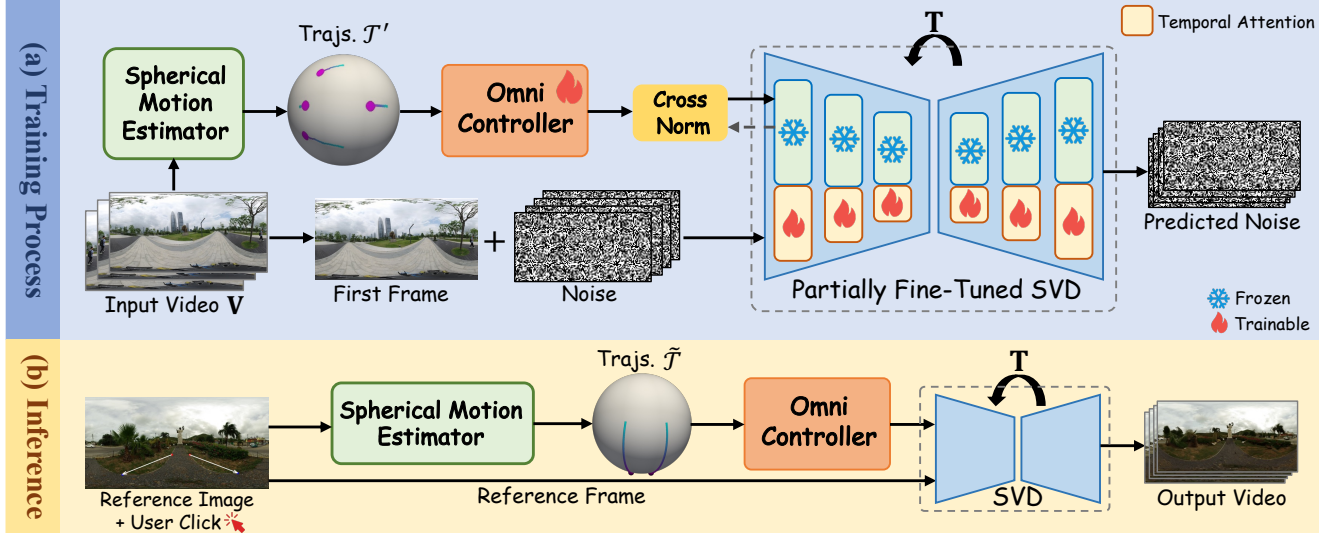


Figure 2. **Overall pipeline of proposed OmniDrag.** (a) During training, spherical motion is extracted by the proposed spherical motion estimator. The Omni Controller and temporal attention layers in the UNet denoiser are jointly fine-tuned. (b) During inference, OmniDrag allows users to simply select handle and target points on the reference image and generates ODVs with the corresponding motion.

proves ControlNet for lightweight image and video control guidance. Effective motion control is essential for producing coherent and dynamic videos. Current strategies employ trajectory-based methods like DragNUWA [71], MotionCtrl [62], DragAnything [65], and Tora [79], as well as box-based techniques [26, 41, 46, 60]. Unlike these 2D video generation methods that achieve desired motion dynamics by training additional motion controllers on frozen pre-trained video diffusion models, OmniDrag focuses on learning complex spherical motion patterns by jointly fine-tuning temporal attention layers in the base SVD model.

2.2. Omnidirectional Image and Video Generation

Generative adversarial network-based methods for producing omnidirectional images (ODIs) have been extensively explored [2, 3, 12, 14, 15, 35, 37, 38, 44, 54, 57, 63]. Recently, diffusion models have significantly advanced ODI generation [9, 31, 32, 34, 39, 58, 59, 64, 68, 69, 72, 74, 76, 77]. Specifically, PanoDiffusion [64] employs a dual-modal diffusion architecture incorporating RGB-D data to capture the spatial patterns of ODIs. PanFusion [72] introduces a dual-branch diffusion model that integrates global panorama and local perspective latent domains. LayerPano3D [68] decomposes a reference ODI into multiple layers at varying depth levels to facilitate explorable panoramic scenes. In the realm of omnidirectional video (ODV) generation, 360DVD [61] utilizes motion modeling modules [20] and 360Adapter to enable text-to-ODV generation with optical flow control. DiffPano [70] introduces a spherical epipolar-aware multi-view diffusion model. However, relying solely on text inputs often leads to inaccuracies and inconsistencies in the generated frames, and acquiring ODV

optical flow poses challenges for users, limiting broader applications. In contrast, our OmniDrag enables control through images and trajectories, providing accurate controllability with a user-friendly interface.

3. Methodology

In this section, we begin with a concise review of the employed base model Stable Video Diffusion (Sec. 3.1). Following this, we provide an overview of our OmniDrag (Sec. 3.2), illustrated in Fig. 2. We then elaborate on the Omni Controller and partial fine-tuning technique (Sec. 3.3) and proposed spherical motion estimator (Sec. 3.4). The proposed Move360 dataset is detailed in Sec. 3.5.

3.1. Preliminaries

Stable Video Diffusion (SVD) [5] is a high-quality and widely used image-to-video generation model. We adopt SVD as the base model for our proposed OmniDrag, to leverage its high-quality video generation capabilities. Specifically, given a reference image c_I , SVD generates a sequence of video frames of length L , starting with given c_I , denoted as $\mathbf{x} = \{x^0, x^1, \dots, x^{L-1}\}$. Following the latent denoising diffusion process in [49], a 3D UNet Φ_θ is used to denoise the sequence iteratively at timestep t :

$$\hat{z}_0 = \Phi_\theta(z_t, t, c_I), \quad (1)$$

where z_t is the latent representation of x_t obtained via an autoencoder [27, 56] as $z_t = \mathcal{E}(x_t)$, and \hat{z}_0 is the model's prediction of $z_0 = \mathcal{E}(x)$. To inject the reference image c_I into the main denoising branch, there are two paths: (1) c_I is embedded into tokens by the CLIP [47] image en-

coder and injected into the diffusion model through a cross-attention [49] mechanism. (2) \mathbf{c}_I is encoded into latent representation by the VAE encoder [27, 56] of the latent diffusion model and concatenated with the latent representations of each frame along the channel dimension. SVD parameterizes the learnable denoiser Φ_θ following the EDM-preconditioning [29] framework, as:

$$\Phi_\theta(\mathbf{z}_t, t, \mathbf{c}_I; \sigma) = c_{skip}(\sigma)\mathbf{z}_t + c_{out}(\sigma)F_\theta(c_{in}(\sigma)\mathbf{z}_t, t, \mathbf{c}_I; c_{noise}(\sigma)), \quad (2)$$

where σ is the noise level, F_θ is the denoising network, and c_{skip} , c_{out} , c_{in} and c_{noise} are hyper-parameters conditioned on σ . Finally, Φ_θ is trained via denoising score matching:

$$\mathbb{E}_{\mathbf{z}_0, t, \mathbf{n} \sim \mathcal{N}(0, \sigma^2)} [\lambda_\sigma \|\Phi_\theta(\mathbf{z}_0 + \mathbf{n}, t, \mathbf{c}_I) - \mathbf{z}_0\|_2^2]. \quad (3)$$

3.2. Overview of OmniDrag

An overview of OmniDrag is illustrated in Fig. 2. Built upon the pretrained SVD model, OmniDrag operates as follows. During training, the proposed spherical motion estimator (SME) first extracts trajectories from the input video. These trajectories are then fed into the Omni Controller, which is jointly fine-tuned with the temporal attention layers of the U-Net denoiser. During inference, users can simply select handle and target points on a reference image. OmniDrag then generates ODVs exhibiting the corresponding motion, enabling intuitive and precise motion control.

3.3. Omni Controller and Partial Fine-Tuning

Motivation. Temporal attention layers play important roles in motion pattern learning for diffusion models [4, 25, 30]. Existing 2D motion-guided methods [62, 65, 71] freeze the main branch of UNet model and utilize a trainable copy of the UNet encoder to inject the motion control. However, different from 2D videos, which generally involve simple motions, the motion patterns in ODVs are often spherical, introducing a significant gap. Consequently, merely training a control module with frozen main UNet branch leads to output videos with spatial distortions (Fig. 6 in Sec. 4). Therefore, we leverage a lightweight Omni Controller and jointly fine-tune the temporal attention layers in the denoising UNet to effectively learn the spherical motion pattern.

Method. Our OmniDrag employs a lightweight Omni controller instead of using a fully trainable copy of the main UNet encoder. Specifically, inspired by the recent ControlNeXt [45], we use a lightweight convolutional module only consisting of multiple ResNet blocks [22] to extract control signals. These controls are then integrated into the main denoising branch at a single selected middle block via an addition operation. Mathematically,

$$\mathbf{y}_m = \mathcal{F}_m(\mathbf{z}, \mathcal{F}_c(\mathbf{c}; \Theta_c); \Theta_m), \quad (4)$$

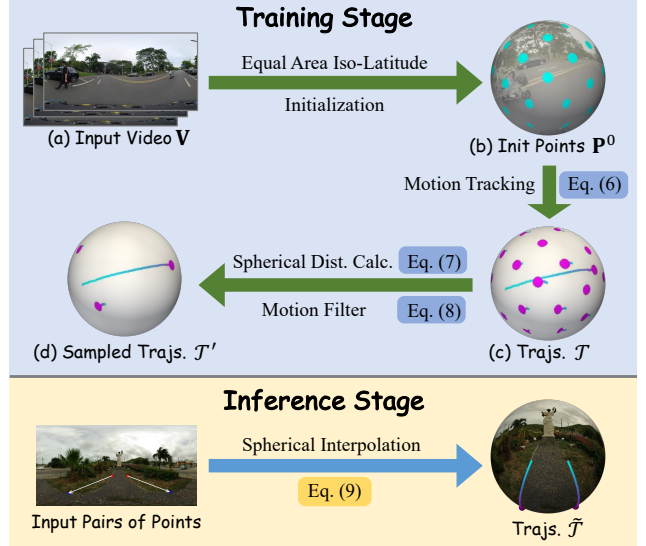


Figure 3. **Illustration of our spherical motion estimator (SME).** In the training stage, given the input video \mathbf{V} , \mathbf{P}^0 is firstly initialized through equal area iso-latitude pixelation. Then trajectories \mathcal{T} are tracked, and finally filtered as \mathcal{T}' according to spherical distance via Eqs. (6-8). During inference, given point pairs by users, the trajectories are estimated through spherical interpolation.

where \mathbf{y}_m represents the updated diffusion feature, \mathcal{F}_m and \mathcal{F}_c denote the main denoising U-Net and the Omni controller with parameters Θ_m and Θ_c , respectively. We propose to jointly fine-tune the temporal attention layers in the main UNet branch, whose parameters are denoted as $\Theta_t \subseteq \Theta_m$, as depicted in Fig. 2. This joint fine-tuning process is crucial for learning spherical motion patterns, resulting in the parameter set of OmniDrag $\Theta = \{\Theta_c, \Theta_t\}$. Additionally, we adopt the cross-normalization [45] technique to efficiently inject motion control signals into the main UNet branch during the fine-tuning process, aligning the distributions of the denoising and control features. Denoting the latent condition signals as $\mathbf{z}_c = \mathcal{F}_c(\mathbf{c}; \Theta_c)$, the final normalized control $\hat{\mathbf{z}}_c$ is calculated as:

$$\hat{\mathbf{z}}_c = \frac{\mathbf{z}_c - \boldsymbol{\mu}_m}{\sqrt{\boldsymbol{\sigma}_m^2 + \epsilon}} * \gamma, \quad (5)$$

where $\boldsymbol{\mu}_m$ and $\boldsymbol{\sigma}_m$ are the mean and variance of the latent \mathbf{z} from the main branch, respectively. γ is a hyper-parameter to scale the normalized value, and ϵ is a small constant for numerical stability. $\hat{\mathbf{z}}_c$ is finally integrated into the main denoising branch through addition. More details of the Omni Controller are provided in the supplementary materials.

3.4. Spherical Motion Estimator

Motivation. Precise motion control signals are essential for both training and inference phases. Existing 2D motion control methods typically initialize tracking points on images using uniform grids [42, 65, 71] and perform prob-

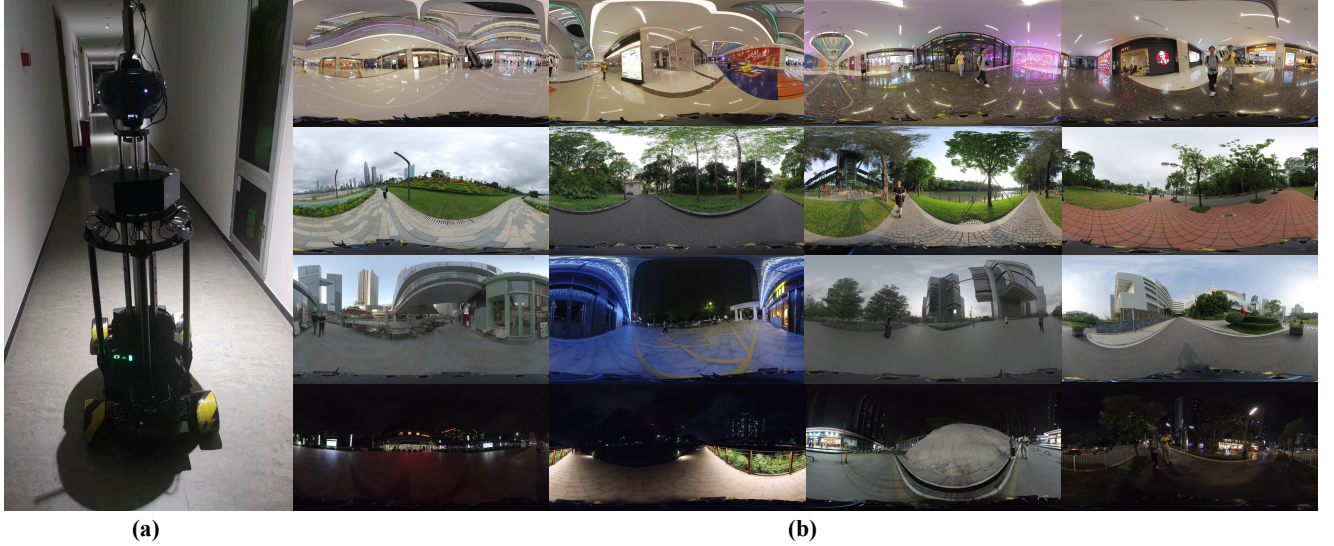


Figure 4. **Our Move360 dataset.** (a) We mount Insta360 Titan on a filming car, enabling its movement along four degrees of freedom. (b) Sample frames from the Move360 dataset showcasing a wide range of scenes, including indoor spaces, green landscapes, urban environments, and nighttime settings. This diversity in motion and environments offers a rich dataset for the community.

ability sampling based on motion distance. However, due to the spatial distortion inherent in equirectangular projection (ERP) [8, 13, 33, 36, 53], pixel density decreases near the poles, resulting in inaccurate and oversampled motion tracking in these areas. Moreover, directly calculating distances on the ERP does not reflect true spherical motion magnitudes, often causing primary motion patterns to be overlooked. Additionally, during inference, current methods require users to manually draw complete trajectories, which is challenging for users to draw reasonable spherical paths on the ERP image. To overcome these limitations, we propose the spherical motion estimator (SME). As illustrated in Fig. 3, SME captures more accurate spherical motion trajectories during training and offers user-friendly control capabilities during inference.

Method. During training, we extract motion trajectories \mathcal{T} from the input video to generate motion conditions. Let N_{init} denote the number of initialized tracking points and L the length of the video. Each trajectory $\mathbf{T}_j \in \mathcal{T}$ is defined as a sequence of spatial positions $\mathbf{T}_j = \{(x_j^i, y_j^i) | i \in \{0, 1, \dots, L-1\}\}$, where (x_j^i, y_j^i) represents the position of the j -th trajectory at frame i . To uniformly sample trajectories on the sphere, we propose to initialize the tracking points using the hierarchical equal area iso-latitude pixelation (HEALPix) grid [19], which provides a uniform distribution of grid points on the sphere, assigning the same area to each pixel, as shown in Fig. 3. Specifically, given a resolution parameter N_{side} , the HEALPix coordinate mapping function outputs a set of initialized points at frame 0, as $\mathbf{P}^0 = \{(x_j^0, y_j^0) | j \in \{0, 1, \dots, N_{init}-1\}\}$, where the total number of points N_{init} is $N_{init} = 12 \times N_{side}^2$. An object tracking function \mathcal{F}_t [28] is then applied to track the

motion of these initialized points \mathbf{P}^0 across the input video $\mathbf{V} \in \mathbb{R}^{L \times C \times H \times W}$, generating the corresponding motion trajectories as:

$$\mathcal{T} = \mathcal{F}_t(\mathbf{P}^0, \mathbf{V}), \quad (6)$$

where $\mathcal{T} \in \mathbb{R}^{N_{init} \times L \times 2} = \{\mathbf{T}_j | j \in \{0, 1, \dots, N_{init}-1\}\}$. Trajectories exhibiting larger motions are particularly beneficial for learning motion controllability. Therefore, we select trajectories with greater motion magnitudes. Instead of measuring motion magnitude in the ERP format, we propose to identify the primary motion of ODVs based on spherical distance, which is calculated as:

$$D(\mathbf{T}_j) = \arccos(\sin(\theta_j^0) \sin(\theta_j^{L-1}) + \cos(\theta_j^0) \cos(\theta_j^{L-1}) \cos(\phi_j^0 - \phi_j^{L-1})), \quad (7)$$

where the angles ϕ and θ are obtained by converting the ERP points (x, y) to spherical coordinates using $\phi = 2\pi x/W - \pi$ and $\theta = \pi y/H - \pi/2$. We then filter \mathcal{T} as:

$$\mathcal{T}' = \{\mathbf{T} \in \mathcal{T} | D(\mathbf{T}_j) > d_{th}\}, \quad (8)$$

where d_{th} is a threshold. Following [42], we use the normalized distances as sampling probabilities to randomly select N_{samp} trajectories from \mathcal{T}' . The final condition map \mathbf{c} is then obtained by applying a Gaussian filter to smooth the sampled trajectories. Consequently, \mathbf{c} serves as the conditional input in Eq. (4) to guide the generation process.

During inference, our objectives are: (1) to provide user-friendly interaction, and (2) to align the control signals with those used during training. Existing methods require users to draw motion trajectories on the reference image, which is feasible on 2D planar images. However, due to the spatial distortions of the ERP format, it is challenging for users

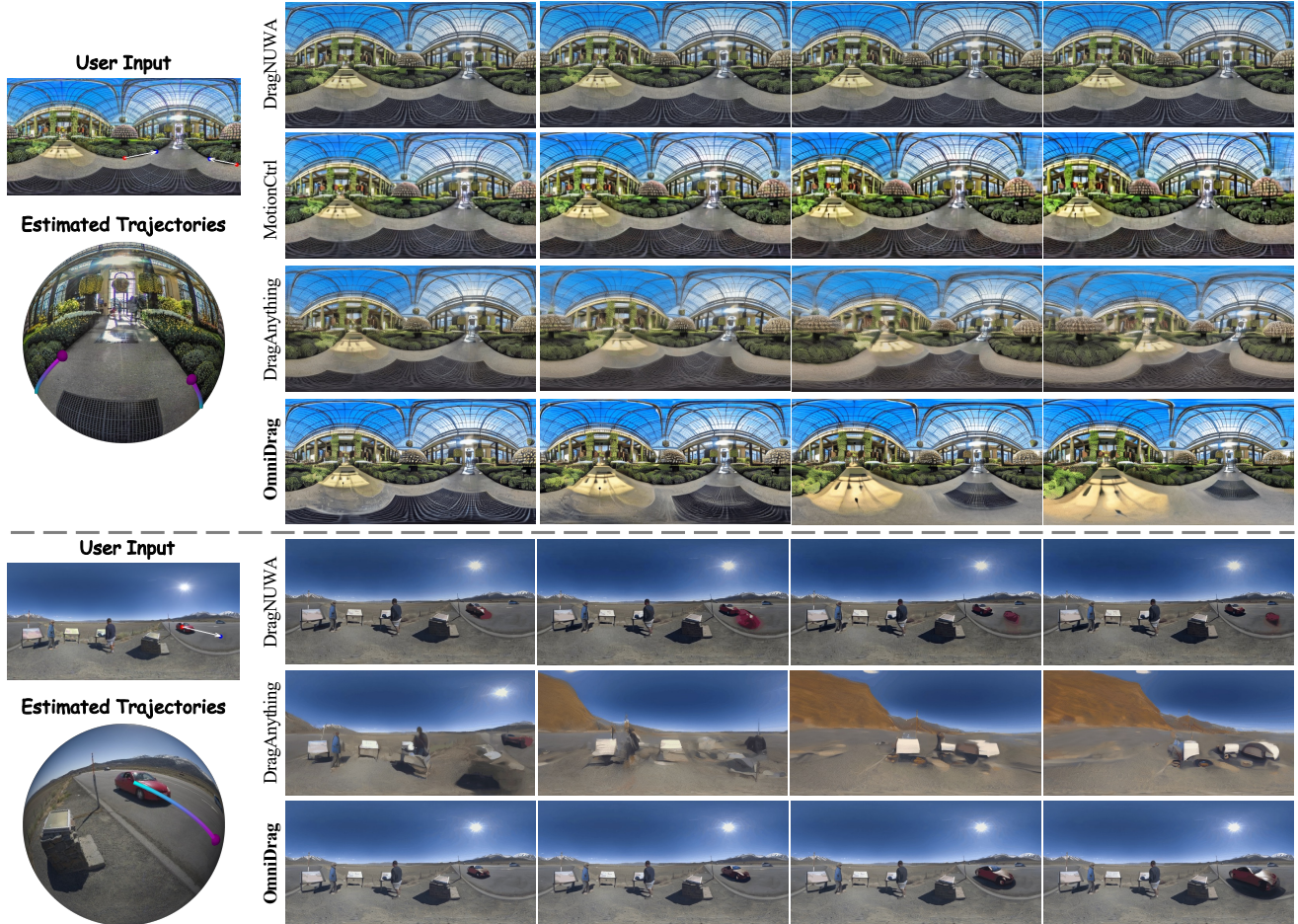


Figure 5. **Visual comparisons** between DragNUWA [71], MotionCtrl [62], DragAnything [65], and our OmniDrag. Our SME estimates reasonable trajectories on the sphere, and OmniDrag achieves precise and stable control under both **scene-level** (the **top** case: go forward on the road) and **object-level** (the **bottom** case: make the car move along the road) motion conditions, outperforming other methods.

to draw accurate spherical paths on the reference ERP image. To address this challenge, we introduce an innovative approach where users only need to specify the handle and target points, and the entire trajectory is then automatically estimated through spherical interpolation. Mathematically, denoting a pair of handle and target points as (x^0, y^0) and (x^{L-1}, y^{L-1}) , these points are firstly transformed to spherical coordinates (θ^0, ϕ^0) and $(\theta^{L-1}, \phi^{L-1})$. Then, the intermediate points (θ^i, ϕ^i) are calculated as:

$$\begin{cases} \theta^i = \arcsin\left(\frac{\sin((1-t_i)\omega)\sin\theta^0 + \sin(t_i\omega)\sin\theta^{L-1}}{\sin\omega}\right) \\ \phi^i = \phi^0 + t_i(\phi^{L-1} - \phi^0), \end{cases} \quad (9)$$

where ω is the spherical distance between these two points calculated as Eq. (7), and $t_i = i/L$ is the interpolation factor. Finally, these points are transformed back to ERP coordinates, and combined with the handle and target points to obtain $\tilde{\mathcal{T}} \in \mathbb{R}^{N_p \times L \times 2}$, where N_p is the number of point

pairs provided by the user. This process aligns inference-time control signals with those used during training.

3.5. Move360 Dataset

Training OmniDrag requires ODV datasets with high-quality motion. However, existing ODV datasets offer limited motion quality and magnitude due to their data acquisition methods. Specifically, videos in WEB360 [61] are primarily collected from the ‘‘AirPano VR’’ YouTube channel. These videos are obtained through aerial photography and contain watermarks, resulting in limited motion patterns and quality. The 360+x dataset [10] includes multiple scenes of ODVs from a third-person perspective. However, most videos in 360+x are filmed with a stationary camera, which is not conducive to learning motion. As shown in Fig. 6, training with existing datasets results in OmniDrag lacking scene-level control capabilities. To address these issues and cover the absence of high-quality panoramic video datasets with large motions, we introduce a new ODV

Table 1. **Quantitative comparisons** between our OmniDrag and other methods. We employ automatic metrics (FVD [55], FID [51] and ObjMC [62]) on both ERP format and final horizontal eight viewports. We also conduct a human evaluation to assess the performance. Throughout this paper, the best and second-best results are highlighted in **bold red** and underlined blue, respectively.

Method	ERP Image			Horizontal 8 viewports		Human Evaluation	
	FID↓	FVD↓	ObjMC↓	FID↓	FVD↓	Overall ↑	Motion Matching ↑
DragNUWA [71]	164.84	<u>1015.32</u>	0.418	<u>96.31</u>	<u>379.16</u>	<u>15.0%</u>	9.5 %
DragAnything [65]	182.63	1113.95	<u>0.085</u>	109.76	401.43	9.3%	<u>14.4%</u>
OmniDrag (Ours)	<u>171.41</u>	933.73	0.044	95.62	322.22	75.7%	76.1%

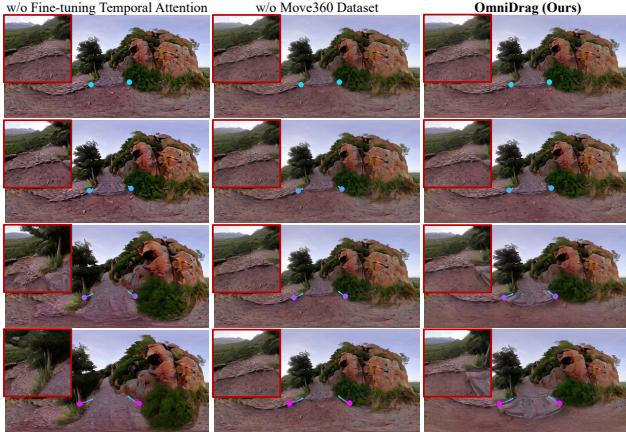


Figure 6. **Ablation study** on jointly fine-tuning temporal attention layers, and training with proposed Move360 dataset. For each ERP image, we show a corresponding viewport at specific perspective.

dataset named Move360. Specifically, we mount an Insta360 Titan camera on a filming car. The Insta360 Titan features eight 200°F3.2 fisheye cameras. The captured circular videos are subsequently de-warped and stitched using optical flow. Ultimately, we obtained ODVs at a resolution of 7680 × 3840 (8K) with a frame rate of 30 FPS. Moreover, our filming car allows the camera to move forward and backward, left and right, up and down, and rotate 360 degrees horizontally. These four degrees of freedom provide flexibility in capturing immersive content from various angles and positions. The original video has a duration of approximately 20 hours with 6TB in size. We curated the data based on scene and content quality, resulting in 1,580 clips from over 60 scenes, each consisting of 100 frames. Move360 contains a wide range of scenes as shown in Fig. 4, offering a rich dataset for training models requiring high-quality ODV content. More video samples from Move360 are provided in the supplementary materials.

4. Experiments

4.1. Experimental Setup

Implementation Details. We choose stable video diffusion (SVD) model [5] as our base model. We use CoTracker [28] as the tracking function \mathcal{F}_t and train the OmniDrag on Move360 and WEB360 [61] dataset. In the training stage, we follow ReVideo [42] to sample the number of trajectory

samples N_{samp} randomly between 1 and 10, and optimize OmniDrag with Adam optimizer [40] for 40K iterations on 8 A100 GPUs, with a batch size of 4 for each GPU. The resolution is downsampled to 640 × 320, the learning rate is set to 1×10^{-5} , and it takes about 2 days for training. Besides, we adopt a latent rotation mechanism [61, 64] to enhance the warp-around consistency of ODVs.

Evaluation metrics. We follow MotionControl [62] to evaluate from the following two aspects: (1) The quality of results is assessed using the Fréchet Inception Distance (FID) [51] and Fréchet Video Distance (FVD) [55], which measure the visual quality and temporal coherence, respectively. (2) The motion control performance (ObjMC) is evaluated by the spherical distance between the trajectories of the generated videos and the user input trajectories. In addition, we conduct a human evaluation, in which thirty volunteers are asked to vote the best method for each sample from two aspects: overall quality and motion matching.

4.2. Comparison with State-of-the-Art Methods

We compare our OmniDrag with state-of-the-art video generation methods incorporating motion control, specifically DragNUWA [71], MotionCtrl [62], and DragAnything [65]. We conduct experiments under both scene-level and object-level control conditions. Because MotionCtrl does not support trajectory control in image-to-video generation, we compare it only under scene-level control, using the corresponding camera pose as conditional input. We select ODIs from ODISR [16] and SUN360 [66] datasets as reference images, and create twelve pairs of input as the test set.

The visual comparisons of some cases are shown in Fig. 5, including scene-level control (top) and object-level control (bottom). Due to the lack of prior knowledge of spherical motion patterns, DragNUWA fails in both cases, producing only slight movements. MotionCtrl generates camera pose transformations in a 2D image manner, while DragAnything produces violent motions that distort image content. In contrast, our OmniDrag performs well in both scene-level and object-level control, conforming to the distortion and motion pattern of ODVs and exhibiting good warp-around continuity. More visual results are provided in the supplementary materials. Quantitative comparisons are presented in Tab. 1. Note that we also compare metrics on horizontal 8 viewports, which represent users’ final content view. It can be seen that our OmniDrag achieves

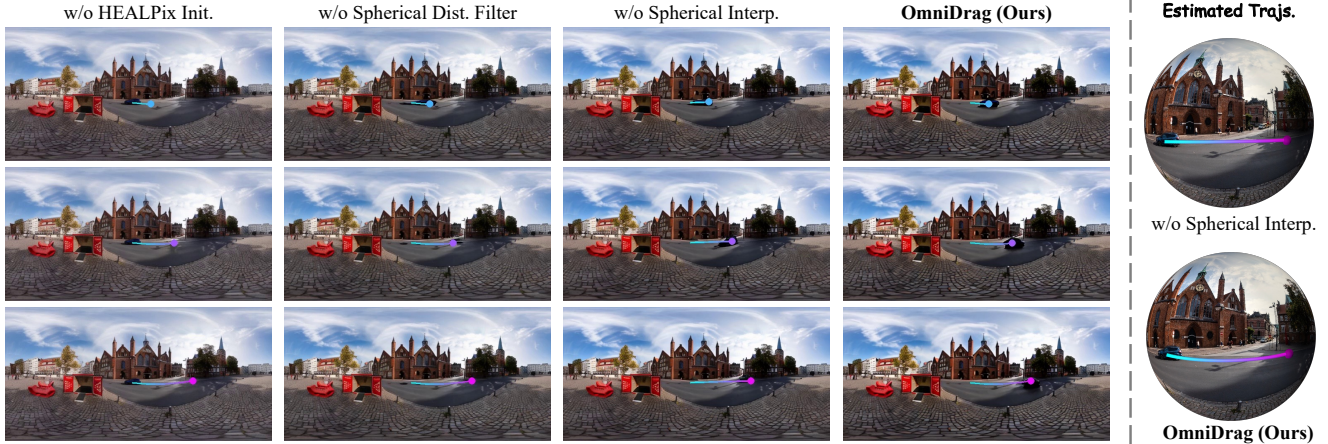


Figure 7. **Ablation study** on proposed spherical motion estimator (SME). The “w/o HEALPix init.” variant fails to control the car, the “w/o spherical dist. filter” variant generates unstable result, and the “w/o spherical interp.” variant leads to unintended path. In contrast, our OmniDrag leverages SME to obtain precise and reasonable trajectories during training and inference, achieving pleasant results.

Table 2. **Ablation study** on five variants of OmniDrag.

Method	ERP Image			Horizontal 8 viewports	
	FID↓	FVD↓	ObjMC↓	FID↓	FVD↓
w/o. Ft Temporal Attn.	182.72	982.41	0.080	97.12	332.97
w/o. Move360 Dataset	167.56	941.58	0.327	95.82	317.73
w/o. HEALPix Init.	<u>170.69</u>	<u>938.18</u>	0.226	95.33	324.05
w/o. Spherical Filter.	174.40	970.06	0.113	96.31	336.81
w/o. Spherical Interp.	174.13	965.47	<u>0.053</u>	96.94	342.18
OmniDrag (Ours)	171.41	933.73	0.044	<u>95.62</u>	<u>322.22</u>

the best FVD on ERP format and the best FID and FVD on the horizontal eight viewports, demonstrating the good quality of our generated results. Notably, DragNUWA typically generates minimal motion, resulting in a lower FID on ERP but a poor ObjMC score, whereas our OmniDrag achieves superior motion consistency. Furthermore, in human evaluations, OmniDrag exhibits clear advantages over other methods, demonstrating its superior performance in video quality and instruction comprehension.

4.3. Ablation Study

To validate the effectiveness of the proposed components in OmniDrag, including the joint fine-tuning strategy, the SME, and the Move360 dataset, we conduct ablation studies, as shown in Figs. 6 and 7, and Tab. 2.

Effect of jointly tuning temporal attention. To demonstrate the importance of jointly tuning the temporal attention layers, we create a variant where we freeze the entire main UNet denoising branch. The results in both ERP format and viewport are shown in Fig. 6. It can be observed that “w/o. fine-tuning” variant generates videos with only trivial 2D zoom-in effects, lacking omnidirectional properties, which leads to distorted viewport quality. This variant also results in higher FID and FVD scores, as shown in Tab. 2.

Effect of training on Move360 dataset. We evaluate another variant by training our OmniDrag only on the WEB360 [61] dataset. Although this variant achieves better

FID results, it exhibits poor motion control performance, as indicated in Tab. 2. The results in Fig. 6 further illustrate that training without datasets containing high-quality motion cannot provide scene-level controllability due to insufficient motion diversity. In contrast, training with our Move360 dataset enables accurate and stable scene-level control, significantly enhancing the model’s capabilities.

Effect of SME. To demonstrate the effectiveness of the proposed SME, we replace the HEALPix initialization, spherical distance calculation and spherical interpolation with 2D grid initialization, Euclidean distance and linear interpolation, respectively. The results are presented in Fig. 7 and Tab. 2. It is evident that removing these components significantly degrades the performance of motion control. Specifically, without HEALPix initialization and spherical distance filtering, the variant fails to control the object or generates unstable results, likely due to the lack of sufficient and accurate control signals during training. During inference, given the user’s input of handle and target points, SME estimates a reasonable spherical trajectory, whereas the linear interpolation variant generates a straight line on the sphere, resulting in ambiguous results, *e.g.*, the car runs off the road in this case, deviating from the intended path.

5. Conclusion

In this paper, we proposed **OmniDrag**, a novel diffusion-based approach for enabling motion control in omnidirectional image-to-video generation. We introduced an Omni Controller, which receives spherical trajectories as input, allowing for easy drag-style control. To effectively learn complex spherical motion patterns, we proposed jointly fine-tuning the controller and temporal layers in the diffusion denoising UNet. Additionally, we designed a spherical motion estimator to capture accurate control signals during training and provide user-friendly interaction during inference.

Furthermore, we collected Move360, a new high-quality ODV dataset featuring significant motion content, which enhances OmniDrag’s scene-level controllability. Experiments manifested that OmniDrag achieves state-of-the-art performance in both scene- and object-level motion control. **Limitations.** Although OmniDrag achieves promising results, its generation quality is constrained by the base SVD model in certain scenarios. Moreover, decoupling camera- and object-level motion is an open problem for future work.

References

- [1] Hao Ai, Zidong Cao, Jinjing Zhu, Haotian Bai, Yucheng Chen, and Lin Wang. Deep learning for omnidirectional vision: A survey and new perspectives. *arXiv preprint arXiv:2205.10468*, 2022. 2
- [2] Hao Ai, Zidong Cao, Haonan Lu, Chen Chen, Jian Ma, Pengyuan Zhou, Tae-Kyun Kim, Pan Hui, and Lin Wang. Dream360: Diverse and immersive outdoor virtual scene creation via transformer-based 360° image outpainting. *IEEE Transactions on Visualization and Computer Graphics (TVCG)*, 2024. 3
- [3] Naofumi Akimoto, Yuhi Matsuo, and Yoshimitsu Aoki. Diverse plausible 360-degree image outpainting for efficient 3ddeg background creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11441–11450, 2022. 3
- [4] Jianhong Bai, Tianyu He, Yuchi Wang, Junliang Guo, Haoji Hu, Zuozhu Liu, and Jiang Bian. Uniedit: A unified tuning-free framework for video motion and appearance editing. *arXiv preprint arXiv:2402.13185*, 2024. 4
- [5] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 2, 3, 7
- [6] 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 (CVPR)*, pages 22563–22575, 2023. 2
- [7] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. <https://openai.com/research/video-generation-models-as-world-simulators>, 2024. 2
- [8] Mingdeng Cao, Chong Mou, Fanghua Yu, Xintao Wang, Yinqiang Zheng, Jian Zhang, Chao Dong, Gen Li, Ying Shan, Radu Timofte, et al. Ntire 2023 challenge on 360deg omnidirectional image and video super-resolution: Datasets, methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop (CVPRW)*, pages 1731–1745, 2023. 5
- [9] Bin Chen, Zhenyu Zhang, Weiqi Li, Chen Zhao, Jiwen Yu, Shijie Zhao, Jie Chen, and Jian Zhang. Invertible diffusion models for compressed sensing. *arXiv preprint arXiv:2403.17006*, 2024. 3
- [10] Hao Chen, Yuqi Hou, Chenyuan Qu, Irene Testini, Xiaohan Hong, and Jianbo Jiao. 360+x: A panoptic multimodal scene understanding dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19373–19382, 2024. 6
- [11] Weifeng Chen, Yatai Ji, Jie Wu, Hefeng Wu, Pan Xie, Jiashi Li, Xin Xia, Xuefeng Xiao, and Liang Lin. Control-a-video: Controllable text-to-video generation with diffusion models. *arXiv preprint arXiv:2305.13840*, 2023. 2
- [12] Zhaoxi Chen, Guangcong Wang, and Ziwei Liu. Text2light: Zero-shot text-driven hdr panorama generation. *ACM Transactions on Graphics (TOG)*, 41(6):1–16, 2022. 3
- [13] Ming Cheng, Haoyu Ma, Qiufang Ma, Xiaopeng Sun, Weiqi Li, Zhenyu Zhang, Xuhan Sheng, Shijie Zhao, Junlin Li, and Li Zhang. Hybrid transformer and cnn attention network for stereo image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop (CVPRW)*, pages 1702–1711, 2023. 5
- [14] Yen-Chi Cheng, Chieh Hubert Lin, Hsin-Ying Lee, Jian Ren, Sergey Tulyakov, and Ming-Hsuan Yang. Inout: Diverse image outpainting via gan inversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11431–11440, 2022. 3
- [15] Mohammad Reza Karimi Dastjerdi, Yannick Hold-Geoffroy, Jonathan Eisenmann, Siavash Khodadadeh, and Jean-François Lalonde. Guided co-modulated gan for 360° field of view extrapolation. In *2022 International Conference on 3D Vision (3DV)*, pages 475–485. IEEE, 2022. 3
- [16] Xin Deng, Hao Wang, Mai Xu, Yichen Guo, Yuhang Song, and Li Yang. Lau-net: Latitude adaptive upscaling network for omnidirectional image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9189–9198, 2021. 7
- [17] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 34:8780–8794, 2021. 2
- [18] Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and content-guided video synthesis with diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7346–7356, 2023. 2
- [19] Krzysztof M Gorski, Eric Hivon, Anthony J Banday, Benjamin D Wandelt, Frode K Hansen, Mstvos Reinecke, and Matthias Bartelmann. Healpix: A framework for high-resolution discretization and fast analysis of data distributed on the sphere. *The Astrophysical Journal*, 622(2):759, 2005. 2, 5
- [20] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024. 2, 3
- [21] Zekun Hao, Xun Huang, and Serge Belongie. Controllable video generation with sparse trajectories. In *Proceedings of*

- the *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7854–7863, 2018. 2
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 4, 13
- [23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 33: 6840–6851, 2020. 2
- [24] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2
- [25] Teng Hu, Jiangning Zhang, Ran Yi, Yating Wang, Hongrui Huang, Jieyu Weng, Yabiao Wang, and Lizhuang Ma. Motionmaster: Training-free camera motion transfer for video generation. *arXiv preprint arXiv:2404.15789*, 2024. 4
- [26] Yash Jain, Anshul Nasery, Vibhav Vineet, and Harkirat Behl. Peekaboo: Interactive video generation via masked-diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8079–8088, 2024. 3
- [27] Xuhui Jia, Yang Zhao, Kelvin CK Chan, Yandong Li, Han Zhang, Boqing Gong, Tingbo Hou, Huisheng Wang, and Yu-Chuan Su. Taming encoder for zero fine-tuning image customization with text-to-image diffusion models. *arXiv preprint arXiv:2304.02642*, 2023. 3, 4
- [28] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. Co-tracker: It is better to track together. *Proceedings of the European Conference on Computer Vision (ECCV)*, 2024. 5, 7, 13
- [29] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 35:26565–26577, 2022. 4
- [30] Max Ku, Cong Wei, Weiming Ren, Harry Yang, and Wenhui Chen. Anyv2v: A tuning-free framework for any video-to-video editing tasks. *arXiv preprint arXiv:2403.14468*, 2024. 4
- [31] Jialu Li and Mohit Bansal. Panogen: Text-conditioned panoramic environment generation for vision-and-language navigation. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 3
- [32] Renjie Li, Panwang Pan, Bangbang Yang, Dejie Xu, Shijie Zhou, Xuanyang Zhang, Zeming Li, Achuta Kadambi, Zhangyang Wang, and Zhiwen Fan. 4k4dgen: Panoramic 4d generation at 4k resolution. *arXiv preprint arXiv:2406.13527*, 2024. 3
- [33] Runyi Li, Xuhan Sheng, Weiqi Li, and Jian Zhang. Omnisr: Zero-shot omnidirectional image super-resolution using stable diffusion model. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 198–216. Springer, 2024. 5
- [34] Weiqi Li, Bin Chen, Shuai Liu, Shijie Zhao, Bowen Du, Yongbing Zhang, and Jian Zhang. D3c2-net: Dual-domain deep convolutional coding network for compressive sensing. *IEEE Transactions on Circuits and Systems for Video Technology*, 2024. 3
- [35] Wenrui Li, Yapeng Mi, Fucheng Cai, Zhe Yang, Wangmeng Zuo, Xingtao Wang, and Xiaopeng Fan. Scenedreamer360: Text-driven 3d-consistent scene generation with panoramic gaussian splatting. *arXiv preprint arXiv:2408.13711*, 2024. 3
- [36] Weiqi Li, Shijie Zhao, Bin Chen, Xinhua Cheng, Junlin Li, Li Zhang, and Jian Zhang. Resvr: Joint rescaling and viewport rendering of omnidirectional images. In *Proceedings of the 32nd ACM International Conference on Multimedia (ACM MM)*, pages 78–87, 2024. 5
- [37] Chieh Hubert Lin, Chia-Che Chang, Yu-Sheng Chen, Da-Cheng Juan, Wei Wei, and Hwann-Tzong Chen. Cocogan: Generation by parts via conditional coordinating. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4512–4521, 2019. 3
- [38] Chieh Hubert Lin, Hsin-Ying Lee, Yen-Chi Cheng, Sergey Tulyakov, and Ming-Hsuan Yang. Infinitygan: Towards infinite-pixel image synthesis. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2022. 3
- [39] Qin Liu, Letong Han, Rui Tan, Hongfei Fan, Weiqi Li, Hongming Zhu, Bowen Du, and Sicong Liu. Hybrid attention based residual network for pansharpening. *Remote Sensing*, 13(10):1962, 2021. 3
- [40] I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 7
- [41] Wan-Duo Kurt Ma, John P Lewis, and W Bastiaan Kleijn. Trailblazer: Trajectory control for diffusion-based video generation. *arXiv preprint arXiv:2401.00896*, 2023. 3
- [42] Chong Mou, Mingdeng Cao, Xintao Wang, Zhaoyang Zhang, Ying Shan, and Jian Zhang. Revideo: Remake a video with motion and content control. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 4, 5, 7
- [43] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 4296–4304, 2024. 2
- [44] Changgyoon Oh, Wonjune Cho, Yujeong Chae, Daehee Park, Lin Wang, and Kuk-Jin Yoon. Bips: Bi-modal indoor panorama synthesis via residual depth-aided adversarial learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 352–371. Springer, 2022. 3
- [45] Bohao Peng, Jian Wang, Yuechen Zhang, Wenbo Li, Ming-Chang Yang, and Jiaya Jia. Controlnext: Powerful and efficient control for image and video generation. *arXiv preprint arXiv:2408.06070*, 2024. 2, 4, 13
- [46] Haonan Qiu, Zhaoxi Chen, Zhouxia Wang, Yingqing He, Menghan Xia, and Ziwei Liu. Freetraj: Tuning-free trajectory control in video diffusion models. *arXiv preprint arXiv:2406.16863*, 2024. 3
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,

- Amanda Askeell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proceedings of the International Conference on Learning Representations (ICLR)*, pages 8748–8763. PMLR, 2021. 3
- [48] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 2
- [49] 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 (CVPR)*, pages 10684–10695, 2022. 2, 3, 4
- [50] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 35: 36479–36494, 2022. 2
- [51] Maximilian Seitzer. pytorch-fid: Fid score for pytorch, 2020. 7
- [52] Hao Shi, Yifan Zhou, Kailun Yang, Xiaoting Yin, Ze Wang, Yaozu Ye, Zhe Yin, Shi Meng, Peng Li, and Kaiwei Wang. Panoflow: Learning 360° optical flow for surrounding temporal understanding. *IEEE Transactions on Intelligent Transportation Systems (TITS)*, 24(5):5570–5585, 2023. 2
- [53] Xiaopeng Sun, Weiqi Li, Zhenyu Zhang, Qiufang Ma, Xuhan Sheng, Ming Cheng, Haoyu Ma, Shijie Zhao, Jian Zhang, Junlin Li, et al. Opdn: Omnidirectional position-aware deformable network for omnidirectional image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop (CVPRW)*, pages 1293–1301, 2023. 5
- [54] Piotr Teterwak, Aaron Sarna, Dilip Krishnan, Aaron Maschinot, David Belanger, Ce Liu, and William T Freeman. Boundless: Generative adversarial networks for image extension. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10521–10530, 2019. 3
- [55] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018. 7
- [56] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017. 3, 4
- [57] Guangcong Wang, Yinuo Yang, Chen Change Loy, and Zhiwei Liu. Stylelight: Hdr panorama generation for lighting estimation and editing. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 477–492. Springer, 2022. 3
- [58] Hai Wang, Xiaoyu Xiang, Yuchen Fan, and Jing-Hao Xue. Customizing 360-degree panoramas through text-to-image diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 4933–4943, 2024. 3
- [59] Jionghao Wang, Ziyu Chen, Jun Ling, Rong Xie, and Li Song. 360-degree panorama generation from few unregistered nfov images. In *Proceedings of the 31th ACM International Conference on Multimedia (ACM MM)*, 2023. 3
- [60] Jiawei Wang, Yuchen Zhang, Jiaxin Zou, Yan Zeng, Guoqiang Wei, Liping Yuan, and Hang Li. Boximator: Generating rich and controllable motions for video synthesis. *arXiv preprint arXiv:2402.01566*, 2024. 3
- [61] Qian Wang, Weiqi Li, Chong Mou, Xinhua Cheng, and Jian Zhang. 360dvd: Controllable panorama video generation with 360-degree video diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 3, 6, 7, 8
- [62] Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping Luo, and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 2, 3, 4, 6, 7, 13
- [63] Songsong Wu, Hao Tang, Xiao-Yuan Jing, Haifeng Zhao, Jianjun Qian, Nicu Sebe, and Yan Yan. Cross-view panorama image synthesis. *IEEE Transactions on Multimedia (TMM)*, 2022. 3
- [64] Tianhao Wu, Chuanxia Zheng, and Tat-Jen Cham. Panodiffusion: 360-degree panorama outpainting via diffusion. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024. 3, 7
- [65] Weijia Wu, Zhuang Li, Yuchao Gu, Rui Zhao, Yefei He, David Junhao Zhang, Mike Zheng Shou, Yan Li, Tingting Gao, and Di Zhang. Draganything: Motion control for anything using entity representation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 331–348. Springer, 2025. 2, 3, 4, 6, 7, 14, 15
- [66] Jianxiong Xiao, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Recognizing scene viewpoint using panoramic place representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2695–2702. IEEE, 2012. 1, 7
- [67] Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Xintao Wang, Tien-Tsin Wong, and Ying Shan. Dynamicrafter: Animating open-domain images with video diffusion priors. *Proceedings of the European Conference on Computer Vision (ECCV)*, 2023. 2
- [68] Shuai Yang, Jing Tan, Mengchen Zhang, Tong Wu, Yixuan Li, Gordon Wetzstein, Ziwei Liu, and Dahua Lin. Layerpano3d: Layered 3d panorama for hyper-immersive scene generation. *arXiv preprint arXiv:2408.13252*, 2024. 3
- [69] Shuzhou Yang, Yu Wang, Haijie Li, Jiarui Meng, Xiangdong Meng, and Jian Zhang. Fourier123: One image to high-quality 3d object generation with hybrid fourier score distillation. *arXiv preprint arXiv:2405.20669*, 2024. 3
- [70] Weicai Ye, Chenhao Ji, Zheng Chen, Junyao Gao, Xiaoshui Huang, Song-Hai Zhang, Wanli Ouyang, Tong He, Cairong Zhao, and Guofeng Zhang. Diffpano: Scalable and consistent text to panorama generation with spherical epipolar-aware diffusion. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 3

- [71] Shengming Yin, Chenfei Wu, Jian Liang, Jie Shi, Houqiang Li, Gong Ming, and Nan Duan. Dragnuwa: Fine-grained control in video generation by integrating text, image, and trajectory. *arXiv preprint arXiv:2308.08089*, 2023. [2](#), [3](#), [4](#), [6](#), [7](#), [14](#), [15](#)
- [72] Cheng Zhang, Qianyi Wu, Camilo Cruz Gambardella, Xiaoshui Huang, Dinh Phung, Wanli Ouyang, and Jianfei Cai. Taming stable diffusion for text to 360 panorama image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6347–6357, 2024. [3](#)
- [73] 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 (ICCV)*, pages 3836–3847, 2023. [2](#)
- [74] Qinsheng Zhang, Jiaming Song, Xun Huang, Yongxin Chen, and Ming-Yu Liu. Diffcollage: Parallel generation of large content with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10188–10198. IEEE, 2023. [3](#)
- [75] Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang, Deli Zhao, and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded diffusion models. *arXiv preprint arXiv:2311.04145*, 2023. [2](#)
- [76] Xuanyu Zhang, Runyi Li, Jiwen Yu, Youmin Xu, Weiqi Li, and Jian Zhang. Editguard: Versatile image watermarking for tamper localization and copyright protection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11964–11974, 2024. [3](#)
- [77] Xuanyu Zhang, Youmin Xu, Runyi Li, Jiwen Yu, Weiqi Li, Zhipei Xu, and Jian Zhang. V2a-mark: Versatile deep visual-audio watermarking for manipulation localization and copyright protection. In *Proceedings of the 32nd ACM International Conference on Multimedia (ACM MM)*, pages 9818–9827, 2024. [3](#)
- [78] Yabo Zhang, Yuxiang Wei, Dongsheng Jiang, Xiaopeng Zhang, Wangmeng Zuo, and Qi Tian. Controlvideo: Training-free controllable text-to-video generation. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023. [2](#)
- [79] Zhenghao Zhang, Junchao Liao, Menghao Li, Long Qin, and Weizhi Wang. Tora: Trajectory-oriented diffusion transformer for video generation. *arXiv preprint arXiv:2407.21705*, 2024. [3](#)
- [80] Michael Zink, Ramesh Sitaraman, and Klara Nahrstedt. Scalable 360 video stream delivery: Challenges, solutions, and opportunities. *Proceedings of the IEEE*, 107(4):639–650, 2019. [1](#)

OmniDrag: Enabling Motion Control for Omnidirectional Image-to-Video Generation

Supplementary Material

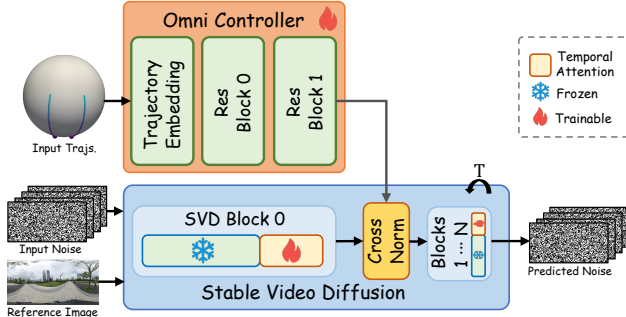


Figure 8. Illustration of our Omni Controller. Our Omni Controller only consists of a trajectory embedding module and two ResBlocks. The output control signal is integrated into the main SVD branch at its first block (SVD Block 0) by adding them to the denoising feature after applying the cross normalization.

Our main paper has outlined the core techniques of our proposed **OmniDrag** method to enable motion control for omnidirectional image-to-video generation. It has also demonstrated the efficacy of our methodological contributions through experiments. This appendix offers further details on our Omni Controller in Sec. A, more details of our Move360 dataset in Sec. B, along with additional experimental results and analyses in Sec. C, which are not included in the main paper due to space constraints. We also provide a [project page](#) to show more results, together with more sample videos from our Move360 dataset.

A. More Details of Omni Controller

An illustration of our Omni Controller is shown in Fig. 8. We further elaborate on the details of Omni Controller in the following. **Trajectory Embedding.** Recall that each sampled trajectory $\mathbf{T}_j \in \mathcal{T}$ is defined as a sequence of spatial positions $\mathbf{T}_j = \{(x_j^i, y_j^i) | i \in \{0, 1, \dots, L-1\}\}$, where (x_j^i, y_j^i) represents the position of the j -th trajectory at frame i . We follow MotionCtrl [62] to explicitly expose the moving speed of the object, as:

$$\{(0, 0), (u_{(x_1, y_1)}, v_{(x_1, y_1)}), \dots, (u_{(x_{L-1}, y_{L-1})}, v_{(x_{L-1}, y_{L-1})})\}, \quad (10)$$

where $u_{(x_i, y_i)} = x_i - x_{i-1}$, $v_{(x_i, y_i)} = y_i - y_{i-1}$, and $i \in \{0, 1, \dots, L-1\}$. The first frame and the other spatial positions in the subsequent frames that the trajectories do not pass are denoted as $(0, 0)$. As a result, $\mathbf{T}' \in \mathbb{R}^{L \times H \times W \times 2}$, where H and W are the height and width of the input ODV, respectively. A Gaussian filter is then applied to smooth the sampled trajectories, and some convolution blocks are adopted to upsample the channel dimension to 320. Finally, the condition $\mathbf{c} \in \mathbb{R}^{L \times H \times W \times 320}$.

Control Injection. We follow ControlNeXt [45] to use a lightweight architecture only composed of two ResBlocks [22]. This pruning maintains the model’s consistency while significantly reduces latency overhead and parameters. For injection of the con-

trol signals, we use cross normalization technique (Eq. 5 in the main paper), which aligns the distribution of the denoising and control features, serving as a bridge to connect the diffusion and control branches. Finally the normed control is integrated into the main branch at the first block (SVD Block 0 in Fig. 8) by addition.

B. More Details of the Move360 Dataset

We have provided the camera parameters and equipment used for data acquisition in the main paper. Here, we detail our data selection strategy and process based on three key criteria: (1) **Scene categories:** The dataset should encompass a wide range of scenes, including schools, parks, markets, landscapes, and more. Besides, the video should not have large-scale obstruction from people or buildings. (2) **Lighting conditions:** The dataset should cover various lighting conditions, such as indoor and outdoor settings, daytime (sunny and cloudy), and nighttime. (3) **Motion magnitude:** To enhance motion controllability, the videos in the dataset should exhibit relatively large motion magnitudes, avoiding examples where the background scene and object are all static.

Specifically, the original video footage has a duration of approximately 20 hours and a size of 6 TB. We first employed optical flow to filter out video segments that are nearly stationary. We then evenly split the remaining video and manually reviewed each clip, considering scene category, lighting conditions, and image quality, and excluded clips with device debugging or occlusions, resulting in 2,100 clips, each consisting of 100 frames. Finally, we utilized CoTracker [28] to calculate the average spherical motion distance of each video and filtered out the lowest 25% based on this metric. After double-checking the filtered videos, we finalized a dataset of 1,580 clips. Fig. 9 presents samples from the Move360 Dataset. We also provide a [project page](#) to showcase more video samples. We hope that the Move360 Dataset will help fill the gap in the field due to the lack of large-scale motion video datasets. We believe that Move360 will serve as a valuable resource for advancing research in omnidirectional video technologies, fostering innovation and collaboration within the community.

C. More Experiment Results

Additional visual comparisons are provided in Figs. 10 and 11. In Fig. 10, the control condition is the inversion of the case presented in Fig. 5 of the main paper; two rendered viewpoints at specific perspectives are also included. It can be observed that OmniDrag achieves stable scene-level control, whereas DragNUWA generates only artifacts. Furthermore, DragAnything distorts the ERP distribution, resulting in deformed user viewpoints. Figure 11 illustrates different drag inputs applied to a single reference image. DragAnything and DragNUWA also affect background regions not intended for control, whereas OmniDrag achieves precise object-level control. Additional results, including those in ERP format and various viewpoints, are provided on our [project page](#).

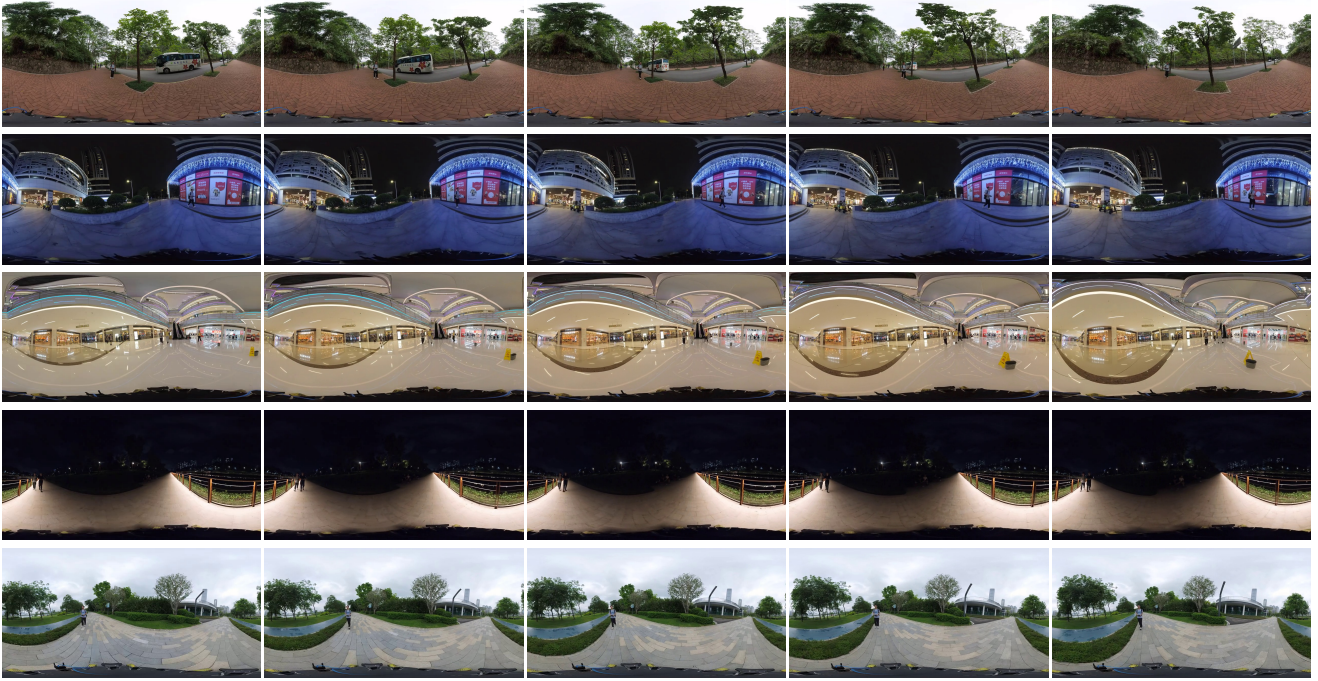


Figure 9. Some sample videos in our Move360 dataset.

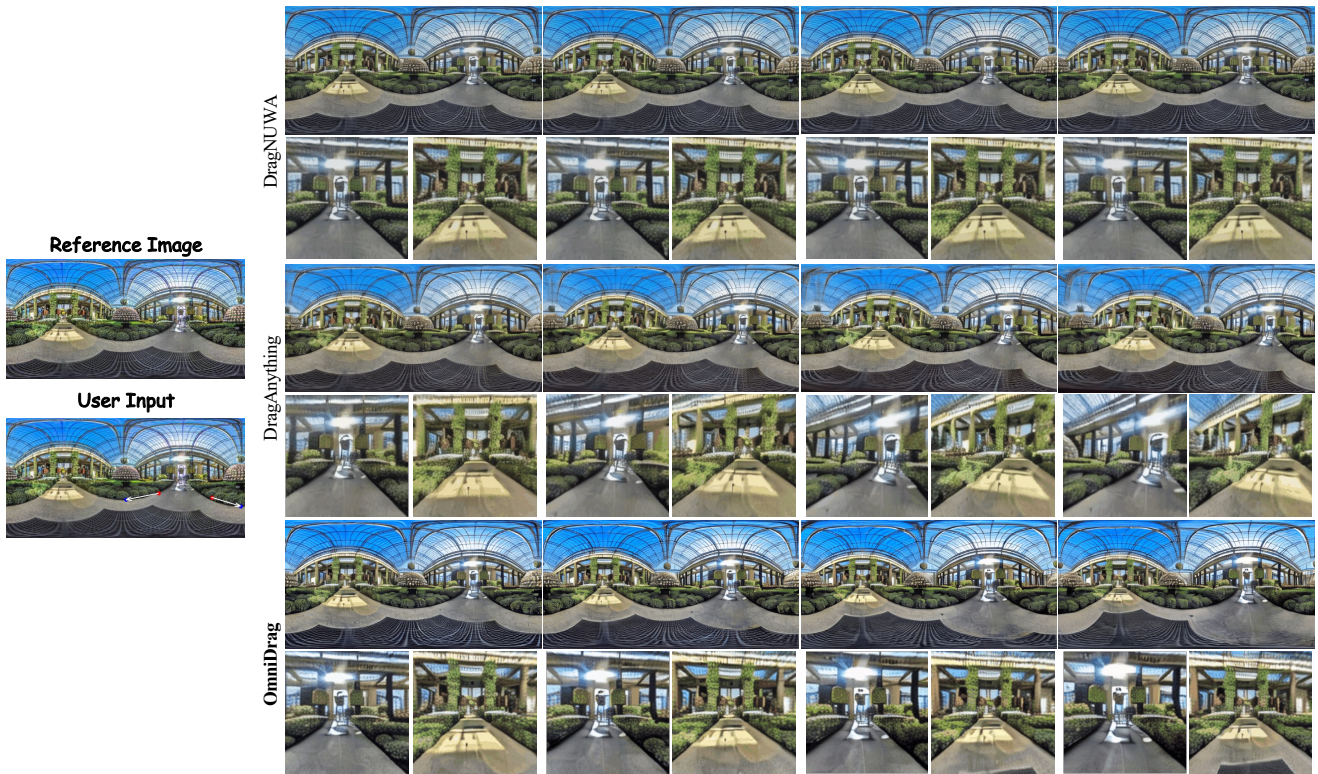


Figure 10. **Visual comparisons** of between DragNUWA [71], DragAnything [65], and our OmniDrag. For each ERP image, we show two corresponding viewports at specific perspectives.

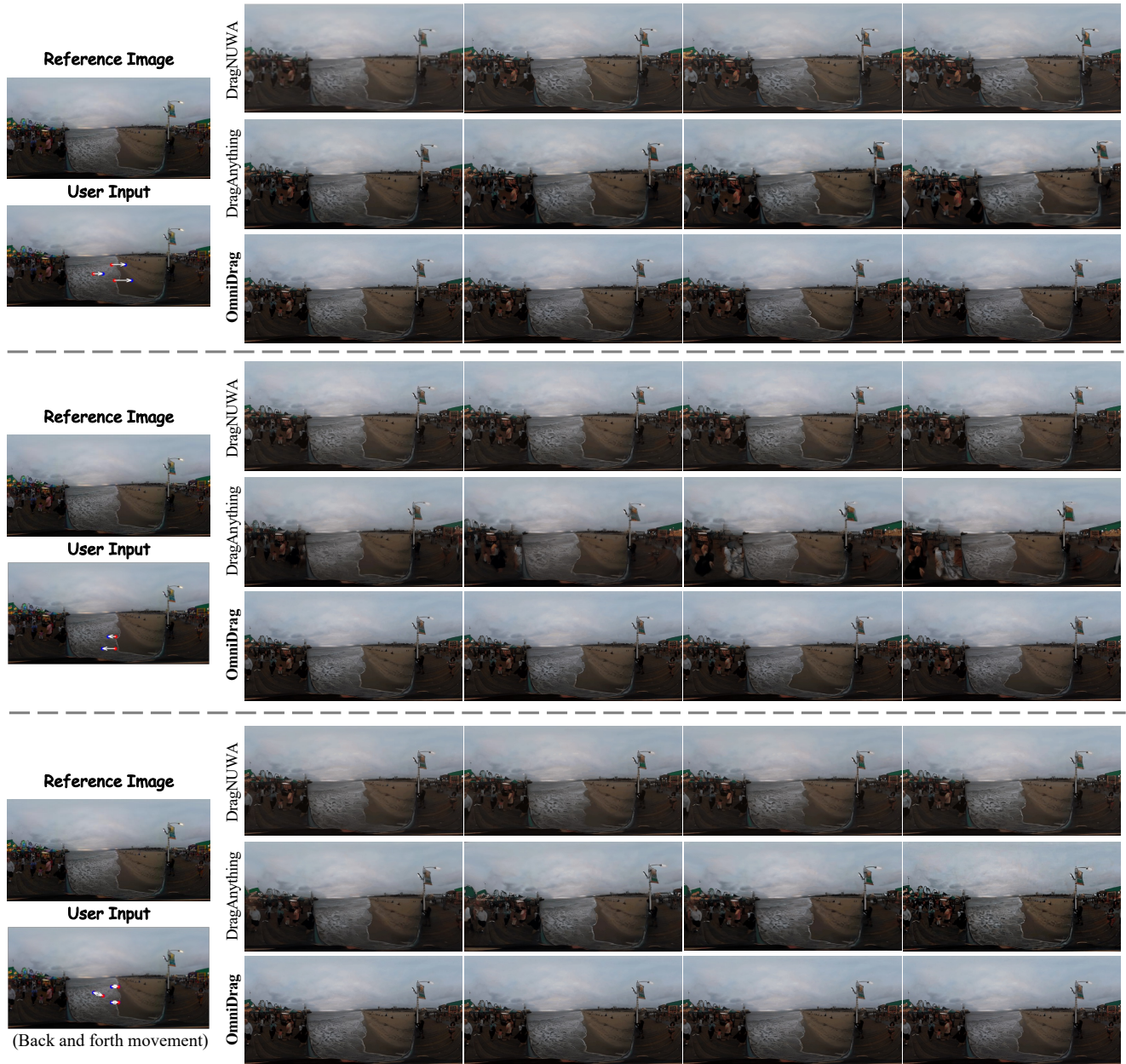


Figure 11. **Visual comparisons** of between DragNUWA [71], DragAnything [65], and our OmniDrag on the same reference image under different drag controls.