000 **REGION-WISE MOTION CONTROLLER FOR IMAGE-TO-**001 002 VIDEO GENERATION 003

Anonymous authors

Paper under double-blind review

ABSTRACT

Animating images with interactive motion control has garnered popularity for image-to-video (I2V) generation. Modern approaches typically regard the condition of Gaussian filtered point-wise trajectory as sole motion control signal. Nevertheless, such flow approximation of trajectory via Gaussian kernel severely limits the controllable capacity of fine-grained movement, and commonly fails to 015 disentangle object and camera moving. To alleviate these, we present ReMoCo, a 016 new recipe of region-wise motion controller that novelly leverages precise regionwise trajectory and motion mask to regulate fine-grained motion synthesis and identify exact target motion category (i.e., object or camera moving), respectively. Technically, ReMoCo first estimates the flow maps on each training video via a tracking model, and then samples the region-wise trajectories from multiple local regions to simulate inference scenario. Instead of approximating flow distribution via Gaussian filtering, our region-wise trajectory preserves original flow information at local area and thus manages to characterize fine-grained movement. A motion mask is simultaneously derived from the predicted flow maps to present holistic motion dynamics. To pursue natural and controllable motion generation, ReMoCo further strengthens video denoising with additional conditions of region-wise trajectory and motion mask in a feature modulation manner. More remarkably, we meticulously construct a benchmark called *ReMoCo-Bench*, which consists of 1.1K real-world user-annotated image-trajectory pairs, for the evaluation of both fine-grained and object-level motion synthesis in I2V generation. Extensive experiments conducted on WebVid-10M and ReMoCo-Bench demonstrate the effectiveness of our ReMoCo for precise motion control.

032 033 034

004

010 011

012

013

014

017

018

019

021

024

025

026

027

028

029

031

INTRODUCTION 1

In recent years, diffusion models (Ho et al., 2022a; Blattmann et al., 2023b; Singer et al., 2023; Ge et al., 2023; Brooks et al., 2024) have shown significant progress in revolutionizing text-to-037 video (T2V) generation. Although promising visual appearance can be attained by these advances, the controllable motion generation is still a grand challenge in video diffusion paradigm. There are several attempts (Esser et al., 2023; Wang et al., 2023; Chai et al., 2023) to enhance controllable 040 capacity of video synthesis with additional guidance (e.g., depth, edge or optical flow). Nevertheless, 041 it might be impractical for users to conveniently provide such signals as input conditions. Hence, the 042 focus of this paper is to capitalize on the user-friendly conditions (i.e., sparse trajectory and region 043 mask) for enabling interactively controllable image-to-video (I2V) generation: given the reference 044 image as the first frame, the motion in the synthesized video should be natural and well-aligned with the provided trajectory.

046 Pioneering practices (Yin et al., 2023; Wu et al., 2024) of controllable I2V generation usually guide 047 video denoising process with the single condition of Gaussian filtered trajectory. In the training 048 stage, the input trajectories are first sparsely sampled from the optical flow maps and then processed by Gaussian filter. The flow approximation brought by Gaussian filtering inevitably results in the inaccuracy of fine-grained motion details and limits the model capability for precise motion control. 051 Therefore, the generated fine-grained movement (e.g., the turning-head of first case in Figure 1) is unnatural. Another issue is that the single condition of trajectory commonly fails to precisely 052 identify the target motion category (i.e., camera or object moving). For instance, as depicted in Figure 1, the trajectory on the planet could be explained as two moving situations, i.e., the camera



Figure 1: An illustration of (a) fine-grained and (b) object-level motion control by using typical 072 Gaussian filtered trajectory and our region-wise motion controller (ReMoCo). The trajectories of generated videos are visualized in the last frame. 074

075 being pulled downwards with relative to static two planets (camera movement) or planet rising corre-076 sponding to static background (object movement). Solely relying on the trajectory might lead to the 077 motion misinterpretation and thus hinder exactly controllable I2V generation. To address the above 078 two issues, we shape a new paradigm of motion controller that capitalizes on region-level trajectory 079 and motion region mask to enhance video denoising for controllable motion synthesis. Specifically, we spatially sample multiple local regions in the video optical flow maps and directly employ the trajectories in the sparse regions as input trajectory condition. In this way, no flow approximation is 081 included in such region-wise trajectory, which manages to adequately reflect the local fine-grained motion details. Meanwhile, a region mask is estimated on the video optical flow maps which aims 083 to globally emphasize the motion area, thereby specifying the target motion category and alleviat-084 ing misinterpretation. To further regulate the motion synthesis in I2V generation, we predict the 085 affine parameters on the collaboration of trajectory and motion mask to modulate the video latent codes during denoising. As shown in Figure 1, our unique region-wise trajectory design and the 087 employment of motion mask complementarily achieves the better fine-grained (e.g., turning-head) 880 and object-level (e.g., planet-moving) motion generation. 089

By materializing the idea of facilitating controllable I2V generation with the proposed conditions, 090 we present a novel framework, namely ReMoCo, to execute Region-wise Motion Control. Specif-091 ically, given the input video, ReMoCo first estimates the sequence of visibility masks and optical 092 flow maps by using an off-the-shelf optical tracking model. Next, the global visibility mask is obtained through computing the intersection of all visibility masks, and further multiplied with the 094 flow map of each frame. Then, ReMoCo splits the masked flow maps into multiple local regions 095 (e.g., the region with the size of 8×8) and employs the trajectories on such sparsely-sampled re-096 gions as region-wise trajectory. Meanwhile, ReMoCo attains the motion mask on the flow maps via thresholding mechanism for representing holistic motion. Given the region-wise trajectory and corresponding motion mask, the multi-scale features are learnt by a motion encoder, and further 098 employed to predict scale and bias for video latent feature modulation. Moreover, ReMoCo finetunes all attention modules in 3D-UNet via utilizing the Low-Rank Adaptation (LoRA) technique to 100 pursue better motion-trajectory alignment. 101

102 103

104

105

071

073

In summary, we have made the following contributions:

- We introduce a new design of region-wise trajectory and motion mask as the complementary control signals in I2V diffusion models for the interactive motion control.
- A novel approach, namely Region-wise Motion Controller (ReMoCo), seamlessly integrates the 106 proposed region-wise trajectory and motion mask into 3D-UNet to guide video denoising for 107 natural and precise motion synthesis in I2V generation.

• We present *ReMoCo-Bench*, to our best knowledge, which is one of the first benchmarks with real-world user-annotated image-trajectory pairs for controllable I2V generation. Extensive experiments on WebVid-10M and ReMoCo-Bench verify the superiority of ReMoCo in terms of both video quality and motion-trajectory alignment.

111 112 113

114

108

109

110

2 RELATED WORK

115 Image-to-Video Diffusion Models. The remarkable progress achieved by text-to-video genera-116 tion (Ho et al., 2022b;a; Blattmann et al., 2023b; Khachatryan et al., 2023; Luo et al., 2023; Singer 117 et al., 2023; Ge et al., 2023; Gupta et al., 2023; Guo et al., 2024; Brooks et al., 2024) encourages 118 the development of image-to-video (I2V) diffusion models. These advances (Girdhar et al., 2024; 119 Blattmann et al., 2023a; Xing et al., 2024; Shi et al., 2024a; Zeng et al., 2024) treat static image as the 120 input condition for temporal coherent video synthesis. VideoComposer (Wang et al., 2023) is one 121 of the earlier works that integrates image condition into 3D-UNet through concatenating the clean 122 image latent with the noisy video latents. Based on this recipe, DynamiCrafter (Xing et al., 2024) 123 and SVD (Blattmann et al., 2023a) additionally inject the CLIP (Radford et al., 2021) feature of reference image into video denoising to enhance the information guidance. To achieve high-resolution 124 I2V generation, I2VGen-XL (Zhang et al., 2023b) introduces a cascading diffusion model to first 125 animate image in the low resolution and further magnifies it via video refinement. Besides, there 126 are several explorations (Chen et al., 2023b; Zeng et al., 2024) that simultaneously utilize two im-127 ages (i.e., the first and last frames) as more powerful references to elevate I2V generation. In this 128 work, we choose the pre-trained I2V diffusion model SVD (Blattmann et al., 2023a) as our base 129 architecture for motion control. 130

Controllable Video Diffusion Models. Despite high-quality video synthesis via I2V diffusion mod-131 els, the controllable motion generation still remains an under-explored problem. The early control-132 lable video diffusion techniques (Wang et al., 2023; Esser et al., 2023; Chen et al., 2023a; Zhang 133 et al., 2024) typically leverage the condition of depth, edge or optical flow, for particular motion gen-134 eration. Nevertheless, it is usually impractical for users to conveniently obtain such kinds of signals. 135 To address this issue, the studies exploring bounding box (Jain et al., 2024; Wang et al., 2024a) or 136 trajectory (Yin et al., 2023; Wu et al., 2024; Niu et al., 2024; Mou et al., 2024; Wang et al., 2024b) 137 as additional condition for motion control start to emerge. One representative of using bounding 138 box as control is PEEKABOO (Jain et al., 2024) which designs the training-free spatial-temporal 139 masked attention for visual-textual alignment in the box. In the direction of utilizing trajectory con-140 dition, pioneering advance DragNUWA (Yin et al., 2023) exploits Gaussian filtered trajectory to regulate motion synthesis via multi-scale feature fusion. Wu et al. (2024) further incorporate the 141 142 entity features of reference image into diffusion to facilitate object-level motion control. Recently, MOFA-Video (Niu et al., 2024) devises a two-stage motion control framework that first densifies 143 input trajectories via conditional motion propagation (CMP), and further regulates video denoising 144 with the estimated dense trajectories. Nevertheless, most of the existing works employ the Gaussian 145 filtered trajectory as the single condition. The Gaussian filtering will lead to flow approximation in 146 local area, which constrains the capacity for fine-grained motion modeling. Solely capitalizing on 147 trajectory could also fail to disentangle object and camera moving in I2V motion synthesis. 148

In short, our work mainly focuses on a new recipe of motion condition, i.e., the region-wise trajectory and motion mask, and the exploitation of these conditions for exact controllable I2V generation.
 The proposal of ReMoCo contributes by studying not only how to express the motion trajectory accurately, but also how to benefit natural and precise motion generation with the synergy of the region-wise trajectory and motion mask.

154

3 OUR APPROACH

155 156

In this section, we introduce our Region-wise Motion Controller (ReMoCo) for controllable I2V
generation. Figure 2 illustrates an overview of our ReMoCo. Given a video clip at training, the
newly-minted region-wise trajectory and motion mask are first extracted as the control signals. Next,
multi-scale features are learnt on the concatenation of the trajectory and mask via a motion encoder.
These features are further injected into the 3D-UNet of SVD (Blattmann et al., 2023a) to regulate
video denoising. In each feature scale of the 3D-UNet, a scale and bias are predicted through



Figure 2: An overview of our Region-wise Motion Controller (ReMoCo) for controllable imageto-video generation. During training, ReMoCo first extracts the proposed region-wise trajectory and motion mask on the input video as the control signals. The multi-scale features are then learnt on these signals by a motion encoder, and further injected into the 3D-UNet of SVD in a feature modulation manner. Meanwhile, LoRA layers are integrated into all attention modules in the transformer blocks to improve the optimization of motion-trajectory alignment. In the inference stage, the region-wise trajectory and motion mask are first derived from the user provided trajectory and brushed region, and then exploited as the guidance to calibrate I2V video generation.

convolutional layers to modulate the feature of video latent codes. Besides, all attention modules are fine-tuned by LoRA (Hu et al., 2022) to attain better alignment between the synthesized motion and input trajectory.

3.1 PRELIMINARIES: STABLE VIDEO DIFFUSION

181

182

183

184

185

186

187 188 189

190

191 192

193

201

206

213

To leverage comprehensive motion prior embedded in the pre-trained diffusion models for video generation, we exploit the advanced I2V generation model, i.e., Stable Video Diffusion (SVD) (Blattmann et al., 2023a) as the base architecture of our ReMoCo. To better understand our proposal, we first revisit the training procedure of SVD. Formally, given an input video clip $\mathbf{x}_0 = \{x_0^i\}_{i=1}^L$ with L frames, the clean video latent codes $\mathbf{z}_0 = \{z_0^i\}_{i=1}^L$ are first extracted via a variational auto-encoder (VAE). Then, the Gaussian noise n is added to \mathbf{z}_0 through forward diffusion procedure as:

$$\mathbf{z} = \mathbf{z}_0 + \mathbf{n}, \quad (\sigma, \mathbf{n}) \sim p(\sigma, \mathbf{n}),$$
 (1)

where z is the noised video latent codes and $p(\sigma, \mathbf{n}) = p(\sigma)\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. σ represents the noise level and $p(\sigma)$ is the pre-determined distribution over σ . Following the training protocol of EDM (Karras et al., 2022), SVD leverages the 3D-UNet F_{θ} (with parameters θ) to predict the clean video latent codes \hat{z}_0 with the condition of input noised latents z, noise level σ and the reference image \mathbf{c}_I :

$$\hat{\mathbf{z}}_0 = c_{\text{skip}}(\sigma) \mathbf{z} + c_{\text{out}}(\sigma) F_{\boldsymbol{\theta}}(c_{\text{in}}(\sigma) \mathbf{z}, \ \mathbf{c}_I; \ c_{\text{noise}}(\sigma)), \tag{2}$$

where $c_{\text{skip}}(\sigma)$, $c_{\text{out}}(\sigma)$, $c_{\text{in}}(\sigma)$ and $c_{\text{noise}}(\sigma)$ are pre-defined hyper-parameters determined by noise level σ . In SVD, the information of reference frame is injected into 3D-UNet along two pathways: a) the channel-wise concatenation of noised video latent codes and first frame latent code; b) the cross-attention between video latent feature and image CLIP (Radford et al., 2021) embedding of first frame. The loss function is formulated via denoising score matching (DSM) as:

$$\mathcal{L} = \mathbb{E}_{(\mathbf{z}_0, \mathbf{c}_I) \sim p_{\text{data}}(\mathbf{z}_0, \mathbf{c}_I), (\sigma, \mathbf{n}) \sim p(\sigma, \mathbf{n})} \left[\lambda_\sigma \| \hat{\mathbf{z}}_0 - \mathbf{z}_0 \|_2^2 \right],$$
(3)

where λ_{σ} is a weighting function. In the scenario of our work, besides the condition of reference first frame, we additionally exploit a new kind of region-wise trajectory and motion mask as the control signals to refine video denoising for motion control.

216 3.2 MOTION CONDITION GENERATION

218 Most existing controllable I2V approaches 219 calibrate the video denoising with the sole guidance of Gaussian filtered point-wise 220 trajectory. Nevertheless, the flow approxi-221 mation brought by Gaussian filtering may 222 result in inaccuracy of fine-grained motion details. Therefore, the ability of precise 224 motion control could be limited. Besides, 225 solely relying on the trajectory for mo-226 tion control might not exactly express tar-227 get motion category (i.e., camera or object 228 moving), leading to motion misinterpreta-229 tion in video generation. To alleviate these 230 issues, we propose to directly sample trajectories from optical flow maps in multi-231 ple local regions as the region-wise trajec-232 tory. Such trajectory preserves the original 233 flow information in local regions, and thus 234 manage to characterize fine-grained move-235 ment. In the meanwhile, a motion mask is 236



Figure 3: Motion condition generation in the training and inference stages of our ReMoCo.

plicitly identify target motion category of the generated videos.

further derived from the flow maps to ex-

242

247

248

252

256

Region-wise Trajectory. As depicted in Figure 3, given the input video clip $\mathbf{x}_0 = \{x_0^i\}_{i=1}^L$ with the size of $L \times H \times W \times 3$, we first employ a dense optical tracking model, i.e., DOT (Moing et al., 2024) to estimate optical flow maps $\mathbf{f} = \{f^i\}_{i=1}^L$ and the sequence of visibility masks $\mathbf{M} = \{M^i\}_{i=1}^L$:

$$f^{i}, M^{i} = \text{DOT}(\mathbf{x}_{0}^{1}, \mathbf{x}_{0}^{i}), \quad i = 1, 2, ..., L,$$
(4)

where $f^i \in \mathbb{R}^{H \times W \times 2}$ and $M^i \in \{0, 1\}^{H \times W}$ is the optical flow map and the visibility mask between the first and the *i*-th frame, respectively. Then, we calculate the intersection on M to attain a global visibility mask $M_g \in \{0, 1\}^{H \times W}$ that indicates the locations having visible optical flow along temporal dimension as:

$$M_g = \prod_{i=1}^{L} M^i.$$
⁽⁵⁾

Next, the masked flow maps $f_m = \{f_m^i\}_{i=1}^L$ are computed by frame-wisely multiplying the flow maps f with the global visibility mask M_g as follows:

$$\mathbf{f}_m = \{ f^i \cdot M_g \}_{i=1}^L.$$
 (6)

253 We split the masked flow maps f_m into multiple local regions and the spatial size of each region is 254 $k \times k$. The region-wise trajectories $\mathbf{T}_s \in \mathbb{R}^{L \times H \times W \times 2}$ are finally sampled from the region-split f_m 255 with a region selection mask $M_{sel} \in \{0, 1\}^{\frac{H}{k} \times \frac{W}{k}}$:

$$\Gamma_{s} = \{f_{m}^{i} \cdot Pad(M_{sel})\}_{i=1}^{L},\tag{7}$$

where M_{sel} is uniformly sampled from $\{0, 1\}$ with the mask ratio r_m , and $Pad(\cdot)$ denotes the padding function which fills the mask value into the $k \times k$ region around each position. Instead of exploiting a constant mask ratio for trajectory selection, we randomly choose r_m in a range of $[r_{min}, 1.0]$ to simulate different real-world motion masking scenarios, which benefits the robust network optimization. In this way, there is no flow approximation of the trajectories in each local region, enhancing the control ability of fine-grained motion in I2V models.

Motion Mask. In addition to the region-wise trajectory for video denoising regulation, the motion mask aims to specify the motion category and benefit the global motion correlation. Given the flow maps $f = \{f^i\}_{i=1}^L$ estimated by DOT, we first calculate the average flow magnitude $f_{avg} \in \mathbb{R}^{H \times W}$ along temporal dimension as: $f_{avg} = \frac{1}{L} \cdot \sum_{i=1}^{L} || f^i ||_2$. Then, we construct the motion mask $M_{mot} \in \{0, 1\}^{H \times W}$ from zero matrix, and set the value of the position where f_{avg} is greater than 1 as True. M_{mot} is finally repeated L times as the motion mask sequence $\mathbf{M}_{mot} \in \{0, 1\}^{L \times H \times W \times 1}$ to align the temporal length of input video for subsequent motion control learning.

270 3.3 MOTION CONTROL LEARNING271

272 With the obtained region-wise trajectory and mo-273 tion mask, we aim to control motion generation with the input signals. Inspired by the recipe of 274 feature adaptation in controllable image genera-275 tion (Zhang et al., 2023a), we propose to exploit a 276 lightweight motion encoder to estimate multi-scale 277 features on the input conditions, and utilize these 278 features to adaptively modulate video latent fea-279 ture in each corresponding scale. To further im-280 prove the alignment between input trajectory and 281 generated video, we fine-tune all attention mod-282 ules in the spatial-temporal transformer blocks of 283 3D-UNet via using LoRA (Hu et al., 2022).



Adaptive Feature Modulation

Figure 4: An illustration of adaptive feature modulation in our ReMoCo.

284 Adaptive Feature Modulation. Given the attained region-wise trajectory T_s and motion mask 285 \mathbf{M}_{mot} , we first concatenate them along channel dimension to form the input condition. As shown 286 in Figure 2, a lightweight motion encoder with a series of convolutional layers first encodes the 287 input condition into multi-scale feature maps. In each scale, the learnt feature map is employed to 288 modulate the video latent feature at the same scale in 3D-UNet. Figure 4 depicts an illustration of the 289 adaptive feature modulation by using the feature map l_s in s-th scale. Particularly, we estimate the 290 scale γ_s and bias β_s on l_s via a spatial-temporal convolutional layer. Then, the normalized feature map of the input video latent feature h_s is modulated via γ_s and β_s , and further added back to itself 291 in a skip-connection manner to form the output feature map h'_s as: 292

$$h'_{s} = GN(h_{s}) \cdot \gamma_{s} + \beta_{s} + h_{s}, \tag{8}$$

where $GN(\cdot)$ denotes the group normalization. Note that we implement zero initialization on temporal convolutional layers to initialize γ_s and β_s as zero at the beginning of training, which guarantees the stability of model optimization.

LoRA Integration. To preserve rich motion prior learnt by the pre-trained video diffusion model and elevate the effectiveness of motion control, we employ LoRA layers in all attention modules of spatial-temporal transformer blocks as demonstrated in Figure 2. Specifically, the LoRA parameters ΔW act as a residue part of the original weights W as follows:

$$\mathcal{V}' = \mathcal{W} + \Delta \mathcal{W} = \mathcal{W} + AB^T,\tag{9}$$

where \mathcal{W}' is the fused weights of attention module. A and B are trainable matrices in LoRA layers.

In the training stage, we fix all parameters in the pre-trained 3D-UNet, and only train the lightweight motion encoder and all introduced LoRA layers of the attention modules.

307 3.4 INFERENCE PIPELINE OF REMOCO

308 Our ReMoCo is a user-friendly I2V generation framework for interactive motion control. In the 309 inference stage, as shown in Figure 3, users can readily brush the motion region on the uploaded 310 reference image and draw the trajectory of moving direction as input control signals. In detail, the 311 motion mask can be directly obtained from the user provided brush mask. Given the user trajectory 312 which generally describes the movement of a single pixel, we pad the trajectory value in the $k \times k$ 313 region around the pixel position to match the training paradigm. The padded trajectory in local 314 region is exploited as the input region-wise trajectory. Finally, ReMoCo regulates video denoising 315 with the guidance of the two collaborative control signals through adaptive feature modulation. Both fine-grained and object-level motion control are facilitated by the synergy of the proposed region-316 wise trajectory and motion mask. 317

318 319

320

322

293

301

306

4 EXPERIMENTS

321 4.1 EXPERIMENTAL SETTINGS

Benchmarks. We empirically verify the merit of ReMoCo on two benchmarks, i.e., WebVid-10M (Bain et al., 2021) and our proposed ReMoCo-Bench. The WebVid-10M dataset consists

324	Input Control	DragNUWA	DragDiffusion	MOFA-Video	ReMoCo
325	14 1 1 1 1 1 1	0	0		
326					
327	i 🧑 🕐 👘				
328					
329					
330					
331					
332					
333					
334					
335					
336					
337					
338					
339					
340					
341	AN AN				
342	SIE				
343	e				
344	Figure 5: Examples of	f fine-grained mo	tion control results	on ReMoCo-Benc	h. The input control
345	signals include the refe	rence image, trai	ectory and motion n	nask. Better viewed	with Acrobat Reader
346	for the animated video	os.	-		

of 10.7M video-caption pairs. There are 5K videos in the validation set and we sample 1K videos for evaluation. For each video, trajectories sampled at a ratio of 15% along with the first frame serve as the input condition for fine-grained I2V motion generation. We follow the protocols in recent controllable I2V advance (Niu et al., 2024) and choose the Frechet Video Distance (FVD) (Unterthiner et al., 2019), Frechet Image Distance (FID) (Heusel et al., 2017), and Frame Consistency (Frame Consis.) (Qi et al., 2023) of CLIP (Radford et al., 2021) features as the evaluation metrics on WebVid-10M.

355 In practical applications, users typically prefer to control video generation through a limited number 356 of representative trajectories, often just one or two. The automatically sampled trajectories employed 357 in WebVid-10M do not adequately represent this scenario, thereby potentially compromising the va-358 lidity of the evaluation. To address this issue, we introduce **ReMoCo-Bench**, a new benchmark with reference images and user-annotated trajectories, which is tailored for the evaluation of controllable 359 I2V generation. Specifically, we meticulously collect 412 high-quality reference images from the 360 internet and construct 1.1K image-trajectory pairs via human annotation. For each reference image, 361 the annotator is required to brush the motion region and draw the trajectory of movement direction 362 according to the user intention, i.e., fine-grained local part moving or global object moving. As 363 such, the motion control performance can be evaluated from both perspectives. Due to the absence 364 of ground-truth video, FVD and FID metrics are not applicable to ReMoCo-Bench. In addition to Frame Consistency, we utilize the Mean Distance (MD) to measure the alignment between generated 366 motion and input trajectory. Two evaluation protocols are exploited for this target, i.e., MD-Img and 367 MD-Vid. MD-Img is proposed by DragDiffusion (Shi et al., 2024b) which estimates the frame-level 368 mean Euclidean distance between trajectories of input and generated frames. To further validate the video-level trajectory accuracy via MD-Vid, we replace the image correspondence detection model 369 DIFT (Tang et al., 2023) in MD-Img with the video tracking model CoTracker (Karaev et al., 2024), 370 which supplies a more precise trajectory reference. 371

Implementation Details. In ReMoCo, we employ SVD (Blattmann et al., 2023a) as our base architecture. Each training sample is 16-frames video clip and the sampling rate is 8 fps. We fix the resolution of each frame as 320×512 , which is centrally cropped from the resized video. The local region size k is set as 8 and the minimal mask ratio r_{min} is set as 0.95 determined by cross validation. We set the rank of LoRA parameters as 32. The motion encoder and LoRA layers are trained via AdamW optimizer with the base learning rate 1×10^{-5} . All experiments are conducted on 6 NVIDIA A800 GPUs with minibatch size 48.

Approach	$FVD\left(\downarrow \right)$	FID (\downarrow)	Frame Consis. (†)	Approach	$\textbf{MD-Img}\left(\downarrow\right)$	$\textbf{MD-Vid}\;(\downarrow)$	Frame Consis. (†)
DragNUWA	96.65	13.19	0.9888	DragDiffusion	14.70	13.84	0.9947
MOFA-Video	87.70	12.18	0.9895	MOFA-Video	13.94	10.50	0.9972
ReMoCo	59.88	10.40	0.9895	ReMoCo	10.56	8.34	0.9962

Table 1: Performances of fine-grained motioncontrol on WebVid-10M.

4.2 EVALUATION ON FINE-GRAINED MOTION CONTROL

We first evaluate ReMoCo on the fine-grained motion control for I2V generation. The performances on WebVid-10M and ReMoCo-Bench are summarized in Table 1 and Table 2, respectively. Our ReMoCo consistently achieves better performances on WebVid-10M across different metrics. In particular, ReMoCo attains the FVD of 59.88, outperforming the best com-

Table 3:Performances of object-level motioncontrol on ReMoCo-Bench.

Table 2: Performances of fine-grained motion

control on ReMoCo-Bench.

Approach	$\textbf{MD-Img}\left(\downarrow\right)$	$\textbf{MD-Vid}~(\downarrow)$	Frame Consis. (\uparrow)
MOFA-Video	15.56	12.04	0.9951
DragAnything	12.30	11.37	0.9917
ReMoCo	10.48	8.59	<u>0.9943</u>

petitor MOFA-Video (Niu et al., 2024) by 27.82. The better FVD indicates the better alignment of 395 data distribution between the generated and ground-truth videos. Such results basically verify the 396 superiority of exploring precise region-wise trajectory to strengthen fine-grained motion dynamic 397 learning. On ReMoCo-Bench, ReMoCo leads to performance boosts against baselines in terms of 398 MD-Img and MD-Vid, showing better alignment between the user input trajectory and synthesized 399 videos. Note that MOFA-Video exploits a two-stage controllable I2V framework that first densifies 400 the input trajectories through conditional motion propagation (CMP), and then calibrates video dif-401 fusion process using the estimated dense trajectories. In contrast, ReMoCo learns precise motion 402 patterns by directly referring region-wise trajectory via adaptive feature modulation, thus enhancing the motion-trajectory alignment, as evidenced by the better MD-Img and MD-Vid performances. 403 Besides, the CMP technique in MOFA-Video generally focuses on flow completion in the local re-404 gion surrounding the input trajectory while neglecting potential movements in other areas. Thus, 405 MOFA-Video tends to synthesize videos with less motion dynamics and obtains slightly higher 406 Frame Consistency (approximately 0.001). To substantiate this, we calculate the average flow mag-407 nitude of videos generated by MOFA-Video, which achieves 4.95. In comparison, ReMoCo attains a 408 higher value of 8.95, verifying that our model achieves greater motion variability while maintaining 409 better motion-trajectory alignment. 410

Figure 5 further showcases three I2V generation results controlled by the user input trajectory and 411 region mask on ReMoCo-Bench. Generally, the videos synthesized by our ReMoCo exhibits more 412 natural movement and better alignment with input trajectory than the baseline methods. For instance, 413 DragNUWA (Yin et al., 2023) suffers from motion misinterpretation issue which wrongly generates 414 videos with camera movement instead of object moving (e.g., the 1st and 2nd cases). The videos 415 generated by MOFA-Video (Niu et al., 2024) usually present unnatural object movement with local 416 part distortion, e.g., the nose of raccoon in the 2nd case. We speculate that such distortion is caused 417 by the lack of global region guidance in MOFA-Video, where the region mask is only employed 418 for flow masking as post-processing. Our ReMoCo, in comparison, integrates the information of 419 motion mask into 3D-UNet on the fly to facilitate the modeling of holistic motion correlation. Thus, 420 the synthesized videos by ReMoCo reflect more rational fine-grained movement.

421 422

423

380

382

384 385

386

4.3 EVALUATION ON OBJECT-LEVEL MOTION CONTROL

Next, we conduct evaluation on object-level motion control for I2V generation. Table 3 lists the performances of different approaches on ReMoCo-Bench. Overall, ReMoCo attains the best performances on the metrics of MD-Img and MD-Vid. Specifically, ReMoCo obtains 10.48 of MD-Img and 8.59 of MD-Vid, reducing the Mean Distance of the best competitor DragAnything (Wu et al., 2024) by 1.82 and 2.78, respectively. The improvements again confirm the merit of leveraging the duet of region-wise trajectory and motion mask for precise motion control. Similar performance trend on Frame Consistency can be also observed in the table.

Figure 6 shows the visual comparison of four object-level motion control results by using different approaches on ReMoCo-Bench. Compared to the baseline methods, videos generated by ReMoCo



Figure 6: Examples of object-level motion control results on ReMoCo-Bench. The input control signals include reference image, trajectory and motion mask. ReMoCo can successfully handle complicated (e.g., the round trip of sun in the 1st case) and counterintuitive (e.g., the train moving back in the 3rd case) motion-trajectory alignment. *Better viewed with Acrobat Reader.*



Figure 7: Performance comparisons of MD-Vid and Frame Consistency on ReMoCo-Bench under the settings of both fine-grained and object-level motion control by using different (a) local region size k and (b) minimal mask ratio r_{min} in ReMoCo.

474 475

476

453

454

455

456 457

458

459

460

461

462

463

4.4 ABLATION STUDY ON REMOCO

In this section, we perform ablation study to delve into the design of ReMoCo for controllable I2Vgeneration. Here, all experiments are conducted on ReMoCo-Bench for performance comparison.

479 **Local Region Size.** We first investigate the choice of local region size k for region-wise trajectory 480 design in our ReMoCo. Figure 7(a) compares the performances of MD-Vid and Frame Consistency 481 on both fine-grained and object-level motion control by using different k. The variation of Frame 482 Consistency is minor (less than 0.01) across different settings, and the MD-Vid decreases when us-483 ing larger k. When k is small (e.g., 1 or 2), the kept trajectories are less in each local region and 484 the control signals are weaken for motion control, leading to the inferior trajectory matching perfor-485 mance. Meanwhile, the improvement of MD-Vid is marginal when increasing k to 16. Specifically, 486 using large k will extend the input trajectory over a large region, which affects the fine-grained



Figure 8: Visualization of controllable I2V generation results with different local region size k in ReMoCo. *Better viewed with Acrobat Reader for the animated videos.*

motion control. Accordingly, we exploit k = 8 to extract the region-wise trajectory as the motion condition. Figure 8 further illustrates the I2V generation results with different k. As shown in this figure, the synthesized videos with k = 8 present more natural motion dynamics and more precise motion-trajectory alignment. Moreover, the unnatural fine-grained motion as shown in the case when k = 16 validates our analysis on the influence of overlarge region size.

499 Minimal Mask Ratio. To explore the effect of minimal mask 500 ratio r_{min} in trajectory selection stage, we then measure the 501 motion control performance by conducting different r_{min} in 502 Figure 7(b). Overall, Frame Consistency is not sensitive when 503 changing r_{min} on both fine-grained and object-level motion 504 control settings. Meanwhile, the performance of MD-Vid be-505 comes better with the increase of the mask ratio at the beginning. The results are expected since using small r_{min} 506 will sample more trajectories for model training, which en-507 larges the gap between training and real-world inference (i.e., 508 only using one or two trajectories). Conversely, employing a 509 large value of mask ratio (e.g., 0.99) could make it difficult to 510 optimize networks with scarce trajectory signals. Therefore, 511 we empirically set r_{min} as 0.95 to obtain the best motion-512 trajectory alignment in the generated videos. 513



Figure 9: MD-Vid (\downarrow) among different multi-scale feature injection approaches on ReMoCo-Bench.

Multi-scale Feature Injection. We also investigate different multi-scale feature injection strategies 514 in ReMoCo. Figure 9 details the MD-Vid performance comparisons among different variants of our 515 ReMoCo. **ReMoCo**^C concatenates the multi-scale features learnt by motion encoder with the video 516 latent features along channel dimension in each scale. **ReMoCo**⁺ replaces the channel-wise feature 517 concatenation in ReMoCo^C with the feature summation. In comparison, our proposal (**ReMoCo**) 518 injects the control signals into 3D-UNet via the adaptive feature modulation. Overall, ReMoCo ex-519 hibits better MD-Vid performances against other two variants. In direct feature aggregation methods 520 such as concatenation or summation, information exchange requires strict spatial-temporal align-521 ment between each other. In contrast, there is no such requirement for feature modulation, as it indirectly utilizes estimated scale and bias for feature regulation. Consequently, such feature injec-522 523 tion approach demonstrates enhanced capacity to extract relevant information from input signals, potentially leading to improved motion control performance. 524

525 526

492

493

5 CONCLUSIONS

527 528

This paper explores the motion condition formulation and the motion-trajectory alignment in diffu-529 sion models for controllable I2V generation. In particular, we study the problem from the viewpoint 530 of integrating accurate motion control signals into video denoising to regulate motion generation. 531 To materialize our idea, we have devised ReMoCo, which leverages the region-wise trajectory and 532 motion mask as the condition to calibrate video generation in a feature modulation manner. The 533 region-wise trajectory preserves the original optical flow information in each local region, char-534 acterizing the fine-grained motion details. The motion mask derived from the optical flow maps presents holistic motion and aims to identify exact target motion category. The collaboration of 536 two signals regulates video denoising for natural motion synthesis with precise motion-trajectory 537 alignment. Moreover, we have carefully construct a new benchmark, i.e., ReMoCo-Bench, with 1.1K real-world user-annotated image-trajectory pairs for the evaluation of both fine-grained and 538 object-level motion control. Extensive experiments on WebVid-10M and ReMoCo-Bench validate the superiority of our proposal over state-of-the-art approaches.

540 **ETHICS STATEMENT** 6 541

542 The primary of this paper is to introduce a controllable image-to-video diffusion model for general 543 individuals to animate particular object motion. It is important to note that the visual contents of 544 our generated video are aligned with those of input reference image. Even though there could be some ethical concerns for input image, employing an additional content safety checker to filter 546 reference image can potentially resolve this issue. We uphold the highest ethical standards in the construction of our ReMoCo-Bench, and believe that all the contents in the dataset are appropriate 547 548 while respecting relevant privacy rights in data collection procedure.

549 550

551 552

553

554

555 556

558

559

560

561

562

563 564

565

566

567

568

571

576

577 578

579

580

581

585

586

587 588

589

590

7 **REPRODUCIBILITY STATEMENT**

We have introduced the model construction and implementation details in the paper. To enhance the reproducibility of our approach, we attach the core code of ReMoCo in the supplementary material with detailed explanations.

- References
- Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in Time: A Joint Video and Image Encoder for End-to-End Retrieval. In ICCV, 2021.
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, Wesam Manassra, Prafulla Dhariwal, Casey Chu, and Yunxin Jiao. Improving Image Generation with Better Captions, 2023.
- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets. arXiv preprint arXiv:2311.15127, 2023a.
- 569 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 570 Models. In CVPR, 2023b.
- 572 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe 573 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video 574 Generation Models as World Simulators. 2024. 575
 - Wenhao Chai, Xun Guo, Gaoang Wang, and Yan Lu. StableVideo: Text-driven Consistency-aware Diffusion Video Editing. In ICCV, 2023.
 - 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 Diffusion Models with Motion Prior and Reward Feedback Learning. arXiv preprint arXiv:2305.13840, 2023a.
- 582 Xinyuan Chen, Yaohui Wang, Lingjun Zhang, Shaobin Zhuang, Xin Ma, Jiashuo Yu, Yali Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. SEINE: Short-to-Long Video Diffusion Model for Generative 583 584 Transition and Prediction. In ICCV, 2023b.
 - Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and Content-Guided Video Synthesis with Diffusion Models. In ICCV, 2023.
 - Songwei Ge, Seungjun Nah, Guilin Liu, Tyler Poon, Andrew Tao, Bryan Catanzaro, David Jacobs, Jia-Bin Huang, Ming-Yu Liu, and Yogesh Balaji. Preserve Your Own Correlation: A Noise Prior for Video Diffusion Models. In ICCV, 2023.
- Rohit Girdhar, Mannat Singh, Andrew Brown, Quentin Duval, Samaneh Azadi, Sai Saketh Ramb-592 hatla, Akbar Shah, Xi Yin, Devi Parikh, and Ishan Misra. Emu Video: Factorizing Text-to-Video Generation by Explicit Image Conditioning. In ECCV, 2024.

626

- Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. AnimateDiff:
 Animate Your Personalized Text-to-Image Diffusion Models without Specific Tuning. In *ICLR*, 2024.
- Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang, and José Lezama. Photorealistic Video Generation with Diffusion Models. *arXiv preprint arXiv:2312.06662*, 2023.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
 GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In
 NeuIPS, 2017.
- Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P.
 Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen Video: High
 Definition Video Generation with Diffusion Models. In *CVPR*, 2022a.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J.
 Fleet. Video Diffusion Models. In *NeurIPS*, 2022b.
- Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
 and Weizhu Chen. LoRA: Low-Rank Adaptation of Large Language Models. In *ICLR*, 2022.
- Yash Jain, Anshul Nasery, Vibhav Vineet, and Harkirat Behl. PEEKABOO: Interactive Video Generation via Masked-Diffusion. In *CVPR*, 2024.
- Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian
 Rupprecht. CoTracker: It is Better to Track Together. In *ECCV*, 2024.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the Design Space of Diffusion-Based Generative Models. In *NeurIPS*, 2022.
- Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang
 Wang, Shant Navasardyan, and Humphrey Shi. Text2Video-Zero: Text-to-Image Diffusion Models are Zero-Shot Video Generators. In *ICCV*, 2023.
- ⁶²³ Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao,
 ⁶²⁴ Jingren Zhou, and Tieniu Tan. VideoFusion: Decomposed Diffusion Models for High-Quality
 ⁶²⁵ Video Generation. In *CVPR*, 2023.
- Guillaume Le Moing, Jean Ponce, and Cordelia Schmid. Dense Optical Tracking: Connecting the Dots. In *CVPR*, 2024.
- Chong Mou, Mingdeng Cao, Xintao Wang, Zhaoyang Zhang, Ying Shan, and Jian Zhang. ReVideo:
 Remake a Video with Motion and Content Control. In *NeurIPS*, 2024.
- Muyao Niu, Xiaodong Cun, Xintao Wang, Yong Zhang, Ying Shan, and Yinqiang Zheng. MOFA Video: Controllable Image Animation via Generative Motion Field Adaptions in Frozen Image to-Video Diffusion Model. In *ECCV*, 2024.
- Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng
 Chen. FateZero: Fusing Attentions for Zero-Shot Text-Based Video Editing. In *ICCV*, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision. In *ICML*, 2021.
- Kiaoyu Shi, Zhaoyang Huang, Fu-Yun Wang, Weikang Bian, Dasong Li, Yi Zhang, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Motion-I2V: Consistent and Controllable Image-to-Video Generation with Explicit Motion Modeling. In ACM SIGGRAPH, 2024a.
- Yujun Shi, Chuhui Xue, Jun Hao Liew, Jiachun Pan, Hanshu Yan, Wenqing Zhang, Vincent Y. F.
 Tan, and Song Bai. DragDiffusion: Harnessing Diffusion Models for Interactive Point-Based Image Editing. In *CVPR*, 2024b.

- ⁶⁴⁸ Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-Video: Text-to-Video Generation without Text-Video Data. In *ICLR*, 2023.
- Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. Emergent
 Correspondence from Image Diffusion. In *NeurIPS*, 2023.
- Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and
 Sylvain Gelly. FVD: A new Metric for Video Generation. In *ICLR DeepGenStruct Workshop*, 2019.
- 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*, 2024a.
- Kiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen,
 Deli Zhao, and Jingren Zhou. VideoComposer: Compositional Video Synthesis with Motion
 Controllability. In *NeurIPS*, 2023.
- Zhouxia Wang, Ziyang Yuan, Xintao Wang, Tianshui Chen, Menghan Xia, Ping Luo, and Ying
 Shan. MotionCtrl: A Unified and Flexible Motion Controller for Video Generation. In ACM
 SIGGRAPH, 2024b.
- 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 *ECCV*, 2024.
- 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. In *ECCV*, 2024.
- Shengming Yin, Chenfei Wu, Jian Liang, Jie Shi, Houqiang Li, Gong Ming, and Nan Duan. Drag NUWA: Fine-Grained Control in Video Generation by Integrating Text, Image, and Trajectory.
 arXiv preprint arXiv:2308.08089, 2023.
- Yan Zeng, Guoqiang Wei, Jiani Zheng, Jiaxin Zou, Yang Wei, Yuchen Zhang, and Hang Li. Make
 Pixels Dance: High-Dynamic Video Generation. In *CVPR*, 2024.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding Conditional Control to Text-to-Image
 Diffusion Models. In *ICCV*, 2023a.
- 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*, 2023b.
 - Yabo Zhang, Yuxiang Wei, Dongsheng Jiang, Xiaopeng Zhang, Wangmeng Zuo, and Qi Tian. ControlVideo: Training-Free Controllable Text-to-Video Generation. In *ICLR*, 2024.
- 689 690

687

688

680

657

- 691
- 692
- 693 694

- 696
- 697
- 698
- 699
- 700

Evaluation Items	Fine-grained Motion Control			Object-level Motion Control		
	DragDiffusion	MOFA-Video	ReMoCo	MOFA-Video	DragAnything	ReMoCo
Motion Quality (^)	3.12	21.88	75.00	12.50	18.75	68.75
Temporal Coherence (↑)	6.25	40.63	53.12	25.00	15.63	59.37
Trajectory Alignment (†)	9.37	18.75	71.88	15.62	21.88	62.50

Table 4: Human evaluation of user preference ratios (%) over both fine-grained and object-level motion control across different approaches on ReMoCo-Bench.

711 A

712

705 706

708 709 710

A APPENDIX: MORE DETAILS OF REMOCO-BENCH

The proposed ReMoCo-Bench consists of 412 high-quality reference images and corresponding 713 1.1K user-annotated trajectories. We collect the reference images with different visual contents, in-714 cluding animal, human, vehicle, etc. There are 72 images sampled from the public DragBench (Shi 715 et al., 2024b) and we further extend it with 340 additional images. Specifically, all the self-collected 716 images about human are automatically generated by DALL·E3 (Betker et al., 2023) to avoid the po-717 tential legal concerns. The remaining self-collected images are real photos which are first crawled 718 on the Pexels platform and then filtered according to the visual quality. For each reference image, 719 the annotator is required to brush the motion region and draw the movement trajectory according to 720 user intention (i.e., fine-grained local part moving or global object moving). During trajectory anno-721 tation, all annotators are encouraged to ensure the trajectory diversity, including some complicated 722 trajectories. Finally, the benchmark is annotated with 460 image-trajectory pairs for fine-grained motion control evaluation, and 680 image-trajectory pairs for object-level motion control evalu-723 ation, respectively. Figure 10 and Figure 11 further illustrate several visual examples (reference 724 image, trajectory and motion mask) from ReMoCo-Bench for the two evaluations. 725

726 727

B APPENDIX: HUMAN EVALUATION

728

In addition to the evaluation over automatic metrics, we also conduct human evaluation to investigate user preferences from three perspectives (i.e., motion quality, temporal coherence and trajectory alignment) across different controllable I2V approaches. In particular, we randomly sample 200 generated videos from both fine-grained and object-level motion control for evaluation. Through the Amazon MTurk platform, we invite 32 evaluators, and ask each evaluator to choose the best one from the generated videos by all models given the same inputs.

Table 4 shows the user preference ratios across different models on ReMoCo-Bench. Overall, our
 ReMoCo clearly outperforms all baselines in terms of the three criteria on both fine-grained and
 object-level motion control. The results demonstrate the advantage of leveraging complementary
 region-wise trajectory and motion mask to benefit video synthesis with natural motion, desirable
 temporal coherence and precise motion-trajectory alignment.

740 741

742

C APPENDIX: OFFLINE PROJECT PAGE

We build an offline project page for our ReMoCo in the "ReMoCo.github.io" folder, and package it into supplementary material. Please click the file of "index.html" in the folder with the Chrome or Firefox browser for more vivid video presentation.

746 747 748

D APPENDIX: CODE RELEASE

Moreover, we package the core code of our ReMoCo in the "ReMoCo-Code" folder of supplementary material. Please refer to the example source code and README in the folder for more details.

751 752

749

- 753
- 754
- 755



Figure 10: Visual examples from ReMoCo-Bench for fine-grained motion control evaluation. Each reference image is annotated with trajectory and motion mask.



Figure 11: Visual examples from ReMoCo-Bench for object-level motion control evaluation. Each reference image is annotated with trajectory and motion mask.