

---

# Track the Noise, Move the World: 3D-Grounded Motion-Consistent Noise for Controllable Video Generation

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Modern image-and-text-to-video diffusion models can synthesize highly realistic  
2 videos by iteratively denoising an initial Gaussian noise tensor conditioned on  
3 reference image and text inputs. However, existing approaches still lack precise  
4 and unified controllability over both object motion and camera motion within a  
5 single generation process. We present **UniCaMo**, a unified framework that en-  
6 ables simultaneous control of object trajectories and camera viewpoints by directly  
7 constructing the input noise of the diffusion model. Specifically, UniCaMo builds  
8 a shared *3D-grounded motion-consistent noise space* across latent video frames.  
9 Sparse 3D point tracks are used to warp the Gaussian noise of the reference frame  
10 along desired object trajectories, while a virtual spherical noise representation  
11 provides globally consistent noise values for newly revealed scene regions under  
12 camera motion. By combining local track-guided noise warping with global sphere-  
13 based noise sampling, UniCaMo maintains geometric and temporal consistency  
14 under both object movement and viewpoint changes. Because UniCaMo modifi-  
15 es only the input noise, it requires no auxiliary adapters, control branches, or  
16 architectural changes to the underlying video diffusion model. With lightweight  
17 LoRA fine-tuning on large pretrained video diffusion models, including Wan 2.1  
18 (14B), UniCaMo achieves state-of-the-art results in both video quality and motion  
19 controllability on standard controllable video generation benchmarks.

## 20 1 Introduction

21 Video is a highly expressive visual medium, and precise control over generated content enables  
22 significant practical value. A scene is defined not only by its objects, but by their motion and  
23 the camera’s movement—e.g., a car on a street looks entirely different when tracked, panned, or  
24 viewed from above. Beyond storytelling, controllable video generation is crucial for robotics and  
25 embodied AI [41, 25], enabling controllable data synthesis during training and predictive imagination  
26 at inference time. Motivated by this, we ask: *can we build a video generation framework that*  
27 *simultaneously and precisely controls object and camera motion without external adapters?*

28 Controllable video generation has gained increasing attention alongside advances in large-scale video  
29 diffusion models [40, 32, 20, 47]. Recent works achieve strong control via point tracks, optical flow,  
30 camera trajectories, or auxiliary networks [13, 16, 4, 8, 10, 31]. However, most methods treat object  
31 and camera motion separately—either animating objects from a fixed viewpoint or controlling camera  
32 motion with unconstrained scene dynamics. When both are specified, this separation often leads to  
33 motion ambiguity and entanglement artifacts. Since object and camera motion are inherently coupled,  
34 independent control limits expressiveness. In this work, we explicitly tackle joint object-and-camera  
35 motion control, enabling a key capability for modern image- and text-to-video frameworks.

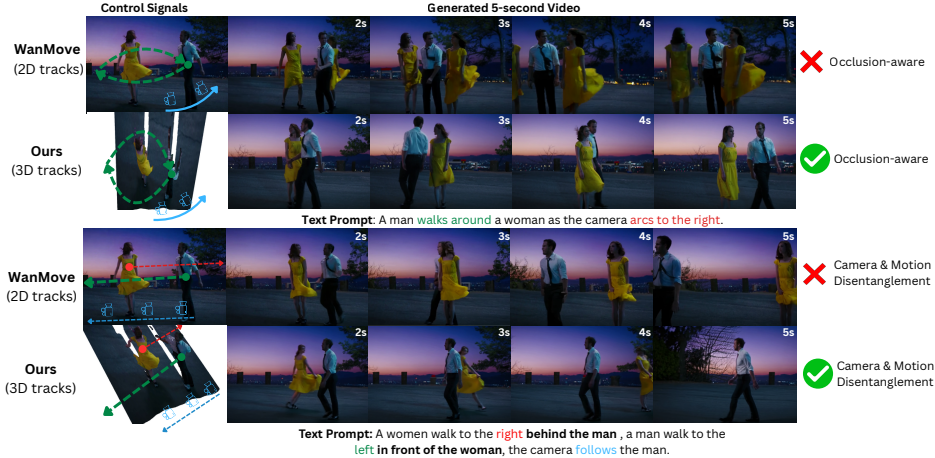


Figure 1: UniCaMo jointly controls object + camera motion with 3D tracks and spherical noise, producing coherent, occlusion-aware videos. WanMove [13] relies on 2D guidance, often causing motion/camera ambiguity and implausible occlusions.

36 Having established the need for joint motion control, we examine how existing methods specify  
 37 motion. Prior work [13, 31, 46] typically relies on sparse 2D point tracks, allowing users to drag  
 38 image-plane points to indicate object or camera motion. While intuitive, this approach has two key  
 39 limitations. First, 2D trajectories provide an incomplete and ambiguous description of underlying 3D  
 40 dynamics: the same image-space motion can arise from multiple combinations of object motion and  
 41 camera movement, and many motions—such as one person moving around another (Fig. 1)—are  
 42 difficult to express without explicit depth or viewpoint information. Second, motion constraints are  
 43 usually injected late in the generation pipeline, when visual content has already emerged. This leaves  
 44 limited flexibility for resolving motion ambiguity in a geometrically coherent way while maintaining  
 45 realism and temporal consistency. Together, these issues often lead to imprecise, entangled, or  
 46 geometrically inconsistent motion control.

47 This observation motivates a different perspective: instead of injecting ambiguous motion guidance  
 48 near the output, can we impose motion constraints directly at the input? Modern image-and-text-  
 49 to-video diffusion models generate videos by iteratively denoising an initial Gaussian noise tensor  
 50 conditioned on a reference image [40]. This suggests a reformulation of controllable video generation  
 51 as the problem of constructing an input noise tensor that already encodes the desired scene dynamics.  
 52 Rather than controlling the generation process itself, we control its input. Our approach is thus related  
 53 to recent noise-warping methods [10, 8], which achieve strong object motion control by warping  
 54 object-associated noise across frames. However, extending these methods to joint object-and-camera  
 55 motion is fundamentally challenging. Under a static camera, object motion creates small, localized  
 56 voids that diffusion models can often infer and inpaint using surrounding context. In contrast, camera  
 57 motion exposes large unseen regions with no corresponding noise from the reference view. As the  
 58 viewpoint diverges, these missing regions grow, making noise warping from a single reference-frame  
 59 fundamentally insufficient for preserving temporal consistency and geometric coherence.

60 To address these challenges, we propose a motion-consistent noise representation that remains  
 61 geometrically coherent under both object and camera motion. Instead of defining noise only on  
 62 the reference image plane, we lift it into a shared spherical 3D representation parameterized by  
 63 Gaussian noise, providing consistent noise values for all viewing directions. As the camera moves,  
 64 each frame samples from this shared sphere based on its viewpoint, enabling temporally coherent  
 65 synthesis of newly revealed regions. Within this space, user-specified point tracks locally warp  
 66 noise along object trajectories, ensuring consistent object motion under changing viewpoints. This  
 67 3D formulation allows users to specify both object and camera motion directly in 3D, avoiding  
 68 image-space ambiguity. Combining sphere-based noise sampling with track-guided warping, our  
 69 method—UniCaMo—constructs a unified 3D-grounded motion-consistent noise space that enables  
 70 precise joint control of object and camera motion while modifying only the input noise. UniCaMo  
 71 requires no architectural changes or auxiliary adapters, supports lightweight LoRA fine-tuning of

72 pretrained models such as Wan 2.1 (14B) [40], adds only 0.2s per sample, and achieves state-of-the-art  
73 video quality and controllability [13].

74 In summary, our main contributions are: **(1)** we propose UniCaMo, a novel motion-consistent noise  
75 representation that combines 3D point tracks with sphere-projected Gaussian noise to jointly encode  
76 object and camera motion in a unified 3D space; **(2)** UniCaMo requires no auxiliary adapters or  
77 trainable control modules, introducing only negligible warping overhead (approximately 0.2s per  
78 sample) while preserving the runtime efficiency of the base diffusion architecture; and **(3)** extensive  
79 experiments on standard controllable video generation benchmarks, including MoveBench [13],  
80 demonstrate state-of-the-art controllability and high visual quality, particularly in scenarios involving  
81 complex coordinated object and camera motion.

## 82 2 Related Work

83 Diffusion-based video generation has advanced from early 3D U-Nets [6, 7, 11, 15, 19, 18, 34]  
84 to large-scale latent diffusion models with stronger spatiotemporal modeling. Transformer-based  
85 foundation video models further improve long-range coherence and fidelity via spatiotemporal  
86 attention [26, 1, 20, 27, 40, 47, 47, 27]. We build on efficient pretrained models such as Wan [40],  
87 and achieve controllable generation through structured latent-noise modeling and lightweight fine-  
88 tuning—without architectural changes.

89 **Motion and camera control for video generation.** Motion-controlled I2V methods differ by motion  
90 cues and how they are injected: masks/boxes/keypoints enables coarse motion cues, while flow or  
91 trajectories provide finer guidance [24, 43, 31, 53, 8, 28, 14, 49, 51]. Among these methods, many rely  
92 on ControlNet-style adapters with higher cost, whereas Wan-Move [13] edits conditioning features  
93 without architecture changes. Camera control in video diffusion is typically achieved by conditioning  
94 on 6DoF poses/trajectories, injected via dedicated modules [16, 46, 3]. Related work enforces cross-  
95 view consistency (e.g., CVD) or performs camera-only V2V “reshoots,” but these methods generally  
96 preserve object dynamics and do not support simultaneous object-motion editing [29, 4, 35, 2].

97 **Warped noise for controllable video generation.** A promising direction is to control diffusion by  
98 using structured initialization noise. Noise-warping methods create temporally correlated noise by  
99 warping Gaussian noise with motion fields (e.g., optical flow), injecting temporal structure while  
100 preserving per-frame spatial Gaussianity. Warped-noise priors such as HIWYN [10] highlight the  
101 importance of this constraint, though large deformations can introduce artifacts and overhead. Go-  
102 with-the-Flow [8] makes warped noise efficient and shows motion control can be added by changing  
103 only the noise/data pipeline, treating the diffusion model as a black box. Recent analysis further finds  
104 that training with warped noise promotes useful equivariances and can enable few-step sampling [33].

105 **Training-free methods.** Training-free methods control video diffusion at inference by editing  
106 noise/latents, attention, or denoising schedules [21, 48, 37]. Examples include CamTrol and DVS  
107 for camera/viewpoint control, and Time-to-Move for motion via region-wise schedules, but such  
108 methods often struggle with long-horizon, multi-object, or precise control [21, 48, 37].

109 **Unified control.** Existing methods with joint camera and object control often relies on 3D adapters  
110 or 3D primitives, increasing runtime cost [43, 31, 14, 36, 13, 30, 52]. Particularly, trajectory-based  
111 methods aim for unified control by using background tracks for camera motion and foreground tracks  
112 for object motion. Edit-by-Track leverages 3D point tracks to resolve depth/occlusion in joint edits,  
113 while other systems use geometry-aware adapters or richer 3D states, typically adding conditioning  
114 modules and runtime [30, 52]. In contrast, we target joint camera+object control via structured noise:  
115 we add no inference-time adapters [44, 46] and avoid dense flow warping [8] by combining sparse 3D  
116 tracks with sphere-projected noise to handle disocclusions while preserving per-frame Gaussianity.

## 117 3 UniCamMo: a Method for Joint Camera and Object Motion Control

118 UniCaMo is built upon an image-and-text-to-video framework such as Wan [40]. It accepts the  
119 following input: (i) a source image  $I$  and text prompt  $y$ , (ii) a sparse set of target 3D point tracks per  
120 video frame  $\mathcal{T} = \{\mathbf{x}_i^t \in \mathbb{R}^3\}_{i=1, N, t=1, T}$ , and (iii) target per-frame camera poses  $\mathcal{C} = \{\mathbf{P}^t\}_{t=1, T}$ ,  
121 with  $N$  being the number of point tracks per frame and  $T$  the number of video frames. Note that  
122 during inference, both  $\mathcal{T}$  and  $\mathcal{C}$  are defined by the user using a graphical interface.

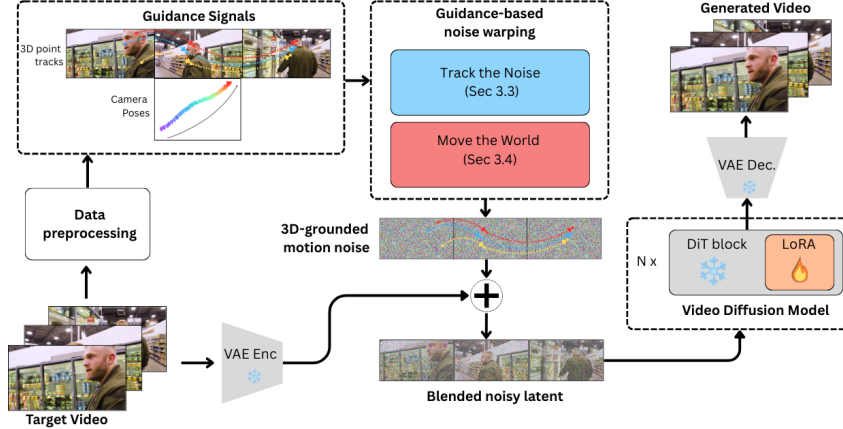


Figure 2: Our framework trains a video diffusion model by warping initialization noise using 3D point tracks and camera poses (via Track-the-Noise and Move-the-World) to create motion-consistent 3D-grounded latents, enabling independent or joint control of object and camera motion at inference.

123 During training, we construct the input based on original video data as follows. For each video, we  
 124 extract the first frame and use it as the source image  $I$ , and apply video captioning to extract text  
 125 prompt  $y$ . We then apply point tracks, camera pose and depth estimation on the video to obtain point  
 126 tracks  $\mathcal{T}$  and pose  $\mathcal{C}$ . During inference, we assume that the target 3D point tracks  $\mathcal{T}$  and target camera  
 127 poses  $\mathcal{C}$  are specified by the user through a graphical user interface. Our goal is to generate a video  
 128 with  $T$  frames initialized from the content of the source image  $I$  with motion dynamics following the  
 129 target object and camera motion encoded by point tracks  $\mathcal{T}$  and camera pose  $\mathcal{C}$ .

130 Fig. 2 summarizes our training pipeline. We first perform data pre-processing to obtain depth  
 131 maps, camera poses, and 3D point tracks of the target video using off-the-shelf 3D models such as  
 132 ViPE [22] or TAPI3D [50]. We then apply a noise warping technique to obtain our 3D-grounded  
 133 motion-consistent noise through two steps: (1) Track-the-Noise (Sec. 3.1), which creates partial  
 134 motion-consistent noise that captures the motion inferred from 3D point tracks, and (2) Move-the-  
 135 World (Sec. 3.2), which samples and propagates Gaussian noise on a virtual 3D sphere to produce  
 136 view-consistent noise along a target camera trajectory. The combined representation results in  
 137 3D-grounded motion-consistent noise that can be used to initialize the latents of a video diffusion  
 138 model [40] to generate the final output video.

### 139 3.1 Track the Noise: Noise Initialization using 3D Point Tracks

140 **Motivation.** Conventional video diffusion models [40, 47] initialize generation with a single Gaussian  
 141 noise latent, which effectively behaves as frame-wise uncorrelated noise, lacking any temporal  
 142 structure. Yet real videos exhibit strong motion correlations, recoverable via dense optical flow [8] or  
 143 sparse point tracks [30]. Leveraging this as an inductive bias, we partially warp the initial noise using  
 144 the estimated motion (Fig. 3(a)), producing a motion-consistent noise prior that guides the model  
 145 toward coherent video dynamics.

146 **Partial Noise Warping.** Let  $\mathcal{T} = \{\mathbf{x}_i^t \in \mathbb{R}^3\}$  denote the set of world-space 3D point tracks across  
 147 video frames. Compared to dense optical flow [38], sparse 3D tracks provide explicit depth cues  
 148 and are robust to occlusion and viewpoint changes. We map the point tracks to latent feature map  
 149 locations by first projecting them onto the image plane using the target camera poses  $\{\mathbf{P}^t\}$ , and  
 150 then rescaling the projected coordinates to the latent feature map resolution. We use the scale ratio  
 151  $W_{lat}/W$  and  $H_{lat}/H$ , where  $W_{lat}, H_{lat}, W, H$  are latent and image dimensions, respectively (for  
 152 Wan 2.1 14B model [40], this scale ratio is 4). Through 3D-to-2D projection, we obtain the 2D pixel  
 153 trajectories be  $\{\mathbf{u}_i^t \in \mathbb{R}^2\}$ , corresponding to the 3D point tracks.

154 By initializing the latent representation of the first frame as i.i.d. Gaussian noise, we now propagate  
 155 the noise to subsequent frames following the pixel trajectories  $\{\mathbf{u}_i^t\}$  constructed from 3D point tracks.  
 156 We simply transfer (i.e., copy and write) noise within a square patch with dimension  $2R - 1$  centered  
 157 at each pixel in the trajectories from the first frame to the target frames. By default, we set  $R = 2$ ,

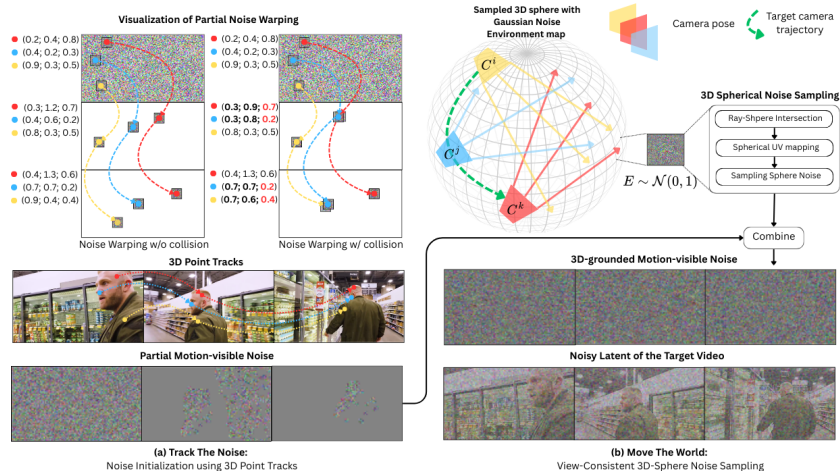


Figure 3: Noise representation in our two key steps, Track the Noise (a) and Move the World (b). (a) We track 3D points (see example triplets of X, Y, Z values) and resolve potential collisions in their trajectories using a depth test, favoring contributions from points closer to the camera. This step results in initial partial motion-visible noise. (b) We create a virtual 3D sphere with Gaussian-textured noise and sample view-consistent noise according to the target camera poses. By combining the outputs of (a) and (b), we obtain 3D-grounded, motion-consistent noise that replaces the randomly sampled noise used in standard video diffusion models.

158 resulting in patch of size  $3 \times 3$ . As the 3D point tracks encode motion in the video, this results in  
 159 a video latent representation that exhibits motion-consistent noise, forming a strong guidance for  
 160 generating motion dynamics during video synthesis.

161 **Depth-aware Collision Tracking.** When multiple point trajectories overlap and project to the same  
 162 latent pixel, naive composition leads to ambiguous assignments and may introduce inconsistent  
 163 motion cues. We therefore resolve collisions using a depth-aware z-buffer rule induced by the  
 164 target tracks  $\mathcal{T}$  and camera poses  $\mathcal{C}$ . Among all candidates mapped to a latent pixel, we select the  
 165 contribution with the smallest camera depth (i.e., the point closest to the camera) and discard the  
 166 rest. This choice enforces a physically plausible occlusion ordering and leverages the explicit depth  
 167 information provided by 3D tracks to disambiguate visibility during noise propagation.

168 In contrast, prior motion-control pipelines often adopt simple but less effective heuristics for handling  
 169 overlaps, such as selecting one contributor without geometric reasoning (e.g., by arbitrary ordering)  
 170 [13], or aggregating multiple contributors through averaging (sometimes followed by re-normalization  
 171 to restore Gaussian statistics) [8]. By explicitly enforcing depth-consistent compositing, our con-  
 172 struction yields geometrically coherent structured noise signals, improving the stability of the noise  
 173 representation under occlusions and viewpoint changes.

### 174 3.2 Move-the-World: View-Consistent 3D-Sphere Noise Sampling

175 **Motivation.** As the Track-the-Noise step transports first-frame noise along point-track trajectories,  
 176 some regions remain unconstrained—especially areas that become newly visible as camera motion  
 177 reveals previously unseen parts of the scene. Filling these gaps with random Gaussian noise provides  
 178 no motion cues, while copying first-frame noise introduces incorrect parallax. To address this, we  
 179 introduce Move-the-World, which represents the scene as a 3D bounding sphere coated with Gaussian  
 180 noise. The camera moves inside this sphere, and projecting the sphere’s noise into each view produces  
 181 view-consistent latents that preserve spatial Gaussianity (Fig. 3(b)).

182 **3D Sphere Parameterization.** Constructing a bounded 3D sphere requires definition on the cen-  
 183 ter position and radius. In particular, we set the sphere center to the camera center of the first  
 184 frame and choose a radius  $r$  such that the sphere encloses the entire camera trajectory across  
 185 frames while satisfying a conservative depth bound  $d_{\max}^0$ . The radius  $r$  is define as follows:  
 186  $r = \max(\tau \cdot \max_t \|\mathbf{C}^t - \mathbf{C}^1\|_2, d_{\max}^0)$  where  $\mathbf{C}^t \in \mathbb{R}^3$  denote the camera center at frame  $t$ ,

187  $d_{\max}^0$  is the furthest depth value of the input image and  $\tau = 1.2$ . We represent the sphere surface by a  
188 latitude-longitude parameterization, which can be captured by a randomly sampled Gaussian noisy  
189 texture map  $\mathbf{E} \in \mathbb{R}^{C \times H_{\text{env}} \times W_{\text{env}}}$ , where  $H_{\text{env}}$  and  $W_{\text{env}}$  are proportional to the video resolution scaled  
190 by a factor  $k$  ( $k = 4$ ), and  $C$  matches the latent channel count of the video diffusion model [40, 47].

191 **Ray-sphere Intersection & Noise Sampling.** To obtain 3D-consistent noisy latents for each camera  
192 view, we cast rays from the camera through its view frustum and compute their intersections with  
193 a bounding sphere using a closed-form ray-sphere intersection formula. For every ray that hits the  
194 sphere surface, we determine the latitude and longitude of the intersection point. These spherical  
195 coordinates are then used to perform nearest-neighbor grid sampling on the previously defined noisy  
196 texture map  $E$ . This procedure produces noise latents that remain spatially Gaussian within each  
197 frame while staying consistent across different viewpoints. Additional details on the ray-sphere  
198 intersection computation can be found in the appendix section.

199 **Noise Composition.** To combine the motion-consistent noises produced independently by  
200 Track-the-Noise and Move-the-World, we apply a composition step that aggregates their outputs into  
201 a single final noise latent. We begin by initializing the final latent with the partial motion-consistent  
202 noise obtained from the Track-the-Noise stage. Then, for regions where this step provides no cover-  
203 age, we fill in the missing areas by copying the corresponding values from the Move-the-World output.  
204 This strategy yields a complete latent that is both 3D-grounded and motion-consistent, making it  
205 well-suited for the subsequent video generation with the diffusion model.

### 206 3.3 Training and Inference

207 **Training data.** We curate 400K training clips from DynPose100K++, WSDG-1M, and 4DNeX [22,  
208 12]. For each clip, we precompute and cache offline the signals needed for noise construction: camera  
209 intrinsics/extrinsics and depth from ViPE [22], and sparse 3D point tracks from TAPIP3D [50]. These  
210 signals are used only to build the motion-consistent noise initializer (not as additional conditioning  
211 inputs). To standardize captions across sources, we re-caption all clips with Qwen2.5-VL [5],  
212 following the Wan-Move prompt-extension protocol to improve prompt consistency [13].

213 **Training details.** We fine-tune a pretrained Wan I2V backbone [40] using flow matching to predict  
214 the velocity field that transports noise samples to the data distribution. We start from Wan2.1-I2V-  
215 A14B weights and apply LoRA only to the DiT denoiser, freezing the VAE and all encoders. We  
216 use LoRA rank 64 with  $\alpha=1$ , and train for 10K steps in bfloat16 with gradient checkpointing on 32  
217 NVIDIA H100s (batch size 1/GPU) for 3 days. To improve robustness to varying guidance density,  
218 we randomly sample 512–1024 points from the extracted 3D tracks during training.

219 **Inference.** Given a single input image, we estimate monocular depth to obtain a 3D point cloud. The  
220 user selects target objects via a mask (manual or SAM [9]), then specifies 3D point tracks for these  
221 objects over time and optionally a camera trajectory (e.g., via a simple 3D GUI). These signals are  
222 passed to our Guidance-based Noise Warping module to construct 3D-grounded, motion-consistent  
223 noise. We generate videos with classifier-free guidance ( $\gamma = 5.0$ ). UniCaMo adds no extra parameters  
224 and incurs only 0.2s warping overhead per sample, preserving the base model’s runtime. Additional  
225 qualitative results and demos are provided in the supplement/project page.

## 226 4 Experiments

### 227 4.1 Experimental Setup

228 **Baselines.** We evaluate UniCaMo on controllable I2V generation, focusing on joint camera-object  
229 motion. We compare against adapter-based controllers (MagicMotion, Tora, ImageConductor,  
230 LeviTor), feature-warping Wan-Move, and flow-based noise-warping Go-with-the-Flow [31, 51, 44,  
231 42, 13, 8]. All methods use the same inputs and we report averages over full benchmark splits.

232 **Benchmarks and Metrics.** We report results on the standard controllable video benchmark, namely  
233 MoveBench [13]. We evaluate along two axes: (i) **visual quality** using FID [17], FVD [39] and  
234 reconstruction-style metrics PSNR, SSIM [45], and (ii) **control fidelity** using motion accuracy  
235 metrics such as end-point error (EPE) for trajectories. MoveBench [13] provides only 2D tracks  
236 and first-frame masks, so we recover the required 3D tracks and camera poses using off-the-shelf  
237 TAPIP3D and ViPE [50, 22]. Using the masks, we sample  $M$  object points (proportional to mask  
238 size) and  $N - M$  background points to match our training track count. To provide an apple-to-apple

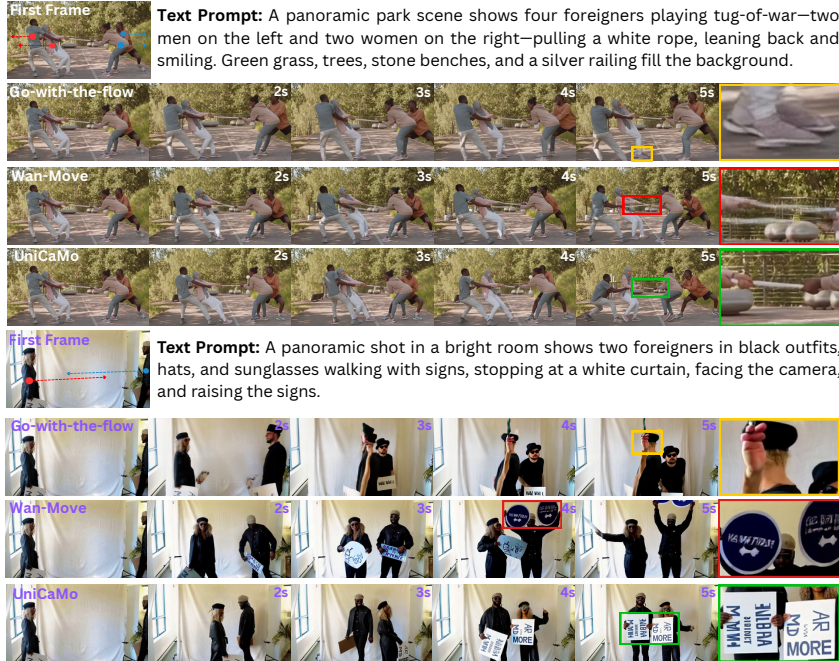


Figure 4: Qualitative results on MoveBench between UniCaMo and recent approaches [13, 31, 51]. Our method consistently outperforms other baselines and manages to follow the target motion and camera. Visual artifacts are shown in red boxes.

239 comparison under matched track sources, we additionally evaluate a *Wan-Move + TAPI3D tracks*  
 240 variant, in which we project the same 3D tracks used by UniCaMo onto the 2D image plane and feed  
 241 them as input to Wan-Move.

## 242 4.2 Qualitative and Quantitative Evaluation

243 Figure 4 shows representative examples in MoveBench under joint camera and object motion control.  
 244 Overall, UniCaMo generates temporally coherent videos that better preserve scene structure while  
 245 following the specified object trajectories and camera motion. Compared to another baselines, our  
 246 results exhibit fewer visuals artifacts and manages to follow the conditioning camera and object  
 247 motions. Go-with-the-flow [8] produces visible artifacts in small detail like human’s heah or feet,  
 248 Wan-Move [13] while exhibits less artifacts in small detail, failed to generate motion that follow  
 249 user’s intention. Please see the included videos in the supplementary material for more scene details.

250 Figure 5 demonstrates the generalization of our method on unseen in-the-wild samples with complex  
 251 geometry and trajectory. In the first example, we expect the white van to move leftward out of the  
 252 video frame while the camera moves to the right. In this case, Wan-Move [13] failed to generate  
 253 plausible results. Our 3D-grounded motion-consistent noise representation represents such motion  
 254 dynamics based on user-provided point tracks and camera trajectories, resulting in high-quality  
 255 motion content. In the second example, similarly, Wan-Move can pan the camera to the right but does  
 256 not generate the cars in the correct moving direction. Contrastively, our method can generate realistic  
 257 flows of the cars in the crowded street.

258 **Single-object motion control.** Table 1 (left block) summarizes quantitative results on the single-  
 259 object split of MoveBench. UniCaMo achieves the best overall performance across both video quality  
 260 and control fidelity metrics. Our model yields substantially improved quality and control compared  
 261 to prior motion-control methods, with lower FID/FVD and higher PSNR/SSIM while reducing  
 262 trajectory error (EPE). In particular, on MoveBench we obtain FID 10.36 and FVD 78.6, improving  
 263 over Wan-Move (FID 12.2, FVD 83.5) and trajectory-based baselines in the same evaluation protocol.  
 264 We also observe consistent gains in SSIM and EPE (e.g., SSIM 0.73 and EPE 2.3 for our model  
 265 versus SSIM 0.64 and EPE 2.6 for Wan-Move). Notably, Go-with-the-Flow [8] did not produce

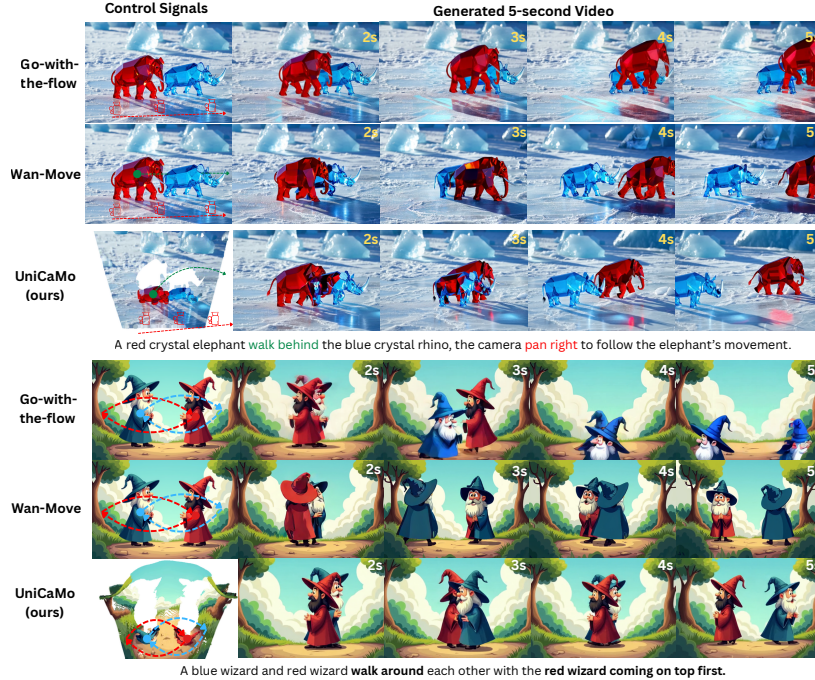


Figure 5: Qualitative results on complex in-the-wild samples between our proposed UniCaMo and Wan-Move [13], Go-with-the-flow [8]

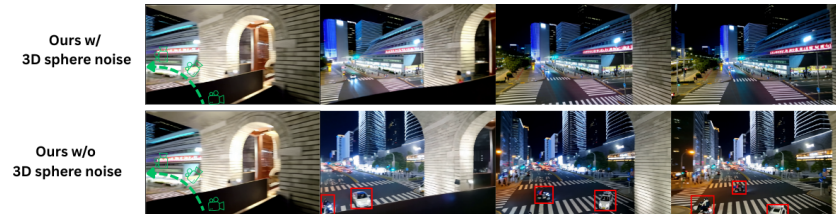


Figure 6: Ablation results on the 3D sphere noise. Without 3D-consistent sphere noise sampling, the model tends to hallucinate random objects (red boxes) with implausible motion trajectories.

Table 1: Evaluation results on MoveBench. We report quantitative comparisons on both the single-object and multi-object motion splits. “-” indicates the baseline does not provide results for that split. UniCaMo achieves best overall performance across. Best in **bold**; second-best underlined.

Methods	<i>MoveBench (Single-object motion)</i>					<i>MoveBench (Multi-object motion)</i>				
	FID ↓	FVD ↓	PSNR ↑	SSIM ↑	EPE ↓	FID ↓	FVD ↓	PSNR ↑	SSIM ↑	EPE ↓
ImageConductor [44]	34.51	424.1	13.4	0.49	15.66	77.5	764.5	13.9	0.51	9.80
LeviTor [42]	18.12	98.8	15.6	0.54	3.40	-	-	-	-	-
MagicMotion [31]	17.53	96.7	14.9	0.56	3.20	-	-	-	-	-
Go-with-the-Flow [8]	12.49	216.9	15.8	0.62	3.04	35.7	399.4	<u>16.9</u>	0.60	3.71
Tora [51]	22.57	100.4	15.7	0.55	3.30	53.2	350.0	14.5	0.54	3.50
Wan-Move [13]	12.23	<u>83.5</u>	17.8	0.64	2.60	28.8	<u>226.3</u>	16.7	0.62	2.20
Wan-Move + TAPIP3D tracks [13, 50]	<u>12.11</u>	87.1	<u>18.1</u>	<u>0.67</u>	<u>2.54</u>	<u>27.6</u>	294.8	16.5	<u>0.65</u>	<u>2.14</u>
UniCaMo (Ours)	<b>10.36</b>	<b>78.6</b>	<b>18.4</b>	<b>0.73</b>	<b>2.30</b>	<b>26.1</b>	<b>215.2</b>	<b>18.15</b>	<b>0.72</b>	<b>1.92</b>

266 meaningful videos under MoveBench’s joint-control protocol despite repeated attempts, indicating  
 267 that flow-based warped-noise priors do not transfer to this evaluation.

268 **Multi-object motion control.** As MoveBench includes 192 cases with annotated multi-object motion,  
 269 we further evaluate UniCaMo against baselines [44, 51, 8, 13] on this challenging setting, as presented  
 270 in Table 1 (right block). Our method achieves the lowest FVD and reduced EPE compared to other  
 271 methods, highlighting its precise adherence to motion constraints in more complex scenarios.

272 **Supplementary videos.** Please check our provided HTML page to see the videos and comparisons.

Table 2: Ablation studies on UniCaMo on the MoveBench [13] single-object motion evaluation set. We report VBench-I2V [23] scores together with PSNR, SSIM, and End Point Error (EPE). (a) compares our method with and without the sphere-projected noise of Move-the-World; (b) sweeps the degradation rate  $\theta$  used to blend constructed motion-noise with pure Gaussian noise.

Setting	PSNR $\uparrow$	SSIM $\uparrow$	EPE $\downarrow$	Subject Consistency $\uparrow$	Background Consistency $\uparrow$	Aesthetic Quality $\uparrow$	Imaging Quality $\uparrow$	Overall Consistency $\uparrow$	Temporal Flickering $\uparrow$	Motion Smoothness $\uparrow$
(a) Effect of sphere-projected noise (Move-the-World)										
Ours w/ sphere	<b>19.41</b>	<b>0.79</b>	<b>2.21</b>	<b>95.29</b>	<b>96.38</b>	<b>53.41</b>	69.58	21.65	<b>98.31</b>	<b>99.12</b>
Ours w/o sphere	19.10	<b>0.78</b>	2.42	94.79	96.11	52.94	<b>70.05</b>	<b>21.72</b>	98.11	98.92
(b) Effect of degradation rate $\theta$										
$\theta = 0.0$	<b>19.46</b>	<b>0.79</b>	<b>2.19</b>	<b>95.40</b>	96.35	53.39	69.51	<b>21.75</b>	<b>98.33</b>	<b>99.13</b>
$\theta = 0.2$	19.41	<b>0.79</b>	2.21	95.29	<b>96.38</b>	53.41	69.58	21.65	98.31	99.12
$\theta = 0.5$	19.02	0.78	2.49	95.08	96.24	<b>53.52</b>	<b>69.59</b>	21.67	98.12	98.92

Table 3: Analysis on number of 3D tracks (left) and patch dimension in Track-the-Noise step (right).

Points	PSNR $\uparrow$	SSIM $\uparrow$	EPE $\downarrow$	CLIP $\uparrow$	Patch	PSNR $\uparrow$	SSIM $\uparrow$	EPE $\downarrow$	CLIP $\uparrow$
$N = 1024$	<b>19.28</b>	<b>0.78</b>	<b>2.26</b>	<b>0.9385</b>	$R = 1.0$	16.95	0.71	5.34	0.9201
$N = 512$	19.04	0.78	2.38	0.9373	$R = 2.0$	19.28	0.78	<b>2.26</b>	<b>0.9385</b>
$N = 256$	18.44	0.77	2.78	0.9347	$R = 3.0$	<b>19.40</b>	<b>0.79</b>	2.28	0.9379

### 273 4.3 Ablation Studies

274 **3D-sphere noise.** We validate 3D-sphere noise sampling in Figure 6 by ablating it and using  
 275 random noise instead. Without sphere-induced cross-frame noise consistency, the model hallucinates  
 276 implausible objects and inconsistent layouts, whereas enabling sphere sampling (Move-the-World)  
 277 yields realistic, coherent motion. Quantitatively on MoveBench single-object (Table 2(a)), sphere  
 278 noise improves PSNR (19.28 vs 19.10) and EPE (2.26 vs 2.42) and consistently boosts VBench-I2V  
 279 metrics (consistency, aesthetics, reduced flicker, smoother motion), confirming it is a key component.

280 **Degradation rate.** We observe that our 3D-grounded motion noise preserves spatial Gaussian  
 281 characteristics but lacks temporal Gaussianity. To mitigate this issue, we blend the constructed motion-  
 282 noise with pure Gaussian noise through alpha blending, using a randomly sampled degradation rate  
 283  $\theta \in [0, 0.5]$ . In Table 2 (b) and Figure 12, a degradation rate of  $\theta = 0.2$  provides the best balance  
 284 between visual quality (PSNR, SSIM, Aesthetics Quality, Imaging Quality) and motion fidelity (EPE,  
 285 Temporal Flickering, Motion Smoothness). We adopt  $\theta = 0.2$  as the default in our experiments.

286 **Point tracks.** As we use TAPI3D [50] to extract 1024 points for each video during preprocessing,  
 287 and then vary the total of point tracks during training, we provide an experiment to validate the  
 288 robustness of our method when the number of point tracks are varied. Table 3 (left) and Fig. 13  
 289 provides the performance of our model on different number of 3D point tracks. The results show  
 290 that using 512 and 1024 point tracks yield similar PSNR, SSIM, and CLIP values, while using 1024  
 291 leads to more favorable EPEs. Using 256 points yields less competitive results, but the performance  
 292 remains close to the 512 and 1024 case.

293 **Patch size.** We analyze the influence of the patch dimension ( $2R - 1$ ) when propagating noise using  
 294 point tracks as described in the Track-the-Noise step (Section 3.1). The results are in Table 3 (right)  
 295 and Fig. 14. It can be seen that using larger patch size ( $R = 2$ ) yields the best EPE and CLIP while  
 296 setting  $R = 3$  yields the best PSNR and SSIM. Given the minor improvement in PSNR and SSIM,  
 297 we opt for  $R = 2$  as our default setting.

## 298 5 Conclusion

299 We presented UniCaMo, an image-to-video generation framework that unifies camera motion and  
 300 object motion through a 3D-grounded motion-consistent Gaussian noise representation. By combin-  
 301 ing noise warping based on 3D point tracks and with 3D-sphere noise sampling, UniCaMo produces  
 302 motion-consistent noise that guides video diffusion without modifying the underlying model archi-  
 303 tecture. Our results on MoveBench suggest that structured noise initialization is an effective and  
 304 efficient way for controllable video generation, enabling independent and joint control over motion.