

Time-to-Move: Training-Free Motion Controlled Video Generation via Dual-Clock Denoising

Anonymous CVPR submission
Paper ID 0000

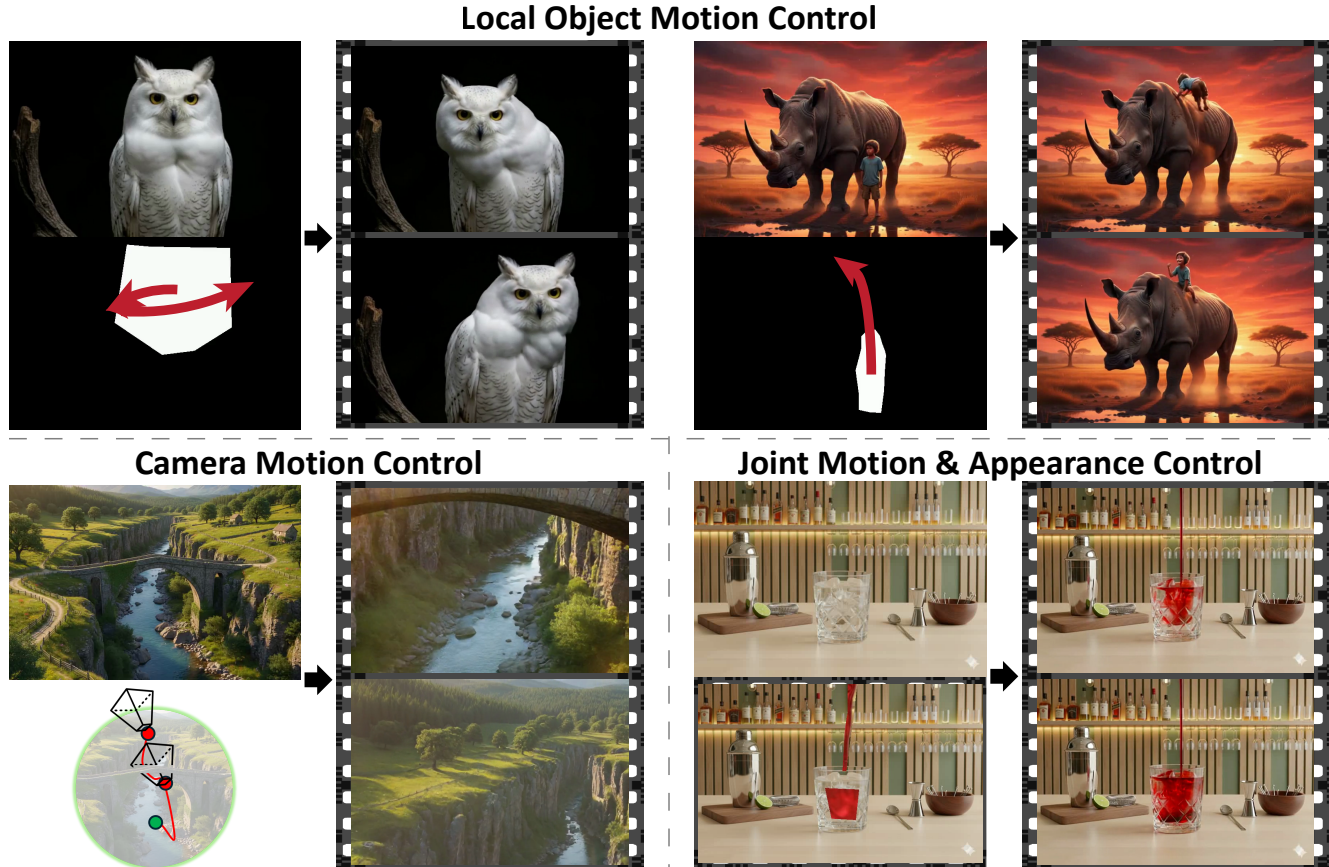


Figure 1. Qualitative results of Time-to-Move on various tasks.

Abstract

001 Diffusion-based video generation can create realistic
 002 videos, yet existing image- and text-based conditioning fails
 003 to offer precise motion control. Prior methods for motion-
 004 conditioned synthesis typically require model-specific fine-
 005 tuning, which is computationally expensive and restrictive.
 006 We introduce Time-to-Move (TTM), a training-free, plug-
 007 and-play framework for motion- and appearance-controlled
 008 video generation with image-to-video (I2V) diffusion mod-
 009 els. Our key insight is to use crude reference anima-
 010 tions obtained through user-friendly manipulations such as
 011 cut-and-drag or depth-based reprojection. Motivated by
 012 SDEdit’s use of coarse layout cues for image editing, we
 013 treat the crude animations as coarse motion cues and adapt

the mechanism to the video domain. We preserve appear-
 014 ance with image conditioning and introduce dual-clock de-
 015 noising, a region-dependent strategy that enforces strong
 016 alignment in motion-specified regions while allowing flex-
 017 ibility elsewhere, balancing fidelity to user intent with natu-
 018 ral dynamics. This lightweight modification of the sampling
 019 process incurs no additional training or runtime cost and
 020 is compatible with any backbone. Extensive experiments
 021 on object and camera motion benchmarks show that TTM
 022 matches or exceeds existing training-based baselines in re-
 023 alism and motion control. Beyond this, TTM introduces
 024 a unique capability: precise appearance control through
 025 pixel-level conditioning, exceeding the limits of text-only
 026 prompting.
 027

028 **1. Introduction**

029 Recent image-to-video (I2V) diffusion models produce im-
 030 pressive visual quality, but prompt-only control remains too
 031 coarse for production-facing tasks such as shot design, ob-
 032 ject blocking, and appearance-directed storyboarding. In
 033 cinematic workflows, users need to specify *what* moves,
 034 *where* it moves, and *what stays consistent* as motion un-
 035 folds. Existing controls often fail to deliver this precision
 036 and frequently introduce drift, static backgrounds, or iden-
 037 tity loss.

038 Most strong motion-control approaches inject trajec-
 039 tories, flow, or tracks through dedicated modules and substan-
 040 tial fine-tuning [1, 2, 18]. While effective, these methods are
 041 expensive to train, tightly coupled to a specific backbone,
 042 and difficult to transfer across rapidly evolving I2V mod-
 043 els. For workshop settings centered on movie-grade proto-
 044 typing, this limits practical use by creators who need fast
 045 adaptation rather than model retraining.

046 A practical assistant for AI filmmaking should support
 047 local edits (moving an actor or prop), global edits (camera
 048 trajectory changes), and style-aware evolution (controlled
 049 appearance changes) under one unified inference process. It
 050 should also preserve temporal plausibility in unconstrained
 051 regions so that secondary motion remains coherent with the
 052 primary edit. This requirement is particularly important in
 053 previsualization and iterative directing, where artists repeat-
 054 edly adjust motion intent while expecting stable visual iden-
 055 tity.

056 We introduce Time-to-Move (TTM), a training-free,
 057 architecture-agnostic inference method for controllable I2V
 058 generation. Our pipeline starts from a crude user-defined
 059 reference animation (e.g., cut-and-drag or depth reprojec-
 060 tion), then uses an SDEdit-inspired initialization to in-
 061 ject motion while first-frame image conditioning preserves
 062 identity and scene appearance. The key novelty is *region-*
 063 *dependent dual-clock denoising*: masked regions receive
 064 stronger motion enforcement, while unmasked regions de-
 065 noise with weaker constraints to maintain natural scene dy-
 066 namics. This design directly addresses a core failure mode
 067 of single-clock sampling, which cannot simultaneously pre-
 068 serve strict local control and global realism.

069 Related methods span trajectory-conditioned training-
 070 based controllers [5, 14–17, 19], geometry/noise-based
 071 camera control [1, 11], and training-free attention manip-
 072 ulations in T2V [4, 7, 10]. Our method differs by en-
 073 forcing region-specific control directly at inference in pre-
 074 trained I2V models, combining motion fidelity with appear-
 075 ance consistency and supporting edits relevant to cinematic
 076 previsualization (Fig. 3).

077 In summary, our contributions are:

- 078 • **Training-free controllable motion from crude refer-**
 079 **ences.** We convert simple user guidance into realistic mo-
 080 tion without architecture-specific retraining.

- **Dual-clock denoising for regional control.** We intro-
 duce a region-dependent schedule that balances strict lo-
 cal adherence with globally coherent dynamics.
- **Unified motion and appearance editing.** By condition-
 ing on full reference frames, TTM supports joint motion
 and visual-attribute control beyond trajectory-only inter-
 faces.

088 **2. Method**

089 Given an input image and a user motion intent, our goal is to
 090 generate a realistic video that follows the intended motion
 091 while preserving identity and scene coherence. TTM com-
 092 bines three components: (1) a crude reference animation
 093 as explicit motion guidance, (2) image-conditioned sam-
 094 pling for appearance anchoring, and (3) dual-clock denois-
 095 ing for region-dependent control strength. Conceptually,
 096 the method treats user edits as structural guidance rather
 097 than as final appearance targets: the reference only specifies
 098 coarse motion and placement, while the diffusion model re-
 099 constructs realistic details under image conditioning. This
 100 makes the approach robust to low-quality warps while re-
 101 maining faithful to intended trajectories.

102 **Problem Formulation.**

103 Inputs are: an image $I \in \mathbb{R}^{3 \times H \times W}$, a coarse warped refer-
 104 ence video $V^w \in \mathbb{R}^{F \times 3 \times H \times W}$, and a binary mask video
 105 $M \in \{0, 1\}^{F \times H \times W}$ marking regions with explicit user
 106 guidance. The objective is a video x_0 that follows the pre-
 107 scribed motion in masked regions while remaining globally
 108 plausible and faithful to the source appearance.

109 **2.1. Motion Signal**

110 We represent user intent as a coarse animation V^w . In
 111 object-motion control, this is produced by cut-and-drag
 112 style warping of selected regions across time; in camera-
 113 motion control, it is produced via depth-based reprojection.
 114 Although crude and artifact-prone, these references encode
 115 the desired trajectory and spatial-temporal layout, making
 116 them effective control carriers during diffusion sampling.
 117 Forward warping naturally introduces tearing and disocclu-
 118 sions, but these artifacts are acceptable because the refer-
 119 ence serves as a control signal rather than a photoreal-
 120 istic target. The key requirement is geometric alignment
 121 with user intent. In practice, this representation is flexi-
 122 ble: users can define motion through trajectories, geometric
 123 transforms, or depth-driven camera changes without retrain-
 124 ing the generator.

125 **2.2. SDEdit Adaptation for Motion Injection**

126 Following SDEdit intuition [8], we initialize denoising from
 127 a noisy version of the reference so coarse dynamics are in-
 128 jected early [12]. To avoid identity drift, we use an *image-*
 129 *conditioned* I2V model anchored to the clean first frame I ,
 130

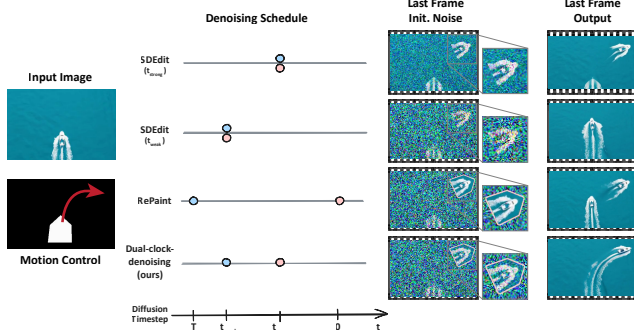


Figure 2. **Single-clock vs. dual-clock denoising.** A single timestep cannot jointly enforce strict local motion and natural global dynamics. Dual-clock denoising applies stronger guidance in masked regions and weaker guidance elsewhere, reducing drift and artifacting.

allowing motion transfer from V^w while preserving appearance. This design separates responsibilities: the noisy reference carries motion layout, and the image condition preserves content identity. Empirically, this is crucial for cinematic scenes where subject consistency and background continuity must be maintained under strong local edits.

2.3. Region-Dependent Dual-Clock Denoising

Single-clock noising creates a hard trade-off: low noise overconstrains the scene (e.g., frozen backgrounds), while high noise weakens motion adherence. We therefore apply different effective schedules across regions. Let $t_{\text{weak}} > t_{\text{strong}}$. We initialize from V^w noised to t_{weak} . For steps $t_{\text{strong}} \leq t < t_{\text{weak}}$, we override masked regions using the reference at the matching noise level, while unmasked regions follow standard denoising. After $t = t_{\text{strong}}$, we stop overriding and denoise jointly to the final sample. Intuitively, masked regions need strong guidance because they encode explicit user intent, whereas unmasked regions should retain flexibility to accommodate physically plausible secondary motion. A single global timestep cannot satisfy both goals simultaneously. Dual-clock denoising enforces this asymmetry directly at inference time, avoiding model changes or additional training.

Let x_t be the current sample and $\hat{x}_{t-1}(x_t, t, I)$ the denoiser prediction. The dual-clock update is:

$$x_{t-1} \leftarrow (1 - M) \odot \hat{x}_{t-1}(x_t, t, I) + M \odot x_{t-1}^w,$$

where x_{t-1}^w is the warped reference noised to step $t - 1$. This yields strong local trajectory control while preserving coherent unconstrained motion and appearance. Operationally, this update can be viewed as a controlled blending between model prediction and reference guidance. Early steps prioritize control where requested; later steps prioritize global harmonization. This staged behavior is what allows TTM to handle both local object manipulation and camera-style motion control within the same framework.

Efficiency and Applicability.

The method modifies only sampling logic, introduces no extra training, and transfers across different I2V backbones.

3. Experiments

We evaluate TTM in three settings: object motion control (Sec. 3.1), camera motion control (Sec. 3.2), and joint motion-appearance control (Sec. 3.3). Unless noted, we follow official baseline settings and each backbone’s native protocol (SVD: 16-frame setup; CogVideoX: 49-frame setup), with a fixed dual-clock schedule per backbone. Extensive ablations (omitted for space) confirm that region-dependent dual-clock denoising is necessary to avoid temporal artifacts while preserving motion adherence. We emphasize evidence most relevant to production use: controllability under large edits, consistency of subject identity, and temporal plausibility in unconstrained regions.

3.1. Object Motion Control

We benchmark on MC-Bench [18] using user-provided masks and trajectories. We compare against DragAnything [15], MotionPro [18], SG-I2V [9], and GWTF [1]. Evaluation uses trajectory adherence (CTD), object-background disentanglement (BG-Obj CTD), and reference-free perceptual metrics from VBench [3].

As summarized in Tab. 1, TTM improves trajectory adherence across both SVD and CogVideoX settings relative to strong training-free alternatives, while remaining competitive with training-based systems on perceptual metrics. Qualitatively, it better preserves first-frame identity under large displacements and reduces common failures such as background co-motion and geometric distortions. These trends support the core design goal: strong local control without retraining-induced specialization. Compared to single-clock behavior, we observe more stable object placement and less temporal freezing in surrounding regions, particularly on long or curved trajectories. Relative to GWTF variants, TTM also improves robustness when trajectory-induced layout changes are large, where noise-warping approaches often drift or deform scene structure.

3.2. Camera Motion Control

Following prior work [1, 6], we evaluate controlled camera trajectories on a DL3DV subset with depth-based re-projection references. Metrics include frame fidelity (MSE, LPIPS, SSIM), flow alignment, FID, and temporal consistency.

As shown in Tab. 2, TTM yields the strongest overall trade-off between motion fidelity and visual quality, with lower pixel/flow errors and better distributional alignment than GWTF variants. In practice, dual-clock guidance better preserves the intended camera path while avoiding drift and tearing artifacts from pure warping. Quantitatively, we

Method	Training Free?	CTD \downarrow	BG-Obj CTD \uparrow	Dynamic Degree \uparrow	Subject Consistency \uparrow	Background Consistency \uparrow	Motion Smoothness \uparrow	Aesthetic Quality \uparrow	Imaging Quality \uparrow
<i>SVD-Based Models</i>									
DragAnything	✗	10.645	50.885	0.981	0.956	0.942	0.983	0.531	0.554
SG-I2V*	✓	5.796	12.042	0.803	0.976	0.953	0.991	0.553	0.621
MotionPro	✗	8.685	24.485	0.422	0.979	0.975	0.993	0.559	0.617
Ours	✓	7.967	35.340	0.427	0.979	0.967	0.993	0.548	0.617
<i>CogVideoX-Based Models with Longer Generated Videos</i>									
GWTF $\gamma=0.7$	✗	32.548	86.614	0.736	0.963	0.965	0.989	0.517	0.539
GWTF $\gamma=0.5$	✗	27.844	87.708	0.764	0.958	0.963	0.988	0.513	0.539
Ours	✓	13.665	70.608	0.357	0.980	0.977	0.995	0.531	0.579

Table 1. Quantitative results on MC-Bench object motion control.

Method	MSE \downarrow	FID \downarrow	LPIPS \downarrow	SSIM \uparrow	CLIP Cons. \uparrow	Optical flow \downarrow
GWTF $\gamma=0.5$	0.033	25.990	0.371	0.526	0.981	76.714
GWTF $\gamma=0.7$	0.042	28.483	0.370	0.410	0.985	81.738
Warped	0.025	33.443	0.339	0.560	0.981	65.494
Ours	0.022	21.966	0.332	0.586	0.983	60.558

Table 2. Quantitative results on DL3DV camera motion control.

215 observe substantial improvements over the strongest GWTF
 216 setting, including notable gains in pixel-level fidelity and
 217 FID, alongside better optical-flow consistency. These gains
 218 align with qualitative behavior: TTM follows target camera
 219 motion more faithfully while preserving scene realism over
 220 time.

221 3.3. Appearance Control

222 Beyond motion, full-frame reference conditioning enables
 223 appearance-directed control relevant to creative produc-
 224 tion tasks: controlled color evolution, object insertion
 225 with scene-consistent blending, and motion-shape coupling.
 226 These edits are difficult to realize with trajectory-only con-
 227 ditioning. For production-style workflows, this is important
 228 because motion and look are often co-designed. In our ex-
 229 amples, users can manipulate both trajectory and per-pixel
 230 attributes in a single inference pass, enabling iterative art-
 231 direction loops without retraining.

232 3.4. Plug-and-Play Model Adaptation

233 TTM transfers without retraining to multiple I2V back-
 234 bones, including SVD, CogVideoX, and WAN 2.2 [13].
 235 This backbone-agnostic behavior is important for produc-
 236 tion workflows where models evolve rapidly and retraining
 237 costs are prohibitive. This portability is central for work-
 238 shop scenarios where teams prototype across heterogeneous
 239 model stacks and need a consistent control interface inde-
 240 pendent of backbone internals.

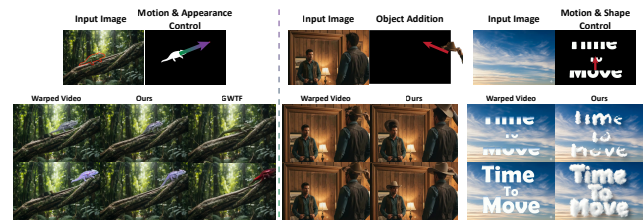


Figure 3. Joint motion and appearance control. TTM uses a user-provided warped reference to jointly control dynamics and per-pixel appearance in diverse cinematic editing scenarios.

241 4. Conclusions, Limitations, Future Directions

242 We presented TTM, a training-free framework for control-
 243 lable I2V generation that combines crude reference anima-
 244 tion with region-dependent dual-clock denoising. The
 245 method improves motion adherence while maintaining
 246 globally coherent dynamics and first-frame appearance, and
 247 transfers across multiple backbones without retraining. This
 248 combination directly targets practical creative workflows:
 249 users can edit local motion intent while preserving global
 250 scene plausibility, and can layer appearance changes with
 251 motion control in a unified inference procedure.

252 Limitations remain: the dual-clock parameters still re-
 253 quire per-backbone tuning, identity anchoring is strongest
 254 for content visible in the first frame, and high-quality con-
 255 trol currently assumes reasonably accurate masks. Fu-
 256 ture work includes multi-region and soft-mask scheduling,
 257 stronger long-horizon temporal coherence, and tighter inte-
 258 gration with film-production interfaces for iterative direct-
 259 ing and storyboarding.

260 References

- 261 [1] Ryan Burgert, Yuancheng Xu, Wenqi Xian, Oliver Pilarski,
 262 Pascal Clausen, Mingming He, Li Ma, Yitong Deng, Lingx-
 263 iao Li, Mohsen Mousavi, et al. Go-with-the-flow: Motion-
 264 controllable video diffusion models using real-time warped

- noise. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13–23, 2025. 2, 3
- [2] Daniel Geng, Charles Herrmann, Junhua Hur, Forrester Cole, Serena Zhang, Tobias Pfaff, Tatiana Lopez-Guevara, Yusuf Aytar, Michael Rubinstein, Chen Sun, et al. Motion prompting: Controlling video generation with motion trajectories. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1–12, 2025. 2
- [3] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 3
- [4] 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*, pages 8079–8088, 2024. 2
- [5] Yaowei Li, Xintao Wang, Zhaoyang Zhang, Zhouxia Wang, Ziyang Yuan, Liangbin Xie, Ying Shan, and Yuexian Zou. Image conductor: Precision control for interactive video synthesis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5031–5038, 2025. 2
- [6] Lu Ling, Yichen Sheng, Zhi Tu, Wentian Zhao, Cheng Xin, Kun Wan, Lantao Yu, Qianyu Guo, Zixun Yu, Yawen Lu, et al. D13dv-10k: A large-scale scene dataset for deep learning-based 3d vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22160–22169, 2024. 3
- [7] Wan-Duo Kurt Ma, John P Lewis, and W Bastiaan Kleijn. Trailblazer: Trajectory control for diffusion-based video generation. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024. 2
- [8] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 2
- [9] Koichi Namekata, Sherwin Bahmani, Ziyi Wu, Yash Kant, Igor Gilitschenski, and David B Lindell. Sg-i2v: Self-guided trajectory control in image-to-video generation. In *International Conference on Learning Representations*, 2025. 3
- [10] 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. 2
- [11] Xuanchi Ren, Tianchang Shen, Jiahui Huang, Huan Ling, Yifan Lu, Merlin Nimier-David, Thomas Müller, Alexander Keller, Sanja Fidler, and Jun Gao. Gen3c: 3d-informed world-consistent video generation with precise camera control. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6121–6132, 2025. 2
- [12] Ariel Shaulov, Itay Hazan, Lior Wolf, and Hila Chefer. Flowmo: Variance-based flow guidance for coherent motion in video generation. *arXiv preprint arXiv:2506.01144*, 2025. 2
- [13] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwei Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 4
- [14] Angtian Wang, Haibin Huang, Jacob Zhiyuan Fang, Yiding Yang, and Chongyang Ma. Ati: Any trajectory instruction for controllable video generation. *arXiv preprint arXiv:2505.22944*, 2025. 2
- [15] 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 *European Conference on Computer Vision*, pages 331–348. Springer, 2024. 3
- [16] 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.
- [17] Zhenghao Zhang, Junchao Liao, Menghao Li, Zuozhuo Dai, Bingxue Qiu, Siyu Zhu, Long Qin, and Weizhi Wang. Tora: Trajectory-oriented diffusion transformer for video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 2063–2073, 2025. 2
- [18] Zhongwei Zhang, Fuchen Long, Zhaofan Qiu, Yingwei Pan, Wu Liu, Ting Yao, and Tao Mei. Motionpro: A precise motion controller for image-to-video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 27957–27967, 2025. 2, 3
- [19] Haitao Zhou, Chuang Wang, Rui Nie, Jinlin Liu, Dongdong Yu, Qian Yu, and Changhu Wang. Trackgo: A flexible and efficient method for controllable video generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 10743–10751, 2025. 2