

PISCO: Precise Video Instance Insertion with Sparse Control

Anonymous CVPR submission

Paper ID

Abstract

001 *AI video generation is moving beyond general generation,*
 002 *which often relies on exhaustive prompt engineering and*
 003 *cherry-picking, toward fine-grained, controllable genera-*
 004 *tion and high-fidelity post-processing. In **professional AI-***
 005 ***assisted filmmaking**, a key capability is **video instance in-***
 006 ***sertion**: placing a specific object into an existing video*
 007 *at a designated spatial-temporal location while preserving*
 008 *scene dynamics and physically consistent interactions such*
 009 *as shadows and reflections. We present PISCO, a video*
 010 *diffusion framework for precise video instance insertion un-*
 011 *der arbitrary sparse keyframe control. PISCO supports*
 012 *single-frame, start-and-end-frame, or sparse keyframes at*
 013 *arbitrary timestamps, and automatically propagates appear-*
 014 *ance, motion, and interaction. To stabilize sparse condi-*
 015 *tioning, we introduce Variable-Information Guidance (VIG)*
 016 *and Distribution-Preserving Temporal Masking (DPTM),*
 017 *together with geometry-aware conditioning. We further con-*
 018 *struct PISCO-Bench, a benchmark with verified instance*
 019 *annotations and paired clean background videos. Experi-*
 020 *ments show that PISCO consistently outperforms represen-*
 021 *tative baselines and improves monotonically as additional*
 022 *control signals are provided.*

023 1. Introduction

024 Recent advances in large-scale video generative models [46]
 025 have enabled high-fidelity video generation with increas-
 026 ingly realistic motion. As these models continue to evolve,
 027 the focus of video generation is shifting beyond producing
 028 visually plausible content toward highly controllable video
 029 generation and editing [41], with the long-term goal of en-
 030 abling **Hollywood-grade AI-assisted filmmaking** [37, 39].
 031 In this new paradigm, the priority is no longer to **cherry-**
 032 **pick** a barely acceptable output from numerous generations,
 033 but to reliably **achieve precise user intent** with minimal
 034 iteration.

035 A particularly demanding yet under-explored capability

in this context is **precise video instance insertion**: inserting 036
 a specific object into an existing video at a precise user- 037
 specified spatial location and temporal position, while pre- 038
 serving the identity and dynamics of the original video. In 039
 professional visual effects and post-production workflows, 040
 instance insertion is not merely about adding a visually plu- 041
 sible object, but about achieving precise control with mini- 042
 mal iteration. In this setting, such precision entails several 043
 requirements: ❶ **Instance-level controllability**, enabling 044
 users to explicitly specify when and where an object appears; 045
 ❷ **Physically plausible temporal propagation**, where the 046
 inserted object automatically evolves over time with coherent 047
 pose and motion, following physically reasonable dynamics; 048
 ❸ **Consistent scene adaptation**, adjusting the background 049
 for insertion-induced effects such as shadows and reflec- 050
 tions; ❹ **Background dynamics preservation**, maintaining 051
 original motions, identities, and temporal patterns without 052
 disruption; ❺ **Low-effort user interaction**, where achieving 053
 the above goals does not require dense per-frame annotations 054
 or manual editing across the entire video. 055

Despite recent advances [57], these requirements re- 056
 main largely unmet by existing paradigms. **Video inpaint-** 057
ing [42, 67] demands exhaustive per-frame mask annota- 058
 tions (❶) and its copy-and-fill nature struggles with holistic 059
 scene adaptation such as realistic shadowing (❸). **Video-** 060
to-video editing [22, 54] can propagate appearance but 061
 lacks fine-grained spatiotemporal controllability for precise 062
 object placement (❶). **Agentic pipelines** (image editing + 063
 I2V) [18, 24] discard temporal context, often causing severe 064
 background drift and disrupted scene dynamics (❹). Finally, 065
geometry-heavy 4D approaches [23, 36] improve spatial 066
 reasoning but incur substantial computational overhead, rely 067
 on fragile geometric priors, and struggle with complex dy- 068
 namic object motions (❷). 069

To overcome these limitations, we propose PISCO 070
 (Precise Instance insertion with Sparse COntrol), a video 071
 diffusion framework designed for professional-grade inser- 072
 tion with minimal user effort. PISCO allows users to specify 073
 instance conditions at arbitrary sparse keyframes (e.g., first 074
 frame only, or first & last frames). Built upon the Wan back- 075
 bone [46], it incorporates a multi-channel context adapter to 076

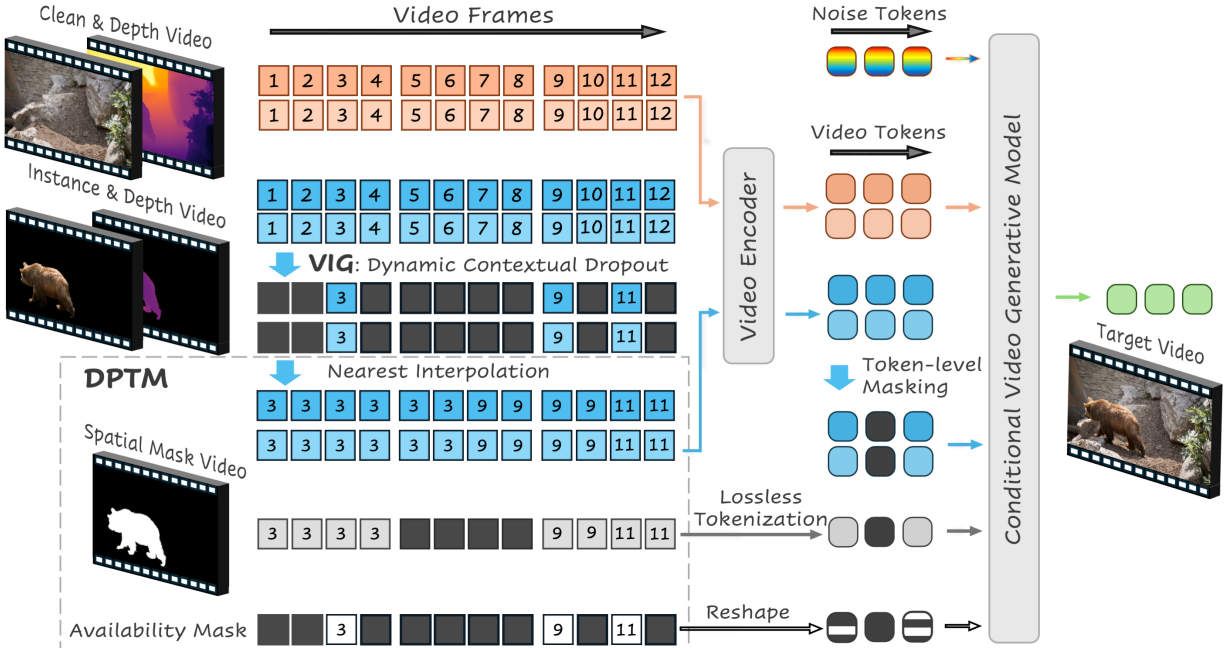


Figure 1. **Overview of PISCO.** PISCO performs precise video instance insertion under sparse keyframe control. VIG exposes the model to diverse conditioning densities, while DPTM stabilizes sparse conditioning for pretrained temporal VAEs. RGB, mask, depth, and availability cues are injected through a multi-channel context adapter.

171 model ϵ_θ with

$$172 \mathcal{L} = \mathbb{E}_{z_0, \epsilon, t, A} \left[\left\| \epsilon - \epsilon_\theta(z_t, t, V, D_V, I^A, D_I^A, M^A, A) \right\|_2 \right]. \quad (2)$$

173 In practice, ϵ_θ is implemented with the Wan backbone [46]
174 and a modified VACE context adapter [22].

175 3.2. Core Design

176 **Variable-Information Guidance.** To ensure robustness
177 across different user-control densities, we introduce
178 **Variable-Information Guidance (VIG)**, which samples
179 availability masks during training. The model is exposed
180 to a mixture of extremely sparse, sparse, dense, and fully
181 supervised regimes, enabling it to propagate motion under
182 limited guidance while retaining appearance fidelity when
183 stronger supervision is available.

184 **Distribution-Preserving Temporal Masking.** Modern
185 video generators often rely on pretrained temporal VAEs [1,
186 10, 32, 46, 49, 58, 65]. Directly zero-masking unavail-
187 able frames causes severe distribution shift. We there-
188 fore introduce **Distribution-Preserving Temporal Masking**
189 **(DPTM)**, which first fills missing frames in pixel space to
190 maintain natural encoder statistics, then applies masking in
191 token space and provides an aligned availability channel to
192 explicitly indicate which conditions are observed.

193 **Geometry- and appearance-aware training.** We further
194 improve insertion realism with depth-aware conditioning,

195 pseudo-amodal instance augmentation [2, 3], and relighting
196 augmentation using IC-Light [61]. These components im-
197 prove depth ordering, occlusion handling, and illumination
198 compatibility without requiring heavy 4D reconstruction. 198

199 4. Experiments

200 **Benchmark and settings.** We introduce PISCO-Bench, a
201 curated benchmark of 100 diverse real-world videos derived
202 from BURST [4]. We manually refine the instance masks
203 and use ROSE [38] to generate paired clean background
204 videos. We evaluate both PISCO-1.3B and PISCO-14B
205 under **First-frame** and **First & Last-frame** control, and
206 additionally report a **Five-Frame** setting for PISCO. Since
207 no open-source method supports identical sparse keyframe
208 conditioning, we adapt representative baselines from three
209 categories: **Image Editing + I2V Generation**, **Video In-**
210 **painting**, and **Reference-guided V2V Editing**. All methods
211 are evaluated on 49-frame videos at 832×480 resolution.

212 **Reference-based evaluation.** Using ground-truth pairs
213 (V, \hat{V}) , we evaluate generated videos with FVD [44],
214 LPIPS [62], PSNR [21], and SSIM [53], measured on both
215 the whole video and the foreground region. As shown in
216 Table 1, PISCO consistently outperforms all baselines.
217 Image Editing + I2V suffers from background hallucination,
218 inpainting methods preserve backgrounds but lag in tempo-
219 ral coherence and object fidelity, and VACE is the strongest
220 baseline overall. Still, PISCO-14B with First & Last control

Table 1. **Reference-based quantitative results.** Comparisons on both whole-video and foreground regions. Gray rows indicate the Five-Frame setting, reported to demonstrate scalability with additional sparse controls.

Method	Whole Video				Foreground			
	FVD ↓	LPIPS ↓	PSNR ↑	SSIM ↑	FVD ↓	LPIPS ↓	PSNR ↑	SSIM ↑
<i>Image Editing + I2V Generation</i>								
First Only	826	0.451	15.47	0.55	297	0.030	30.24	0.97
First + Last	624	0.392	16.44	0.56	250	0.030	30.38	0.98
<i>Video Inpainting</i>								
CoCoCo	590	0.191	23.62	0.80	398	0.031	30.26	0.97
VideoPainter	524	0.154	23.11	0.78	384	0.035	29.27	0.97
<i>Video-to-Video Editing</i>								
VACE _{14B}	371	0.103	25.55	0.88	273	0.028	30.55	0.98
UniVideo	485	0.211	19.22	0.61	310	0.031	29.21	0.97
<i>PISCO 1.3B (Ours)</i>								
First Only	398	0.121	25.35	0.87	243	0.029	30.51	0.98
First + Last	269	0.103	27.01	0.88	171	0.024	32.99	0.98
Five Frames	172	0.089	28.53	0.89	104	0.018	35.07	0.98
<i>PISCO 14B (Ours)</i>								
First Only	337	0.116	24.81	0.88	222	0.029	30.61	0.97
First + Last	204	0.097	26.58	0.89	138	0.022	33.58	0.98
Five Frames	136	0.084	28.01	0.90	75	0.015	35.94	0.98

Table 2. **Reference-free VBench comparison.** PISCO achieves the strongest overall perceptual quality among standard settings and continues to improve as additional sparse keyframes are provided.

Method	Background Consistency	Subject Consistency	Aesthetic Quality	Imaging Quality	Motion Smoothness	Overall Consistency	Temporal Flickering	Temporal Style	Average
<i>Image Editing + I2V Generation</i>									
First Only	91.43	83.86	48.83	59.01	97.89	15.11	95.89	15.11	63.39
First + Last	92.28	85.09	49.47	59.71	98.18	15.58	96.26	15.58	64.02
<i>Video Inpainting</i>									
CoCoCo	94.63	89.55	47.81	55.76	98.67	14.61	97.69	14.61	64.17
VideoPainter	94.51	89.14	47.82	57.68	98.97	13.85	97.77	13.85	64.20
<i>Video-to-Video Editing</i>									
VACE _{14B}	94.21	90.29	48.69	60.95	98.90	14.87	97.56	14.87	65.04
UniVideo	94.04	89.88	49.41	60.84	98.80	15.51	97.01	15.92	65.18
<i>PISCO 1.3B (Ours)</i>									
First Only	93.72	87.16	47.70	60.67	98.85	14.92	97.54	14.92	64.43
First + Last	94.07	91.33	49.48	61.18	98.86	15.58	97.51	15.58	65.45
Five Frames	94.23	91.45	50.01	62.09	98.86	15.61	97.67	15.61	65.89
<i>PISCO 14B (Ours)</i>									
First Only	93.84	87.26	48.24	60.80	98.79	15.14	97.26	15.14	64.56
First + Last	94.20	91.57	50.08	62.00	98.79	15.64	97.21	15.64	65.64
Five Frames	94.42	91.98	51.45	62.87	98.89	15.82	97.34	15.57	66.04

221 reduces whole-video FVD from 371 to 204 and foreground
 222 FVD from 273 to 138. Additional sparse controls further
 223 improve performance in the Five-Frame setting.

224 **Reference-free evaluation.** For reference-free evaluation,
 225 we employ VBench [20] and adapt it with instance masks
 226 to separately assess subject and background consistency.
 227 As shown in Table 2, PISCO achieves the strongest over-
 228 all perceptual quality among standard settings. In particu-
 229 lar, it yields better subject consistency and visual quality
 230 than prior baselines, while remaining competitive on back-
 231 ground preservation. Performance further improves when
 232 more sparse controls are provided.



Figure 2. **Qualitative comparison.** Existing methods struggle with background preservation, instance fidelity, or spatial-temporal alignment. In contrast, PISCO produces stronger object fidelity, more consistent scale and motion, and better preservation of original scene dynamics, especially under sparse First & Last control.

233 **Qualitative results.** Figure 2 shows that prior methods often
 234 fail in different ways: agentic pipelines tend to regenerate
 235 the scene and drift in the background, inpainting methods
 236 weaken or blur the inserted object, and V2V baselines may
 237 lose precise long-term control. In contrast, PISCO achieves
 238 stronger spatiotemporal alignment and scene consistency,
 239 especially when sparse endpoint constraints are provided.

5. Conclusion

240 We presented PISCO, a video diffusion framework for
 241 precise video instance insertion under sparse user control.
 242 By combining Variable-Information Guidance, Distribution-
 243 Preserving Temporal Masking, and geometry-aware condi-
 244 tioning, PISCO addresses key challenges in sparse-control
 245 video generation and significantly outperforms represen-
 246 tative inpainting, video editing, and agentic baselines on
 247 PISCO-Bench. Moreover, it improves monotonically as
 248 more sparse control frames are provided, supporting a con-
 249 trollable and scalable paradigm for instance-centric video
 250 editing.
 251

252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307

References

- [1] Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. *arXiv preprint arXiv:2501.03575*, 2025. 3
- [2] Jiayang Ao, Qiuhong Ke, and Krista A Ehinger. Image amodal completion: A survey. *Computer Vision and Image Understanding*, 229:103661, 2023. 3
- [3] Jiayang Ao, Yanbei Jiang, Qiuhong Ke, and Krista A Ehinger. Open-world amodal appearance completion. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 6490–6499, 2025. 3
- [4] Ali Athar, Jonathon Luiten, Paul Voigtlaender, Tarasha Khurana, Achal Dave, Bastian Leibe, and Deva Ramanan. Burst: A benchmark for unifying object recognition, segmentation and tracking in video. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 1674–1683, 2023. 2, 3
- [5] Chen Bai, Zeman Shao, Guoxiang Zhang, Di Liang, Jie Yang, Zhuorui Zhang, Yujian Guo, Chengzhang Zhong, Yiqiao Qiu, Zhendong Wang, et al. Anything in any scene: Photorealistic video object insertion. *arXiv preprint arXiv:2401.17509*, 2024. 2
- [6] Yuxuan Bian, Zhaoyang Zhang, Xuan Ju, Mingdeng Cao, Liangbin Xie, Ying Shan, and Qiang Xu. Videopainter: Any-length video inpainting and editing with plug-and-play context control. In *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers*, pages 1–12, 2025. 2
- [7] Huanqia Cai, Sihan Cao, Ruoyi Du, Peng Gao, Steven Hoi, Shijie Huang, Zhaohui Hou, Dengyang Jiang, Xin Jin, Liangchen Li, et al. Z-image: An efficient image generation foundation model with single-stream diffusion transformer. *arXiv preprint arXiv:2511.22699*, 2025. 2
- [8] Jiayin Cai, Changlin Li, Xin Tao, Chun Yuan, and Yu-Wing Tai. Devit: Deformed vision transformers in video inpainting. In *Proceedings of the 30th ACM international conference on multimedia*, pages 779–789, 2022. 2
- [9] Ya-Liang Chang, Zhe Yu Liu, Kuan-Ying Lee, and Winston Hsu. Free-form video inpainting with 3d gated convolution and temporal patchgan. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9066–9075, 2019. 2
- [10] Shoufa Chen, Chongjian Ge, Yuqi Zhang, Yida Zhang, Fengda Zhu, Hao Yang, Hongxiang Hao, Hui Wu, Zhichao Lai, Yifei Hu, et al. Goku: Flow based video generative foundation models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 23516–23527, 2025. 3
- [11] Yiyang Chen, Xuanhua He, Xiujun Ma, and Yue Ma. Contextflow: Training-free video object editing via adaptive context enrichment. *arXiv preprint arXiv:2509.17818*, 2025. 2
- [12] Gang Cheng, Xin Gao, Li Hu, Siqi Hu, Mingyang Huang, Chaonan Ji, Ju Li, Dechao Meng, Jinwei Qi, Penchong Qiao, et al. Wan-animate: Unified character animation and replacement with holistic replication. *arXiv preprint arXiv:2509.14055*, 2025. 2
- [13] Chen Gao, Ayush Saraf, Jia-Bin Huang, and Johannes Kopf. Flow-edge guided video completion. In *European Conference on Computer Vision*, pages 713–729. Springer, 2020. 2
- [14] Dylan Green, William Harvey, Saeid Naderiparizi, Matthew Niedoba, Yunpeng Liu, Xiaoxuan Liang, Jonathan Lavington, Ke Zhang, Vasileios Lioutas, Setareh Dabiri, et al. Semantically consistent video inpainting with conditional diffusion models. *arXiv preprint arXiv:2405.00251*, 2024. 2
- [15] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yao-hui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 322
- [16] Yuwei Guo, Ceyuan Yang, Anyi Rao, Chenlin Meng, Omer Bar-Tal, Shuangrui Ding, Maneesh Agrawala, Dahua Lin, and Bo Dai. Keyframe-guided creative video inpainting. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13009–13020, 2025. 2
- [17] Yingqing He, Menghan Xia, Haoxin Chen, Xiaodong Cun, Yuan Gong, Jinbo Xing, Yong Zhang, Xintao Wang, Chao Weng, Ying Shan, et al. Animate-a-story: Storytelling with retrieval-augmented video generation. *arXiv preprint arXiv:2307.06940*, 2023. 2
- [18] Yaosi Hu, Chong Luo, and Zhenzhong Chen. Make it move: controllable image-to-video generation with text descriptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18219–18228, 2022. 1, 2
- [19] Yuan-Ting Hu, Heng Wang, Nicolas Ballas, Kristen Grauman, and Alexander G Schwing. Proposal-based video completion. In *European Conference on Computer Vision*, pages 38–54. Springer, 2020. 2
- [20] Ziqi Huang, Yanan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 4
- [21] Bernd Jähne. *Digital image processing*. Springer, 2005. 3
- [22] Zeyinzi Jiang, Zhen Han, Chaojie Mao, Jingfeng Zhang, Yulin Pan, and Yu Liu. Vace: All-in-one video creation and editing. *arXiv preprint arXiv:2503.07598*, 2025. 1, 2, 3
- [23] Hoiyeong Jin, Hyojin Jang, Jeongho Kim, Junha Hyung, Kinam Kim, Dongjin Kim, Huijin Choi, Hyeonji Kim, and Jaegul Choo. Insertanywhere: Bridging 4d scene geometry and diffusion models for realistic video object insertion. *arXiv preprint arXiv:2512.17504*, 2025. 1, 2
- [24] Yang Jin, Zhicheng Sun, Ningyuan Li, Kun Xu, Hao Jiang, Nan Zhuang, Quzhe Huang, Yang Song, Yadong Mu, and Zhouchen Lin. Pyramidal flow matching for efficient video generative modeling. *arXiv preprint arXiv:2410.05954*, 2024. 1, 2
- [25] Dahun Kim, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Deep video inpainting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5792–5801, 2019. 2

- 366 [26] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Fred-
367 eric Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn,
368 Jack English, Zion English, Patrick Esser, et al. Flux. 1 kon-
369 text: Flow matching for in-context image generation and
370 editing in latent space. *arXiv preprint arXiv:2506.15742*,
371 2025. 2
- 372 [27] Minhyeok Lee, Suhwan Cho, Chajin Shin, Jungho Lee,
373 Sunghun Yang, and Sangyoun Lee. Video diffusion mod-
374 els are strong video inpainter. In *Proceedings of the AAAI*
375 *Conference on Artificial Intelligence*, pages 4526–4533, 2025.
376 2
- 377 [28] Ang Li, Shanshan Zhao, Xingjun Ma, Mingming Gong,
378 Jianzhong Qi, Rui Zhang, Dacheng Tao, and Ramamoha-
379 narao Kotagiri. Short-term and long-term context aggregation
380 network for video inpainting. In *European Conference on*
381 *Computer Vision*, pages 728–743. Springer, 2020. 2
- 382 [29] Xiaowen Li, Haolan Xue, Peiran Ren, and Liefeng Bo. Dif-
383 fueraaser: A diffusion model for video inpainting. *arXiv*
384 *preprint arXiv:2501.10018*, 2025. 2
- 385 [30] Zhen Li, Cheng-Ze Lu, Jianhua Qin, Chun-Le Guo, and Ming-
386 Ming Cheng. Towards an end-to-end framework for flow-
387 guided video inpainting. In *Proceedings of the IEEE/CVF*
388 *conference on computer vision and pattern recognition*, pages
389 17562–17571, 2022. 2
- 390 [31] Ziyi Li, Hao Luo, Xincheng Shuai, and Henghui Ding.
391 Anyi2v: Animating any conditional image with motion con-
392 trol. *arXiv preprint arXiv:2507.02857*, 2025. 2
- 393 [32] Jinlai Liu, Jian Han, Bin Yan, Hui Wu, Fengda Zhu, Xing
394 Wang, Yi Jiang, Bingyue Peng, and Zehuan Yuan. Infini-
395 tystar: Unified spacetime autoregressive modeling for visual
396 generation. *arXiv preprint arXiv:2511.04675*, 2025. 3
- 397 [33] Rui Liu, Hanming Deng, Yangyi Huang, Xiaoyu Shi, Lewei
398 Lu, Wenxiu Sun, Xiaogang Wang, Jifeng Dai, and Hong-
399 sheng Li. Decoupled spatial-temporal transformer for video
400 inpainting. *arXiv preprint arXiv:2104.06637*, 2021. 2
- 401 [34] Rui Liu, Hanming Deng, Yangyi Huang, Xiaoyu Shi, Lewei
402 Lu, Wenxiu Sun, Xiaogang Wang, Jifeng Dai, and Hongsheng
403 Li. Fuseformer: Fusing fine-grained information in transform-
404 ers for video inpainting. In *Proceedings of the IEEE/CVF*
405 *international conference on computer vision*, pages 14040–
406 14049, 2021. 2
- 407 [35] Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang,
408 Wei Cheng, Jiaqi Liao, Yingming Wang, Honghao Fu, Chun-
409 rui Han, et al. Step1x-edit: A practical framework for general
410 image editing. *arXiv preprint arXiv:2504.17761*, 2025. 2
- 411 [36] Ziling Liu, Jinyu Yang, Mingqi Gao, and Feng Zheng. Place
412 anything into any video. *arXiv preprint arXiv:2402.14316*,
413 2024. 1, 2
- 414 [37] Lev Manovich. The rise of generative media. Manovich.net,
415 2023. From the collection "Artificial Aesthetics". 1
- 416 [38] Chenxuan Miao, Yutong Feng, Jianshu Zeng, Zixiang Gao,
417 Hantang Liu, Yunfeng Yan, Donglian Qi, Xi Chen, Bin Wang,
418 and Hengshuang Zhao. Rose: Remove objects with side
419 effects in videos. *arXiv preprint arXiv:2508.18633*, 2025. 2,
420 3
- 421 [39] Justine Moore. Why 2023 was ai video’s breakout year, and
422 what to expect in 2024, 2024. Accessed: 2024-05-20. 1
- 423 [40] Koichi Namekata, Sherwin Bahmani, Ziyi Wu, Yash Kant,
424 Igor Gilitschenski, and David B Lindell. Sg-i2v: Self-
425 guided trajectory control in image-to-video generation. *arXiv*
426 *preprint arXiv:2411.04989*, 2024. 2
- 427 [41] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra,
428 Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-
429 Yao Ma, Ching-Yao Chuang, et al. Movie gen: A cast of
430 media foundation models. *arXiv preprint arXiv:2410.13720*,
431 2024. 1
- 432 [42] Weize Quan, Jiayi Chen, Yanli Liu, Dong-Ming Yan, and Pe-
433 ter Wonka. Deep learning-based image and video inpainting:
434 A survey. *International Journal of Computer Vision*, 132(7):
435 2367–2400, 2024. 1, 2
- 436 [43] Shuwei Shi, Biao Gong, Xi Chen, Dandan Zheng, Shuai Tan,
437 Zizheng Yang, Yuyuan Li, Jingwen He, Kecheng Zheng, Jing-
438 dong Chen, et al. Motionstone: Decoupled motion intensity
439 modulation with diffusion transformer for image-to-video
440 generation. In *Proceedings of the Computer Vision and Pat-
441 tern Recognition Conference*, pages 22864–22874, 2025. 2
- 442 [44] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elho-
443 seiny. Stylegan-v: A continuous video generator with the
444 price, image quality and perks of stylegan2. In *Proceedings*
445 *of the IEEE/CVF conference on computer vision and pattern*
446 *recognition*, pages 3626–3636, 2022. 3
- 447 [45] Yuanpeng Tu, Hao Luo, Xi Chen, Sihui Ji, Xiang Bai, and
448 Hengshuang Zhao. Videoanydoor: High-fidelity video object
449 insertion with precise motion control, 2025. 2
- 450 [46] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao,
451 Chen-Wei Xie, Di Chen, Feiwei Yu, Haiming Zhao, Jianxiao
452 Yang, et al. Wan: Open and advanced large-scale video
453 generative models. *arXiv preprint arXiv:2503.20314*, 2025.
454 1, 2, 3
- 455 [47] Zhen Wan, Chenyang Qi, Zhiheng Liu, Tao Gui, and Yue Ma.
456 Unipaint: Unified space-time video inpainting via mixture-
457 of-experts. In *Proceedings of the IEEE/CVF International*
458 *Conference on Computer Vision*, pages 1861–1871, 2025. 2
- 459 [48] Chuan Wang, Haibin Huang, Xiaoguang Han, and Jue Wang.
460 Video inpainting by jointly learning temporal structure and
461 spatial details. In *Proceedings of the AAAI conference on*
462 *artificial intelligence*, pages 5232–5239, 2019. 2
- 463 [49] Junke Wang, Yi Jiang, Zehuan Yuan, Bingyue Peng, Zuxuan
464 Wu, and Yu-Gang Jiang. Omnitokenizer: A joint image-
465 video tokenizer for visual generation. *Advances in Neural*
466 *Information Processing Systems*, 37:28281–28295, 2024. 3
- 467 [50] Wenhao Wang and Yi Yang. Tip-i2v: A million-scale real
468 text and image prompt dataset for image-to-video generation.
469 In *Proceedings of the IEEE/CVF International Conference*
470 *on Computer Vision*, pages 14898–14908, 2025. 2
- 471 [51] Xiaojuan Wang, Boyang Zhou, Brian Curless, Ira
472 Kemelmacher-Shlizerman, Aleksander Holynski, and
473 Steven M Seitz. Generative inbetweening: Adapting
474 image-to-video models for keyframe interpolation. *arXiv*
475 *preprint arXiv:2408.15239*, 2024. 2
- 476 [52] Yukun Wang, Longguang Wang, Zhiyuan Ma, Qibin Hu, Kai
477 Xu, and Yulan Guo. Videodirector: Precise video editing via
478 text-to-video models. In *Proceedings of the Computer Vision*
479 *and Pattern Recognition Conference*, pages 2589–2598, 2025.
480 2

- 481 [53] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P
482 Simoncelli. Image quality assessment: from error visibility to
483 structural similarity. *IEEE transactions on image processing*,
484 13(4):600–612, 2004. 3
- 485 [54] Cong Wei, Quande Liu, Zixuan Ye, Qiulin Wang, Xintao
486 Wang, Pengfei Wan, Kun Gai, and Wenhui Chen. Univideo:
487 Unified understanding, generation, and editing for videos.
488 *arXiv preprint arXiv:2510.08377*, 2025. 1, 2
- 489 [55] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan
490 Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei
491 Chen, et al. Qwen-image technical report. *arXiv preprint*
492 *arXiv:2508.02324*, 2025. 2
- 493 [56] Rui Xu, Xiaoxiao Li, Bolei Zhou, and Chen Change Loy.
494 Deep flow-guided video inpainting. In *Proceedings of the*
495 *IEEE/CVF conference on computer vision and pattern recog-*
496 *niton*, pages 3723–3732, 2019. 2
- 497 [57] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Run-
498 sheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-
499 Hsuan Yang. Diffusion models: A comprehensive survey of
500 methods and applications. *ACM computing surveys*, 56(4):
501 1–39, 2023. 1, 2
- 502 [58] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu
503 Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaoh-
504 an Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video
505 diffusion models with an expert transformer. *arXiv preprint*
506 *arXiv:2408.06072*, 2024. 3
- 507 [59] Zuhao Yang, Jiahui Zhang, Yingchen Yu, Shijian Lu, and
508 Song Bai. Versatile transition generation with image-to-video
509 diffusion. In *Proceedings of the IEEE/CVF International*
510 *Conference on Computer Vision*, pages 16981–16990, 2025.
511 2
- 512 [60] Guy Yariv, Yuval Kirstain, Amit Zohar, Shelly Sheynin,
513 Yaniv Taigman, Yossi Adi, Sagie Benaim, and Adam
514 Polyak. Through-the-mask: Mask-based motion trajecto-
515 ries for image-to-video generation. In *Proceedings of the*
516 *Computer Vision and Pattern Recognition Conference*, pages
517 18198–18208, 2025. 2
- 518 [61] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Scaling
519 in-the-wild training for diffusion-based illumination harmo-
520 nization and editing by imposing consistent light transport. In
521 *The Thirteenth International Conference on Learning Repre-*
522 *sentations*, 2025. 3
- 523 [62] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman,
524 and Oliver Wang. The unreasonable effectiveness of deep
525 features as a perceptual metric. In *Proceedings of the IEEE*
526 *conference on computer vision and pattern recognition*, pages
527 586–595, 2018. 3
- 528 [63] Zhixing Zhang, Bichen Wu, Xiaoyan Wang, Yaqiao Luo,
529 Luxin Zhang, Yinan Zhao, Peter Vajda, Dimitris Metaxas,
530 and Licheng Yu. Avid: Any-length video inpainting with dif-
531 fusion model. In *Proceedings of the IEEE/CVF conference on*
532 *computer vision and pattern recognition*, pages 7162–7172,
533 2024. 2
- 534 [64] Zhongwei Zhang, Fuchen Long, Zhaofan Qiu, Yingwei Pan,
535 Wu Liu, Ting Yao, and Tao Mei. Motionpro: A precise motion
536 controller for image-to-video generation. In *Proceedings of*
537 *the Computer Vision and Pattern Recognition Conference*,
538 pages 27957–27967, 2025. 2
- [65] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen,
Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang
You. Open-sora: Democratizing efficient video production
for all. *arXiv preprint arXiv:2412.20404*, 2024. 3
- [66] Shangchen Zhou, Chongyi Li, Kelvin CK Chan, and
Chen Change Loy. Propainter: Improving propagation
and transformer for video inpainting. In *Proceedings of*
the IEEE/CVF international conference on computer vision,
pages 10477–10486, 2023. 2
- [67] Bojia Zi, Shihao Zhao, Xianbiao Qi, Jianan Wang, Yukai Shi,
Qianyu Chen, Bin Liang, Rong Xiao, Kam-Fai Wong, and
Lei Zhang. Cococo: Improving text-guided video inpainting
for better consistency, controllability and compatibility. In
Proceedings of the AAAI Conference on Artificial Intelligence,
pages 11067–11076, 2025. 1, 2
- [68] Xueyan Zou, Linjie Yang, Ding Liu, and Yong Jae Lee. Pro-
gressive temporal feature alignment network for video inpaint-
ing. In *Proceedings of the IEEE/CVF Conference on Com-*
puter Vision and Pattern Recognition, pages 16448–16457,
2021. 2