

Physics-Aware Video Instance Removal Benchmark

Anonymous CVPR submission

Paper ID

Abstract

001 **Video Instance Removal (VIR)** requires removing target
002 objects while maintaining background integrity and phys-
003 ical consistency, such as specular reflections and illumi-
004 nation interactions. Despite advancements in text-guided
005 editing, current benchmarks primarily assess visual plau-
006 sibility, often overlooking the physical causalities—such as
007 lingering shadows—triggered by object removal. We intro-
008 duce the **Physics-Aware Video Instance Removal (PVIR)**
009 benchmark, featuring 95 high-quality videos annotated with
010 instance-accurate masks and removal prompts. PVIR is
011 partitioned into Simple and Hard subsets, the latter ex-
012 plicitly targeting complex physical interactions. We eval-
013 uate four representative methods—PISCO-Removal, Uni-
014 Video, DiffuEraser, and CoCoCo—using a decoupled hu-
015 man evaluation protocol across three dimensions to isolate
016 semantic, visual, and spatial failures: instruction following,
017 rendering quality, and edit exclusivity. Our results show
018 that **PISCO-Removal** and **UniVideo** achieve state-of-the-
019 art performance, while **DiffuEraser** frequently introduces
020 blurring artifacts and **CoCoCo** struggles significantly with
021 instruction following. The persistent performance drop on
022 the Hard subset highlights the ongoing challenge of recov-
023 ering complex physical side effects.

024 1. Introduction

025 Removing a target instance from a video is a foundational
026 editing capability with direct utility in post-production, pri-
027 vacy protection, robotics simulation, and synthetic data cu-
028 ration. Compared with image object removal, video in-
029 stance removal requires consistency over time and consis-
030 tency with scene physics. When an object is removed, not
031 only should the object disappear, but its side effects should
032 be updated as well: reflections in windows, mirror appear-
033 ances, indirect occlusion patterns, and local illumination
034 cues should all remain plausible. These requirements ex-
035 pose a major gap between qualitative demos and reliable,
036 comparable evaluation.

037 Recent video editing and inpainting systems have im-

proved temporal coherence and controllability [11, 23, 28].
At the same time, methods designed for side-effect-aware
removal demonstrate that physics-aware editing is now an
explicit research target [14]. Despite these algorithmic ad-
vances, the field lacks a common yardstick: existing evalu-
ations are typically performed on private subsets, inconsis-
tent prompts, and varying resolution constraints, leading to
a “closed-world” comparison that obscures true progress.
As a result, it remains difficult to answer basic questions:
Which model follows removal instructions best? Which
model preserves visual quality over time? Which model
minimizes unintended edits outside the target region?

To address this gap, we introduce **Physics-Aware Video
Instance Removal Benchmark**, a task-focused benchmark
that standardizes data, protocols, and evaluation. Our
benchmark contains 95 high-quality videos with per-video
target segmentation and removal prompts. We explicitly
partition data into *Simple* and *Hard* subsets, where *Hard*
clips include stronger interactions with real-world physics
(e.g., specular reflections, mirror appearance, and pro-
nounced scene coupling).

A second key contribution is a decoupled human eval-
uation protocol. Instead of a single holistic score, we as-
sess three independent dimensions: (1) **Instruction Follow-
ing**, i.e., whether the target is correctly removed; (2) **Ren-
dering Quality**, i.e., whether the inpainted result is tem-
porally stable and visually plausible; (3) **Edit Exclusivity**,
i.e., whether non-target content remains unchanged. Each
dimension uses a 1–4 rubric with explicit criteria, enabling
interpretable diagnosis rather than one-number ranking.

We evaluate four representative models under a unified
setting: CoCoCo [28], UniVideo [23], DiffuEraser [11],
and PISCO-Removal [3]. For PISCO-Removal, we eval-
uate the variant trained in the PISCO paper [3] using the
ROSE dataset [14]; this version supports more demanding
practical configurations, including 720p resolution, portrait
orientations, and sequences up to 120 frames. This cross-
model comparison is designed to expose the trade-off be-
tween general-purpose inpainting stability and task-specific
physical fidelity.

078	Contributions.	Our contributions are summarized as fol-	126
079		low:	127
080	• We present a new benchmark for <i>physics-aware video</i>		128
081	<i>instance removal</i> , with 95 high-quality videos, instance-		129
082	level masks, and removal prompts.		130
083	• We define a decoupled, interpretable human evaluation		131
084	protocol with three independent 1–4 metrics and an ex-		132
085	PLICIT aggregation rule.		133
086	• We provide a unified benchmark of four representative		134
087	methods and establish analysis protocols for overall, per-		135
088	difficulty, and failure-mode evaluation.		136
089	Scope and current status.	This paper focuses on bench-	
090	mark construction and standardized evaluation rather than		
091	proposing a new removal architecture. Our analysis un-		
092	covers a "performance ceiling" where even state-of-the-art		
093	models fail to resolve secondary physical interactions, re-		
094	vealing that the primary bottleneck in VIR has shifted from		
095	temporal flickering to physical incoherence.		
096	2. Related Works		
097	Evolution of Video Inpainting and Instance Removal.		
098	Modern video object removal pipelines inherit core spa-		
099	tial advances from image inpainting [15, 20, 26] and		
100	have evolved through flow-guided propagation [2, 12] and		
101	transformer-based temporal modeling [13, 27] to ensure		
102	long-horizon consistency. Recent diffusion-based formu-		
103	lations, such as DiffuEraser [11], further enhance realism		
104	under complex motion, while ROSE [14] specifically for-		
105	mulates physical interactions like reflections as first-order		
106	constraints. However, while these methods focus on the		
107	architectural challenge of texture propagation and physical		
108	coupling, our work shifts the focus toward providing a stan-		
109	dardized evaluation framework to quantify how well these		
110	models actually recover such complex physical causalities.		
111	Text-guided Video Generation and Editing.	The field	
112	has shifted from general video generation [6, 18] to highly		
113	controllable editing pipelines [4, 10, 24]. Recent large-		
114	scale models like UniVideo [23] and CoCoCo [28] leverage		
115	strong generative priors to achieve high-quality restoration		
116	and better workflow compatibility. Despite their impressive		
117	zero-shot capabilities, these generative models often strug-		
118	gle with "semantic leakage" or fail to strictly respect local		
119	masking constraints in instance removal tasks. Unlike these		
120	generative frameworks that prioritize visual plausibility, our		
121	benchmark emphasizes the decoupling of instruction fol-		
122	lowing and spatial exclusivity to expose these specific fail-		
123	ure modes.		
124	Benchmarking and Perceptual Assessment.	While	
125	general-purpose video benchmarks [1, 7, 22] and visual		
		metrics like FID [5] and FVD [21] have improved cov-	126
		erage, they often fail to capture high-level physical logic,	127
		such as lingering shadows or inconsistent illumination. Ex-	128
		isting data infrastructures from segmentation [17, 25] pro-	129
		vide annotation principles but do not isolate the entangled	130
		dimensions of target compliance and background preserva-	131
		tion. Our PVIR benchmark addresses this gap by introduc-	132
		ing a physics-aware dataset and a decoupled human evalu-	133
		ation protocol with a 1–4 scoring rubric, providing a more	134
		interpretable assessment than conventional automatic met-	135
		rics.	136
	3. Dataset: Physics-Aware Video Instance Re-		137
	moval		138
	Design goal.	The dataset is designed for one core task: re-	139
		move a designated instance from a real video while preserv-	140
		ing non-target content and maintaining physically plausible	141
		side effects. To enable robust benchmarking, we prioritize	142
		high visual quality, diverse scenes, and explicit interaction	143
		complexity.	144
	Scale and split.	The benchmark comprises 95 high-	145
		quality sequences, ensuring a balanced distribution between	146
		foundational removal tasks and advanced physics-aware	147
		challenges. We organize them into two subsets: Simple	148
		and Hard . Simple videos (57 videos) usually contain ob-	149
		jects with simpler geometry and weaker coupling to scene	150
		physics. They serve as a baseline to evaluate a model's fun-	151
		damental ability to maintain spatial-temporal coherence in	152
		the absence of complex physics. Hard videos (38 videos)	153
		include stronger physics interactions, such as mirror reflec-	154
		tions, specular highlights, and complex motion-appearance	155
		coupling. These cases require the model to not only fill the	156
		disoccluded pixels but also maintain physical causality by	157
		updating or removing secondary side effects, such as reflec-	158
		tions and shadows, that are anchored to the target. Each	159
		video includes two annotations:	160
		• Target segmentation mask: a high-quality instance	161
		mask indicating the object to remove.	162
		• Removal prompt: a natural-language instruction that	163
		unambiguously refers to the target instance.	164
		These two elements define a standardized input interface for	165
		all benchmarked models.	166
	3.1. Collection and Annotation Pipeline		167
	Data sourcing and pre-filtering.	Candidate videos are	168
		collected from diverse real-world scenes and filtered by	169
		minimal requirements on spatial quality, temporal smooth-	170
		ness, and compression artifacts. To avoid overly syn-	171
		thetic bias, we prioritize clips with natural motion and	172
		illumination variation. Our primary sources include	173
		the Inter4k dataset [19], selected for its ultra-high-	174

175 definition (UHD) clarity and high frame rates, and the
176 **DAVIS2016** dataset [16], chosen for its diversity in object-
177 to-background interactions. Specifically, we curate 45 se-
178 quences from Inter4k to ensure high-fidelity rendering as-
179 sessment, and 50 sequences from DAVIS2016 to leverage
180 its complex motion patterns. All source videos are either li-
181 censed under Creative Commons (CC) or are explicitly per-
182 mitted for non-commercial research use, ensuring a clear
183 and ethical release metadata.

184 **Mask annotation workflow.** To ensure pixel-level preci-
185 sion, we employ a multi-stage annotation pipeline. Annota-
186 tors first utilize the Segment Anything Model 2 (SAM 2) [9]
187 to generate initial object tracks across each sequence. By
188 providing sparse point or box prompts on keyframes, SAM
189 2’s memory-based propagation produces coarse masks for
190 the entire video. Subsequently, annotators perform manual
191 frame-level refinement to clean up boundaries, particularly
192 in challenging cases involving motion blur or occlusion. To
193 ensure temporal smoothness, we conduct a final quality-
194 control pass focusing on reducing “shape jitter” (flickering
195 boundaries). Any masks exhibiting temporal instability are
196 manually corrected or re-propagated using flow-based con-
197 sistency checks. This hybrid workflow combines the effi-
198 ciency of foundation models with the rigor of human ver-
199 ification, yielding the instance-accurate masks required for
200 high-fidelity removal.

201 **Prompt writing workflow.** Prompts are designed to be
202 specific, concise, and target-disambiguating. To ensure
203 consistency across the benchmark, we adopt a structured
204 template: [Action] + [Target Attributes] +
205 [Spatial/Contextual Qualifiers]. For in-
206 stance, a PVIR prompt specifies “Remove the silver sedan
207 parked under the flickering streetlamp” instead of a generic
208 command, providing unique semantic grounding that mini-
209 mizes instruction ambiguity.

210 About the language policy and disambiguation, all
211 prompts are authored in English using standard descriptive
212 vocabulary. Ambiguous references—such as multiple simi-
213 lar objects in a single frame—are strictly disallowed unless
214 unique qualifiers (e.g., “the person on the far left”) are in-
215 cluded. For the *Hard* subset, prompts are intentionally aug-
216 mented with physical context, such as “including its reflec-
217 tion on the water surface,” to explicitly signal the expected
218 physics-aware behavior. Each prompt undergoes a cross-
219 verification pass by a second annotator to ensure that the
220 textual description uniquely identifies the instance masked
221 in the ground truth.

222 3.2. Difficulty Taxonomy

223 **Simple subset.** Simple clips typically contain targets with
224 limited appearance variation and weak side-effect coupling.

Examples include matte surfaces, stable backgrounds, and
225 short-term motion without heavy occlusion. 226

Hard subset. Hard clips emphasize complex light-
227 surface coupling and geometric reconstruction challenges. 228

Table 1. **Summary of the PVIR Benchmark Dataset.** The benchmark comprises 95 high-definition videos, categorized into Simple and Hard subsets based on the complexity of physical interactions (e.g., reflections, shadows, and dynamic fluid wakes).

Subset	#Videos	Avg. Frames	Key Properties
Simple	57	81	weaker interaction and simpler geometry
Hard	38	81	reflection/mirror/specular and stronger coupling
Total	95	81	mixed scenes and motions

4. Benchmark Protocol 229

Task definition. Given an input video V , a target mask sequence M , and a removal prompt P , the model outputs an edited video \hat{V} where the target instance is removed. The benchmark requires: (1) complete target removal alongside its associated physical derivatives, (2) high-quality temporally consistent rendering, (3) minimal unintended changes outside the edited region. 230
231
232
233
234
235
236

Evaluated models. We evaluate four representative meth-
237 ods: **CoCoCo** [28], **UniVideo** [23], **DiffuEraser** [11],
238 and **PISCO-Removal** [3], where PISCO-Removal is
239 a WACE [8]-like model fine-tuned with the ROSE
240 dataset [14] that support support 720p, portrait videos, and
241 up to 120 frames [3]. 242

4.1. Unified Inference Setup 243

Input/output interface normalization. Each model re-
244 ceives the same semantic input triplet (V, M, P) and out-
245 puts a completed video sequence. We standardize video de-
246 coding, prompt formatting, and export codecs to minimize
247 evaluation noise from non-model factors. 248

Method-specific adaptation policy. To eliminate perfor-
249 mance bias stemming from disparate input constraints, all
250 evaluated baselines are unified under a standardized con-
251 figuration: 720p resolution at 81 frames. Unlike previous
252 benchmarks that often resort to downsampling or temporal
253 truncation, our protocol ensures that every model is tested
254 at its maximum practical capacity. We utilize the officially
255



Figure 1. **Qualitative comparison on the PVIR benchmark.** Each row presents a specific instance removal scenario with its corresponding textual prompt. As a unified model, **UniVideo** demonstrates strong physics-awareness, successfully removing coupled side effects like ground shadows (e.g., rows b and c); however, it suffers from severe semantic hallucinations, occasionally generating unprompted artifacts to fill the void (e.g., the distorted figure generated in row d). **PISCO-Removal** consistently achieves clean erasure of both the target and its physical side effects while maintaining high background fidelity. **DiffuEraser** reliably masks the correct instance but frequently leaves unnatural residual shadows and spatial blurring (rows c and d). Finally, **CoCoCo** struggles significantly with instruction following, often leaving obvious “ghosting” silhouettes or failing to remove the object entirely (e.g., the white duck in row a).

256 released checkpoints for each method, ensuring that any ob-
 257 served performance gaps are intrinsic to the models’ archi-
 258 tectures rather than artifacts of suboptimal parameter tun-
 259 ing. This high-resolution, long-duration setup serves as
 260 a rigorous stress test for temporal-physical consistency in
 261 Video Instance Removal.

262 4.2. Human Evaluation Protocol

263 **Decoupled scoring principle.** All three dimensions are
 264 scored *independently*; a score in one dimension must not in-
 265 fluence another. Each dimension uses a 1–4 ordinal rubric,
 266 where higher is better.

(1) **Instruction Following (IF).** **Core question:** Does the
 267 edited video correctly satisfy the removal instruction? For
 268 this benchmark, judges check whether the specified target
 269 instance is removed, and whether visible remnants (edges,
 270 fragments, obvious traces) remain. 271

(2) **Rendering Quality (RQ).** **Core question:** Is the
 272 filled content visually plausible over space and time? 273
 274 Judges inspect naturalness, sharpness, temporal sta-
 275 bility (flicker/jitter), and physical plausibility of mo-
 276 tion/appearance.

277 **(3) Edit Exclusivity (EE). Core question:** Did the
278 model only perform the requested removal? Judges verify
279 that non-target regions preserve original content, including
280 background structures, lighting appearance, and unrelated
281 objects.

282 **Aggregation.** Let $s_d^{(i,r)} \in \{1, 2, 3, 4\}$ denote the score for
283 video i , rater r , and dimension $d \in \{\text{IF}, \text{RQ}, \text{EE}\}$. We first
284 average over raters:

$$\bar{s}_d^{(i)} = \frac{1}{R} \sum_{r=1}^R s_d^{(i,r)}. \quad (1)$$

285 Then report per-dimension dataset mean:

$$S_d = \frac{1}{N} \sum_{i=1}^N \bar{s}_d^{(i)}. \quad (2)$$

288 The aggregated benchmark score is

$$S_{\text{overall}} = \frac{1}{3} (S_{\text{IF}} + S_{\text{RQ}} + S_{\text{EE}}). \quad (3)$$

290 To ensure unbiased assessment, videos are assigned
291 to raters using a balanced, randomized sampling strategy.
292 Each video-model pair is evaluated by at least 2 indepen-
293 dent raters, and the presentation order of the four methods
294 is shuffled for every trial to eliminate model-specific or se-
295 quential bias.

296 4.3. Reporting Protocol

297 **Primary Metrics.** We report the mean scores for Instruc-
298 tion Following (IF), Rendering Quality (RQ), and Edit Ex-
299 clusivity (EE) across the entire benchmark. An Overall
300 Score is computed as the unweighted arithmetic mean of
301 these three dimensions, providing a holistic measure of in-
302 stance removal performance.

303 **Split-wise Analysis.** To isolate model robustness under
304 varying physical complexities, we additionally report per-
305 formance partitioned by the Simple and Hard subsets.
306 This granular reporting highlights the "performance decay"
307 models experience when transitioning from basic scenarios
308 to those with strong physical coupling (e.g., reflections and
309 wakes).

310 **Uncertainty and Significance.** To ensure the reliabil-
311 ity of our rankings, we report 95% Confidence Intervals
312 (CIs) for all primary metrics. For pairwise model com-
313 parisons, we employ Bootstrap Significance Testing (with
314 $N = 10,000$ iterations). Our analysis confirms that the
315 performance superiority of PISCO-Removal and UniVideo
316 is statistically significant ($p < 0.05$) across all dimensions,
317 particularly on the Hard subset where simpler baselines suf-
318 fer from severe physical artifacts.

5. Experiments 319

5.1. Experimental Setup 320

Benchmark setting. We evaluate all methods on the full
95-video benchmark and report: (1) overall scores, (2) split-
wise scores on Simple and Hard subsets, (3) qualitative fail-
ure analysis. Unless otherwise noted, all scores are human
ratings following Sec. 4.2. 321
322
323
324
325

Models. We benchmark CoCoCo [28], UniVideo [23],
DiffuEraser [11], and PISCO-Removal [3, 14]. 326
327

Implementation Details. All experiments are conducted
on a workstation with three NVIDIA A100 (80GB) GPUs.
To ensure a high-fidelity evaluation, we process all se-
quences at a native 720p (81 frames) resolution. The infer-
ence latency varies significantly across baselines: UniVideo
requires ~ 3.5 hours per video, while PISCO-Removal,
DiffuEraser, and CoCoCo complete in ~ 30 minutes. For
models with limited temporal receptive fields (e.g., Co-
CoCo), we apply a sliding window inference with a 4-
frame overlap to maintain coherence. All outputs are ex-
ported in a lossless format to avoid secondary compression
artifacts during human evaluation. 328
329
330
331
332
333
334
335
336
337
338
339

5.2. Main Results 340

Overall comparison. Tab. 3 summarizes the overall per-
formance. We observe a clear performance hierarchy:
PISCO-Removal and UniVideo consistently define the
state-of-the-art across all metrics, forming a high-fidelity
tier. In contrast, while DiffuEraser provides competitive
efficiency, it suffers from a significant "fidelity gap" com-
pared to the leaders. The most striking finding is the univer-
sal performance decay on the *Hard* subset, where even the
top-performing models struggle to maintain physical con-
sistency, highlighting the diagnostic value of our bench-
mark. 341
342
343
344
345
346
347
348
349
350
351

Instruction following. This dimension evaluates the
model's ability to ground textual instructions into pixel-
level removal, where we observe a clear trade-off between
semantic autonomy and execution reliability. UniVideo
and PISCO-Removal, as representative end-to-end archi-
tectures, demonstrate superior visual integration in the ma-
jority of cases; however, they exhibit occasional "ground-
ing drift" where the target subject is either ignored or only
partially erased. This suggests that while their latent se-
mantic alignment is powerful, it can occasionally fail to lo-
calize the instance accurately amidst cluttered backgrounds.
In contrast, DiffuEraser achieves the highest reliability in
complete removal across the benchmark. Since its pipeline
is explicitly constrained by the input mask, it bypasses the
352
353
354
355
356
357
358
359
360
361
362
363
364
365

Table 2. **Decoupled Human Evaluation Rubric.** Each of the three dimensions—Instruction Following (IF), Rendering Quality (RQ), and Edit Exclusivity (EE)—is scored independently on a 1–4 scale, enabling granular diagnosis of model failure modes.

Score	Instruction Following	Rendering Quality	Edit Exclusivity
Score 1	<i>Target not removed or unrelated edit</i>	<i>Severe artifacts; unusable output</i>	<i>Uncontrolled edits across scene</i>
Score 2	<i>Partial removal with obvious residual traces</i>	<i>Obvious distortion/flicker; poor temporal consistency</i>	<i>Multiple non-target regions altered</i>
Score 3	<i>Target mostly removed with minor artifacts</i>	<i>Moderate quality degradation but viewable result</i>	<i>Minor non-target changes but structure mostly preserved</i>
Score 4	<i>Target precisely removed with no obvious residuals</i>	<i>High visual quality and stable temporal consistency</i>	<i>Non-target regions preserved with only negligible differences</i>

Table 3. Performance comparison on the comprehensive benchmark (95 videos), including overall results and the Simple/Hard split breakdown. Metric scores range from 1 to 4, where higher values indicate superior performance. **Dark green** and **light green** denote the best and second-best results within each subset, respectively.

Subset	Method	Instruction Following \uparrow	Rendering Quality \uparrow	Edit Exclusivity \uparrow	Overall \uparrow
Overall	CoCoCo	1.60	1.84	3.07	2.17
	UniVideo	3.06	3.45	3.53	3.35
	DiffuEraser	2.89	2.63	3.52	3.01
	PISCO-Removal	3.62	3.28	3.58	3.49
Simple	CoCoCo	1.75	1.98	3.33	2.35
	UniVideo	3.21	3.34	3.53	3.36
	DiffuEraser	2.73	2.73	3.73	3.06
	PISCO-Removal	3.75	3.32	3.57	3.55
Hard	CoCoCo	1.52	1.76	2.92	2.07
	UniVideo	2.96	3.53	3.53	3.34
	DiffuEraser	3.00	2.56	3.38	2.98
	PISCO-Removal	3.45	3.23	3.59	3.42

366 semantic localization errors inherent in end-to-end mod- 384
 367 els, ensuring the designated regions are always processed. 385
 368 Notably, **CoCoCo** consistently fails this dimension; as 386
 369 a general-purpose inpainting model, it lacks the specialized 387
 370 inductive bias for large-scale instance removal, often pro- 388
 371 ducing "ghosting" artifacts that retain the original object's 389
 372 silhouette or failing to initiate the removal command alto- 390
 373 gether in favor of local texture synthesis. 391

374 **Rendering quality.** This dimension evaluates visual fi- 392
 375 delity and temporal stability, where we observe a stark con- 393
 376 trast in spatial resolution and coherence. **UniVideo** and 394
 377 **PISCO-Removal** consistently produce the most visually 395
 378 pleasing results, maintaining high-frequency textures that 396
 379 blend seamlessly with the original background. While they 397
 380 exhibit occasional flickering in complex dynamic scenes, 398
 381 their overall rendering remains stable at 720p. In con- 399
 382 trast, **DiffuEraser** suffers from pervasive *spatial blurring* 400
 383 within the inpainted regions. Despite its ability to reliably 401

remove the target, the synthesized textures often lack the 384
 sharpness of the surrounding environment, creating a no- 385
 ticeable "patchwork" effect that disrupts the scene's visual 386
 harmony. **CoCoCo** performs the worst in this category, fre- 387
 quently generating severe "ghosting" artifacts—where rem- 388
 nants of the original object reappear as semi-transparent 389
 textures—and significant temporal flickering. These fail- 390
 ures suggest that while general inpainting models can fill 391
 small holes, they lack the structural priors necessary to re- 392
 construct large-scale, 81-frame backgrounds with the requi- 393
 site physical and temporal plausibility. 394

Edit exclusivity. This dimension assesses the models' 395
 precision in localized editing, specifically focusing on 396
 whether the transformation "leaks" into non-target back- 397
 ground regions. Overall, all four baselines demonstrate a 398
 respectable ability to preserve the surrounding scene con- 399
 text. **UniVideo** and **PISCO-Removal** exhibit high spatial 400
 fidelity, maintaining the original pixel values of the static 401

402 background with minimal drift. However, **DiffuEraser** oc- 452
403 casionally suffers from *blurring leakage*, where the spa- 453
404 tial smoothing intended for the erased region inadvertently 454
405 spreads beyond the mask boundaries into the neighboring 455
406 textures. This creates a subtle but perceptible halo of re- 456
407 duced sharpness in the background, a phenomenon that is 457
408 particularly visible at 720p. For the *Hard* subset involving 458
409 complex reflections, we observe that models often strug- 459
410 gle to disentangle the target instance from its environmental 452
411 "side effects," sometimes leading to unintended modifica- 453
412 tions of the global lighting or shadow maps in areas that 454
413 should remain untouched. 455

414 5.3. Simple vs. Hard Split Analysis

415 The split-wise analysis isolates model robustness against 452
416 physics-coupled interactions. As detailed in Tab. 3, we ob- 453
417 serve a noticeable performance degradation across multiple 454
418 dimensions when transitioning from the *Simple* to the *Hard* 455
419 subset. Looking beyond the aggregate scores, the multi- 456
420 dimensional breakdown reveals that the most significant 457
421 drops typically occur in **Instruction Following** and **Ren- 458
422 dering Quality**. For instance, top-tier models like **PISCO- 459
423 Removal** experience a drop in IF from 3.75 on the Simple 452
424 set to 3.45 on the Hard set, illustrating the increased diffi- 453
425 culty of executing clean removals amidst complex physical 454
426 constraints. While all four baselines can reliably erase an 455
427 isolated object against a static background in the *Simple* set, 456
428 they frequently struggle to maintain structural consistency 457
429 when the task scales in physical complexity. 458

430 The interaction types involving **specular reflections** and 452
431 **dynamic wakes/fluid ripples** prove to be the most chal- 453
432 lenging. In many *Hard* cases, even when the primary 454
433 instance is successfully erased by leaders like **UniVideo** 455
434 or **PISCO-Removal**, its corresponding physical "side ef- 456
435 fects"—such as a moving shadow on a textured wall or a 457
436 mirror reflection on a car’s surface—remain stubbornly vis- 458
437 ible. This severely impacts the **Edit Exclusivity** scores for 459
438 methods like **DiffuEraser**, which drops from 3.73 (Simple) 452
439 to 3.38 (Hard). Its local blurring strategy fails to properly 453
440 propagate the background’s global illumination, leading 454
441 to physically implausible "ghost reflections." These find- 455
442 ings underscore that current Video Instance Removal (VIR) 456
443 models generally treat removal as a 2D texture-filling task 457
444 rather than a 3D-aware scene reconstruction, highlighting a 458
445 critical direction for future physics-augmented research. 459

446 5.4. Cross-Metric Trade-off Analysis

447 To analyze whether methods sacrifice one dimension for 452
448 another, we examine the pairwise relationships in Fig. 2. 453
449 This visualization exposes the inherent tension between a 454
450 model’s "semantic fluidity" and its adherence to local con- 455
451 straints. 456

High-Fidelity Clustering. **PISCO-Removal** and **Uni- 452
453 Video** (blue and green) exhibit tight clustering in the (4,4) 454
455 quadrants of the IF-RQ and IF-EE plots, defining a high- 456
457 fidelity tier that balances task completion with background 458
459 preservation. However, a clear shift toward lower RQ scores 452
is evident for the **Hard** subset (\times), confirming that complex 453
physical coupling remains the primary performance bottle- 454
neck even for top-tier models. 455

Reliability vs. Precision. **DiffuEraser** (red) achieves 452
453 high IF scores (Score 3–4) due to its mask-guided nature, 454
455 yet shows significant dispersion on the RQ and EE axes. 456
457 This reflects a trade-off where reliable object removal often 458
459 introduces spatial blurring or "halo" artifacts in non-target 452
regions. Conversely, **CoCoCo** (purple) is heavily skewed 453
454 toward low IF scores (Score 1–2), often failing to execute 455
456 the removal command while maintaining high EE scores 457
458 simply by leaving the scene unedited. 459

Metric Correlations. Statistical analysis reveals that IF 452
453 and RQ are moderately correlated ($r \approx 0.65$), suggesting 454
455 that models with better instruction grounding typically syn- 456
457 thesize more plausible textures. However, the weak cor- 458
459 relation between RQ and EE in lower-performing models 452
underscores the necessity of our decoupled protocol: vi- 453
454 sual plausibility alone does not guarantee that the rest of 454
455 the scene remains untouched. 456
457

458 5.5. Discussion

Why decoupled evaluation matters. Traditional VIR 452
453 evaluation often relies on a single visual plausibility score, 454
455 which masks critical failure modes. Our three-axis proto- 456
457 col reveals that a model can score highly on *instruction fol- 458
459 lowing* while simultaneously failing *edit exclusivity*, or pro- 452
453 duce visually plausible textures while completely missing 454
455 the target removal. By decoupling these dimensions, we ex- 456
457 pose the inherent trade-offs in current architectures—such 458
459 as the balance between the strict spatial constraints of mask- 452
453 guided models like **DiffuEraser** and the semantic fluidity of 454
455 end-to-end generative models like **UniVideo**. This granular 456
457 mapping is essential for diagnosing whether a model’s fail- 458
459 ure stems from poor instruction grounding, low rendering 452
453 fidelity, or unintended scene drift. 454
455

Current limitations. Current benchmark evaluations re- 452
453 main heavily human-centered, making them both cost- 454
455 sensitive and difficult to scale. While we provide a compre- 456
457 hensive human study, the development of robust automatic 458
459 metrics for physics-aware side effects—such as detecting 452
453 residual reflections, lingering shadows, or fluid inconsisten- 454
455 cies—remains highly challenging. Existing feature-space 456
457 proxies fail to capture this high-level physical logic, indicat- 458
459

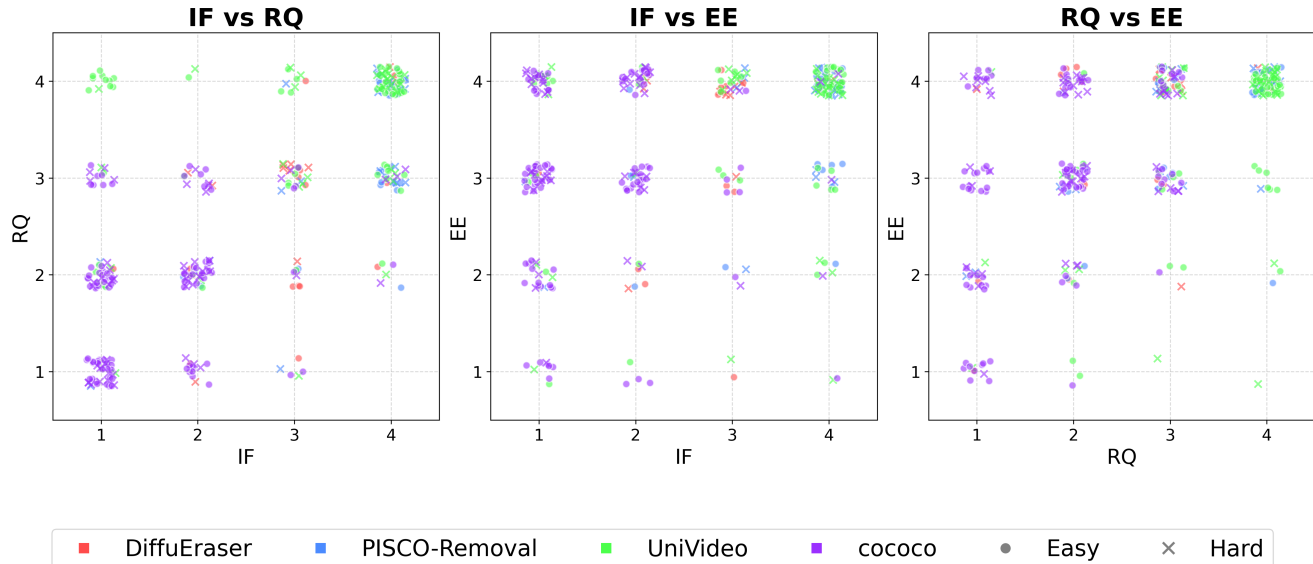


Figure 2. **Cross-metric trade-off analysis across Instruction Following (IF), Rendering Quality (RQ), and Edit Exclusivity (EE).** Each point represents a single video-model pair. Markers distinguish between the *Simple* (●) and *Hard* (×) subsets. The tight clustering in the top-right quadrant for **PISCO-Removal** and **UniVideo** indicates high-fidelity, balanced performance. Conversely, the wide dispersion of **DiffuEraser** and **CoCoCo** reveals inherent trade-offs between reliable target removal and background preservation.

500 ing an urgent need for automated, physics-aware evaluation
501 models.

502 **Future directions.** Our empirical results highlight a clear
503 paradigm shift: large-scale generative models significantly
504 outperform smaller, traditional inpainting networks. Fur-
505 thermore, within the regime of large models, specialized ar-
506 chitectures fine-tuned for precise editing demonstrate no-
507 ticeable superiority over unified, general-purpose video
508 models in handling strict spatial constraints. This sug-
509 gests that while foundational priors are necessary, they are
510 not sufficient for instance-level physical accuracy. Con-
511 sequently, future advancements in video instance removal
512 should focus on two key pillars: the continued scaling of
513 large video foundation models, and the rigorous curation
514 of high-quality, domain-specific datasets designed to inject
515 precise physical and spatial constraints into these models.

516 **Benchmark extension roadmap.** As an evolving plat-
517 form, the PVIR roadmap includes expanding the dataset to
518 encompass greater category diversity, extended sequences
519 (e.g., 10+ seconds) to stress-test long-term temporal con-
520 sistency, and more complex physical phenomena such as
521 multi-object interactions and fluid/smoke dynamics.

522 6. Conclusion

523 In this work, we introduced the Physics-Aware Video In-
524 stance Removal (PVIR) benchmark, a dedicated evaluation

525 infrastructure designed to bridge the gap between visual
526 plausibility and physical consistency in video editing. Our
527 benchmark contributes three key pillars: a high-quality 95-
528 video dataset with dense annotations, a Simple/Hard dif-
529 ficulty split driven by complex physical interactions, and
530 a decoupled evaluation protocol spanning instruction fol-
531 lowing, rendering quality, and edit exclusivity. By bench-
532 marking four representative models under a unified, high-
533 resolution inference setup, we revealed a significant perfor-
534 mance hierarchy and a universal “physics blindness” in cur-
535 rent state-of-the-art methods.

536 Crucially, our analysis demonstrates that while exist-
537 ing models excel in static background synthesis, they suffer
538 from severe performance when encountering optical reflec-
539 tions, shadows, and fluid interactions in our Hard sub-
540 set. This finding underscores that current video instance
541 removal is still predominantly treated as a 2D texture-filling
542 task rather than a 3D-aware physical reconstruction. We
543 offer PVIR as a rigorous yardstick to encourage the com-
544 munity to move beyond surface-level aesthetic metrics and
545 toward physically grounded video intelligence.

546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601**References**

- [1] Yinan Chen, Jiangning Zhang, Teng Hu, Yuxiang Zeng, Zhucun Xue, Qingdong He, Chengjie Wang, Yong Liu, Xiaobin Hu, and Shuicheng Yan. Ivebench: Modern benchmark suite for instruction-guided video editing assessment. *arXiv preprint arXiv:2510.11647*, 2025. 2
- [2] Chen Gao, Ayush Saraf, Jia-Bin Huang, and Johannes Kopf. Fgvc: Flow-guided video completion. In *CVPR*, 2019. 2
- [3] Xiangbo Gao, Renjie Li, Xinghao Chen, Yuheng Wu, Suofei Feng, Qing Yin, and Zhengzhong Tu. Pisco: Precise video instance insertion with sparse control. *arXiv preprint arXiv:2602.08277*, 2026. 1, 3, 5
- [4] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arXiv:2307.10373*, 2023. 2
- [5] 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 *NeurIPS*, 2017. 2
- [6] Jonathan Ho, William Chan, and Pieter Abbeel. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022. 2
- [7] Weizhe Huang, Xiaofeng Liu, Yifan Wang, Xin Li, et al. Vbench++: Comprehensive and versatile benchmark for video understanding and generation. *arXiv preprint*, 2024. 2
- [8] 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. 3
- [9] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 3
- [10] Xuan Li, Yujie Wang, Chuhang Zhang, Bin Zhao, and Ying Shan. Video-p2p: Video editing with cross-attention control. *arXiv preprint arXiv:2303.04761*, 2023. 2
- [11] Xiaowen Li, Haolan Xue, Peiran Ren, and Liefeng Bo. Diffuseraser: A diffusion model for video inpainting. *arXiv preprint arXiv:2501.10018*, 2025. 1, 2, 3, 5
- [12] Zhen Li, Cheng Xie, Weidi Zhang, Yebin Liu, Qi Tian, and Ying Shan. Towards an end-to-end framework for flow-guided video inpainting. In *CVPR*, 2022. 2
- [13] Rui Liu, Hanming Deng, Yangyi Huang, Xiaoyu Shi, Lewei Lu, Wenxiu Sun, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. Fuseformer: Fusing fine-grained information in transformers for video inpainting. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14040–14049, 2021. 2
- [14] Chenxuan Miao, Yutong Feng, Jianshu Zeng, Zixiang Gao, Hantang Liu, Yunfeng Yan, Donglian Qi, Xi Chen, Bin Wang, and Hengshuang Zhao. Rose: Remove objects with side effects in videos. *arXiv preprint arXiv:2508.18633*, 2025. 1, 2, 3, 5
- [15] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 2
- [16] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 724–732, 2016. 3
- [17] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, 2016. 2
- [18] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2
- [19] Alexandros Stergiou and Ronald Poppe. Adapool: Exponential adaptive pooling for information-retaining downsampling. 2021. 2
- [20] Roman Suvorov, Ekaterina Logacheva, Anton Mashikhin, Mikhail Melnikov, Mikhail Kaigorodov, Sergey Yudin, Denis Davydov, Anastasia Molchanova, Artem Malkov, Alexander Ilin, et al. Resolution-robust large mask inpainting with fourier convolutions. In *WACV*, 2022. 2
- [21] Thomas Unterthiner, Bernhard Nessler, Guenter Klambauer, Martin Heusel, Hubert Ramsauer, and Sepp Hochreiter. Towards accurate generative models of video: A new metric and challenges. *arXiv preprint arXiv:1812.01717*, 2018. 2
- [22] Jiayu Wang, Yicheng Zhang, Kai Liu, Xin Sun, Lijuan Ye, Yunchao Wei, et al. Vbench: Comprehensive benchmark suite for video generative models. *arXiv preprint arXiv:2311.17982*, 2024. 2
- [23] Cong Wei, Quande Liu, Zixuan Ye, Qiulin Wang, Xintao Wang, Pengfei Wan, Kun Gai, and Wenhui Chen. Univideo: Unified understanding, generation, and editing for videos. *arXiv preprint arXiv:2510.08377*, 2025. 1, 2, 3, 5
- [24] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Z. Lei, Yuchao Gu, Bolei Huo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. *arXiv preprint arXiv:2212.11565*, 2023. 2
- [25] Ning Xu, Linjie Yang, Yuchen Fan, Dingakang Yue, Yuchen Liang, James Yang, and Thomas Huang. Youtube-vos: Sequence-to-sequence video object segmentation. In *ECCV*, 2018. 2
- [26] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas Huang. Free-form image inpainting with gated convolution. In *ICCV*, 2019. 2
- [27] 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
- [28] 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, 3, 5

658 A. Supplementary Material**659 A.1. Evaluation Guide for Human Raters**

660 This section provides the detailed 1–4 criteria used by annotators. All dimensions are strictly independent.

661 **Instruction Following (IF).** **Core question:** Does the edited video correctly satisfy the removal instruction?

- 662 • **4:** Perfect adherence. Target instance is removed cleanly with no obvious remnants.
- 663 • **3:** High adherence. Target removed with minor residual imperfections.
- 664 • **2:** Low adherence. Target partly removed, or obvious residual traces remain.
- 665 • **1:** Failure. Target remains or output is largely unrelated to the instruction.

666 **Rendering Quality (RQ).** **Core question:** Is the filled content visually plausible and temporally stable?

- 667 • **4:** Excellent quality, negligible artifacts, temporally stable.
- 668 • **3:** Acceptable quality, moderate artifacts but viewable.
- 669 • **2:** Poor quality, obvious distortion/flicker affecting readability.
- 670 • **1:** Unusable quality, severe artifacts or physical implausibility.

671 **Edit Exclusivity (EE).** **Core question:** Are non-target regions preserved?

- 672 • **4:** Non-target regions are effectively unchanged.
- 673 • **3:** Minor non-target drift but scene remains mostly preserved.
- 674 • **2:** Noticeable over-editing in multiple non-target areas.
- 675 • **1:** Major uncontrolled edits; scene is substantially altered.

676 A.2. Annotation and Prompt Templates

677 **[TODO: Placeholder: annotation UI screenshots and mask quality checklist.] [TODO: Placeholder: prompt template**
678 **examples and ambiguity handling rules.]**

679 A.3. Additional Qualitative Cases

680 **[TODO: Placeholder: extended qualitative comparison for each method.] [TODO: Placeholder: side-effect-focused**
681 **examples (mirror, reflection, lighting).]**

682 A.4. Extended Experimental Details

683 **[TODO: Placeholder: full inference hyper-parameter table and environment details.] [TODO: Placeholder: per-**
684 **model runtime/memory statistics and sensitivity analysis.]**