

The Evaluation Imperative for Video Generative Models: A Survey on Metrics, Benchmarks, and Trustworthiness

Anonymous CVPR submission

Paper ID 16

Abstract

001 *The rapid evolution of video generative models has shifted*
 002 *the research focus from basic synthesis to high-fidelity,*
 003 *physically grounded content. However, traditional evalu-*
 004 *ation frameworks remain largely insufficient, often failing*
 005 *to capture the spatiotemporal consistency and semantic*
 006 *nuances required for verifiable quality. This survey provides*
 007 *a critical analysis of the AI-generated video evaluation*
 008 *(AIGVE) field, categorizing methodologies into metric-*
 009 *based, human-involved, and model-centered paradigms. We*
 010 *examine how architectural shifts primarily toward diffusion*
 011 *transformers and autoregressive frameworks have redefined*
 012 *evaluation requirements. We trace the evolution of metrics*
 013 *from frame-level heuristics to perceptually grounded spa-*
 014 *tiotemporal measures and analyze benchmarks spanning*
 015 *text-to-video, specialized conditional generation, and long-*
 016 *form storytelling. The review also addresses the critical*
 017 *dimensions of human preference alignment, physical world*
 018 *simulation, and safety-aware assessment, while exploring*
 019 *the impact of generative priors on downstream applica-*
 020 *tions like video editing and compression. By identifying current*
 021 *limitations, we propose a path toward unified and trustwor-*
 022 *thy evaluation frameworks for the next generation of video*
 023 *foundation models.*

024 1. Introduction

025 Establishing rigorous evaluation is the cornerstone of reli-
 026 able video generation.

027 1.1. From Visual Synthesis to Verifiable Quality

028 The rapid evolution of video generative models has shifted
 029 the research frontier from basic visual synthesis to the pur-
 030 suit of high-fidelity, temporally consistent, and physically
 031 grounded content. However, existing evaluation frame-
 032 works such as VBench and EvalCrafter often fall short in
 033 assessing dynamic scenarios, primarily due to their reliance
 034 on subject-centric, static prompts and frame-level metrics

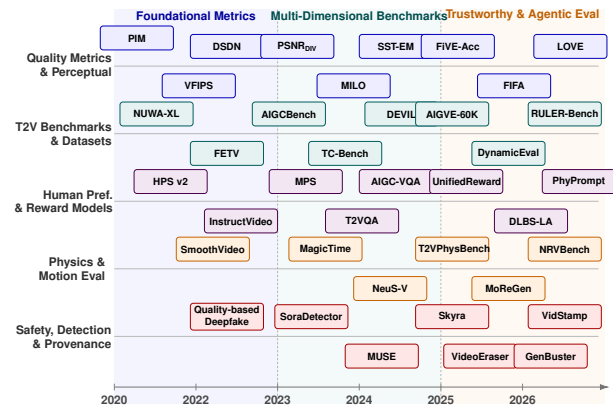


Figure 1. Chronological development timeline of the video generative model evaluation landscape (2020–2026), organized across five research lanes: quality metrics, benchmarks, human preference alignment, physics-aware evaluation, and safety/provenance. Three evolutionary phases are identified: foundational per-frame metrics, multi-dimensional benchmark expansion, and the current shift toward trustworthy, agentic evaluation frameworks.

[1]. These benchmarks frequently prioritize aesthetic fidelity over cinematic camera motion and temporal causality, leading to technical failures in consistency metrics when faced with complex occlusions or disocclusions [1].

Furthermore, traditional automated metrics exhibit a significant disconnect from human perception. Standard measures like Fréchet Inception Distance (FID) and Fréchet Video Distance (FVD) are loosely coupled with subjective quality, while CLIP-based scores often fail to capture temporal nuances [16, 22, 32]. This misalignment is particularly evident in temporal compositionality tasks; for instance, state-of-the-art models achieve less than 20% success on the TC-Bench benchmark [13]. The imperative for verifiable quality thus necessitates a move beyond simple image-based heuristics toward metrics that can rigorously evaluate physical consistency, complex motion dynamics, and long-term temporal coherence [56, 64].

Table 1. Key papers advancing evaluation across metrics, benchmarks, alignment, and trustworthiness.

Category	Paper	Main Contribution
Evaluation Metrics	PSNR _{Div} [9]	Enhances PSNR with motion-field divergence weighting, achieving +0.09 PLCC over FIoLPIPS at 2.5× speed.
	VFIPS [18]	Swin Transformer-based spatiotemporal perceptual metric for video frame interpolation quality.
	NeuS-V [41]	Neuro-symbolic text-video alignment via Temporal Logic specifications; 5× higher human correlation.
	FIFA [21]	Spatio-Temporal Semantic Dependency Graph for unified faithfulness evaluation of T2V and V2T.
	DEVIL [26]	Multi-granularity dynamics protocol achieving >90% Pearson correlation with human ratings.
T2V Benchmarks	AIGVE-60K [47]	Largest AIGV evaluation dataset: 58,500 videos, 2.6M human ratings, 30 T2V models, 20 task dims.
	DynamicEval [1]	45K human annotations for dynamic camera motion; proposes MS-Debias and Tracker-FG metrics.
	TC-Bench [13]	Temporal compositionality benchmark; reveals <20% success rate for state-of-the-art generators.
	RULER-Bench [17]	40 cognitive reasoning tasks across 6 categories; rule coherence scores often below 50%.
	T2VPhysBench [14]	First-principles physics benchmark evaluating 12 core physical laws via human judgment.
	LoCoT2V-Bench [65]	Long-form generation benchmark with hierarchical metadata for character and camera consistency.
Human Alignment	HPS v2 [53]	Fine-tunes CLIP on 798K human preference choices; superior generalization across distributions.
	UnifiedReward [49]	First unified reward model supporting pairwise ranking & pointwise scoring for image and video.
	DLBS-LA [36]	Inference-time beam search with VLM-calibrated rewards; training-free alignment improvement.
	PhyPrompt [52]	RL-based prompt refinement with GRPO; +8.6pp joint success on VideoPhy2 with zero-shot transfer.
Safety & Trust	Skyra [24]	MLLM for grounded artifact reasoning with ViF-CoT-4K dataset and two-stage SFT+RL training.
	GenBuster-200K [50]	200K+ video forensics dataset with three-track benchmark and MLLM-as-Judge protocol.
	VideoEraser [55]	Training-free concept erasure achieving 46% average reduction in undesirable content via SPEA+ARNG.
	SafeVid [48]	DPO-based safety alignment for VLMs; up to 42.39% safety improvement on SafeVidBench.
	VidStamp [45]	Temporally-aware watermarking: 48 bits/frame, 0.96 tamper localization accuracy, robust to distortions.

052 1.2. Scope, Taxonomy, and Organization

053 This review provides a structured overview of the AI-
054 generated video evaluation (AIGVE) landscape, catego-
055 rized into three primary paradigms: metric-based, human-
056 involved, and model-centered evaluations [31]. Metric-
057 based approaches utilize quantitative formulas to assess
058 objective dimensions such as visual fidelity and temporal
059 consistency [20]. Human-involved protocols validate these
060 metrics through user studies and preference tests, while
061 emerging model-centered paradigms leverage Multimodal
062 Large Language Models (MLLMs) and agent-based sys-
063 tems to automate perceptually grounded assessments [56,
064 67]. A summary of key papers across these dimensions is
065 provided in Tab. 1.

066 The unique challenges of video evaluation stem from
067 the dual consideration of spatial and temporal dynam-
068 ics. Unlike static image assessment, video evaluation
069 must account for motion direction, event ordering, and
070 temporal compositionality, including attribute transitions
071 and object relationship shifts over time [13, 32]. We
072 categorize these challenges across four primary quality
073 dimensions: visual fidelity, temporal consistency, semantic

alignment, and physical consistency [11, 64]. Furthermore,
we examine the taxonomy of applications based on input
modalities, including text-to-video (T2V), image-to-video
(I2V), video-to-video (V2V), and audio-conditioned gen-
eration [35]. The remainder of this paper is organized as
follows: Section 2 discusses technical foundations; Section
3 details evaluation metrics; Sections 4 and 5 explore
general-purpose and specialized benchmarks; Section 6
examines human preference alignment; Section 7 addresses
safety and trustworthiness; and Section 8 concludes with
open challenges and future directions.

2. Foundations: Video Gen. Model Taxonomy

Architectural innovations drive the spatiotemporal com-
plexity of modern video models.

2.1. Diffusion Transformer Architectures

Diffusion Transformers (DiT) have emerged as the dom-
inant backbone for high-fidelity video synthesis, lever-
aging scalable attention mechanisms to model complex
spatiotemporal dependencies. CogVideoX [57] exem-
plifies this trend by utilizing a 3D causal VAE and expert
transformer blocks to facilitate deep text-video fusion. To
optimize the tradeoff between semantic fidelity and compu-
tational cost, FlashVideo [63] employs a two-stage frame-
work that decouples motion alignment from high-resolution
detail refinement, strategically allocating model capacity
across stages. Efficiency remains a critical focus in DiT
design, with methods like ADAPTOR [38] introducing
adaptive token reduction to exploit temporal redundancy.
Similarly, SNED [25] utilizes neural architecture search to
address the cost-resolution tradeoff. For temporal stability,
the Dual-Stream Diffusion Net (DSDN) [27] separates
content and motion into distinct streams, mitigating flicker
through cross-transformer synchronization. These archi-
tectural advancements necessitate evaluation frameworks
that can specifically target motion-detail decoupling and
adaptive compute efficiency.

2.2. Autoregressive Architectures

Autoregressive (AR) models leverage the predictive power
of Large Language Models (LLMs) to synthesize video
as a sequence of discrete tokens. Lumos-1 [60] adapts
this paradigm by introducing spatiotemporal position em-
beddings and discrete diffusion forcing to address frame-
wise loss imbalance. This approach achieves performance
comparable to state-of-the-art diffusion models on bench-
marks like VBench and GenEval. Further innovations
in AR frameworks focus on long-range consistency and
multimodal integration. ART•V [51] decomposes video
synthesis into a sequential process, utilizing anchored con-
ditioning to handle extended sequences. For interactive

123	applications, Midas [7] integrates audio, pose, and text encodings into an LLM backbone, utilizing high-ratio spatial compression to maintain real-time efficiency. These models highlight the importance of evaluating long-term coherence and cross-modal synchronization in token-based synthesis.	
124		
125		
126		
127		
128	3. Evaluation Metrics for Video Generation	
129	Quantitative measures bridge the gap between pixel-level fidelity and human perception.	
130		
131	3.1. Perceptual and Full-Reference Quality Metrics	
132	Traditional full-reference metrics such as PSNR and SSIM often fail to capture video-specific artifacts like ghosting and temporal flickering, as they operate on a per-frame basis without considering motion coherence. To address these limitations, PSNRDIV [9] enhances standard PSNR by weighting the mean squared error with motion-field divergence, effectively identifying temporal singularities. Similarly, VFIPS [18] leverages Swin Transformer blocks to extract spatiotemporal features, significantly outperforming image-based perceptual metrics like LPIPS in correlating with human judgment on video frame interpolation tasks.	
133		
134		
135		
136		
137		
138		
139		
140		
141		
142		
143	Efficiency and unsupervised learning are also key research directions. MILO [69] introduces a lightweight multiscale CNN architecture that utilizes pseudo-MOS supervision from an ensemble of existing metrics, enabling training without large-scale human-labeled datasets. From an information-theoretic perspective, the Perceptual Information Metric (PIM) [3] employs unsupervised contrastive learning on video frame pairs. Grounded in human visual system physiology, PIM extracts temporally persistent features that demonstrate competitive performance with supervised metrics while exhibiting superior robustness to geometric distortions.	
144		
145		
146		
147		
148		
149		
150		
151		
152		
153		
154		
155	3.2. Temporal Consistency and Motion Metrics	
156	Assessing temporal consistency requires metrics that can distinguish between intentional motion and technical artifacts. DynamicEval [1] identifies critical failures in standard motion smoothness metrics, which often yield large errors near occlusions caused by camera movement. To correct this, they propose MS-Debias, which utilizes object segmentation masks to isolate problematic regions, and Tracker-FG, which tracks points within object instances to assess fidelity under dynamic camera motion. These methods significantly improve correlation with human preferences by focusing on pixel-level temporal smoothness rather than global features.	
157		
158		
159		
160		
161		
162		
163		
164		
165		
166		
167		
168	Comprehensive dynamics evaluation is further addressed by the DEVIL protocol [26], which measures video dynamics across multiple temporal granularities, inter-frame, inter-segment, and video-level. DEVIL defines metrics for	
169		
170		
171		
	dynamics range, controllability, and dynamics-based quality, achieving a Pearson correlation exceeding 90% with human ratings. For adjacent frame transitions, SmoothVideo [37] introduces the Video Latent Score (VL score), which processes frames through a fine-grained autoencoder to ensure temporal order and capture granular details. Additionally, the SST-EM framework [4] integrates semantic extraction, object tracking, and temporal consistency for video editing assessment, with component weights derived from regression analysis against human evaluations.	172 173 174 175 176 177 178 179 180 181
	3.3. Semantic Alignment and Text-Video Metrics	182
	Achieving high semantic alignment is a core challenge in text-to-video generation, yet traditional metrics like CLIP-Score and FVD often correlate poorly with human judgment due to their lack of deep vision-language understanding and sensitivity to sample size [32]. To address this, the FETV benchmark introduces UMTScore and FVD-UMT, which leverage the Unbiased Multimodal Transformer (UMT) to provide more consistent rankings and stronger correlation with human perception. For formal verification, NeuS-V [41] converts text prompts into Temporal Logic (TL) specifications and translates videos into automaton representations, demonstrating a correlation with human evaluations that is over five times higher than existing heuristics.	183 184 185 186 187 188 189 190 191 192 193 194 195
	Emerging metrics also utilize Large Multimodal Models (LMMs) for diagnostic assessment. FiVE-Acc [19] leverages VLMs to perform automated semantic querying, using Yes/No and multiple-choice questions to verify fine-grained object-level editing success. Similarly, the FIFA framework [21] models semantic dependencies via a Spatio-Temporal Semantic Dependency Graph (STSDG). By structuring descriptive facts as a directed acyclic graph, FIFA can systematically assess hallucinations through VideoQA, ensuring that prerequisite entities are verified before their dependent attributes, thereby providing a more faithful evaluation of text-video correspondence.	196 197 198 199 200 201 202 203 204 205 206 207
	3.4. Human Preference and Multidimensional Assessment	208 209
	Aligning automated metrics with human subjective preferences is essential for evaluating the real-world impact of generative models. The Human Preference Score v2 (HPS v2) [53] fine-tunes CLIP on a large-scale dataset of human choices, demonstrating superior generalization across diverse image distributions compared to prior metrics. Extending this to multiple dimensions, the Multidimensional Preference Score (MPS) [62] incorporates a preference condition module to capture aesthetics, semantic alignment, and detail quality separately, ensuring that the model attends to visual features relevant to specific human criteria.	210 211 212 213 214 215 216 217 218 219 220 221
	For holistic video quality assessment, T2VQA [22]	222

223 utilizes a transformer-based architecture that combines a
224 BLIP-based alignment encoder with a Video Swin Trans-
225 former backbone. By leveraging cross-attention fusion
226 and a Large Language Model (Vicuna v1.5) for score
227 regression, T2VQA achieves significant improvements in
228 subjective alignment on the T2VQA-DB dataset. Simi-
229 larly, AIGC-VQA [33] integrates three functional branches
230 for technical quality, aesthetics, and alignment. Utilizing
231 a Spatial-Temporal Adapter (ST-Adapter) for image-to-
232 video knowledge transfer and a divide-and-conquer training
233 strategy, AIGC-VQA achieves state-of-the-art performance,
234 providing predictions that are more closely aligned with
235 human Mean Opinion Scores (MOS) than single-branch
236 baselines.

237 4. Text-to-Video Benchmarks

238 Standardized benchmarks provide the testing grounds for
239 general-purpose video synthesis.

240 4.1. General-Purpose T2V Evaluation Benchmarks

241 The evaluation of text-to-video (T2V) models has transi-
242 tioned from simple visual quality metrics to comprehensive,
243 multi-dimensional frameworks, as illustrated in the unified
244 evaluation pipeline (Fig. 2). AIGVE-60K [47] represents
245 the largest effort to date, comprising 58,500 videos and
246 2.6 million human ratings. It introduces the LOVE met-
247 ric, which leverages a Large Multimodal Model (LMM)
248 architecture to disentangle perceptual quality, text-video
249 correspondence, and task-specific accuracy. Similarly,
250 AIGCBench [11] employs 11 metrics across four dimen-
251 sions, control-video alignment, motion effects, temporal
252 consistency, and video quality, to provide a unified assess-
253 ment of image-to-video and text-to-video generation.

254 To improve alignment with human perception, recent
255 benchmarks have integrated advanced reasoning capabili-
256 ties. Video-Bench [16] utilizes MLLMs with a Chain-of-
257 Query (CoQ) technique and few-shot scoring to achieve
258 superior correlation with human preferences. Addressing
259 the limitations of model-level aggregated scores, Dynam-
260 icEval [1] focuses on video-level evaluation through 45,000
261 human annotations, specifically targeting background and
262 foreground consistency in dynamic scenes. Furthermore,
263 agent-based systems like VideoGen-Eval [56] integrate
264 LLM-based content structuring with MLLM judgment and
265 specialized patch tools to evaluate temporal-dense dimen-
266 sions across thousands of videos from both commercial and
267 open-source models.

268 4.2. Dynamic Motion and Camera Benchmarks

269 Assessing motion dynamics and camera control is a critical
270 challenge, as many existing benchmarks exhibit a bias
271 toward static or low-dynamic content. DynamicEval [1]
272 addresses this by curating 100 prompts focused on dynamic

camera motion and operationalizing background and fore-
ground consistency through debiased optical flow and point
tracking. Similarly, the DIVE benchmark [30] mitigates the
low-dynamic bias in image-to-video (I2V) frameworks by
utilizing GPT-4o to quantify dynamic range and controlla-
bility, penalizing static outputs that might otherwise achieve
high quality scores.

Temporal granularity is further explored in the DEVIL
protocol [26], which defines dynamics scores across inter-
frame, inter-segment, and video-level scales. By establish-
ing a benchmark of 50,000 prompts across five dynamics
grades, DEVIL enables a more nuanced evaluation of mo-
tion complexity versus visual quality. For multi-shot sce-
narios, the TalkCuts dataset [6] provides a comprehensive
taxonomy of camera shot categories (e.g., close-up, full-
body) for human speech video generation. The associated
Orator framework leverages LLM-guided orchestration to
ensure cinematographic coherence, evaluating shot transi-
tions and emotional flow across extended narratives.

4.3. Physics, Reasoning, and Semantic Benchmarks

As video generative models advance, evaluating their ad-
herence to physical laws and logical reasoning becomes
paramount. The MoReSet benchmark [2] spans multi-
ple classes of Newtonian phenomena to assess physical
validity. Similarly, T2VPhysBench [14] shifts focus to-
ward first-principles physics, utilizing human judgment to
evaluate model adherence to core physical laws such as
conservation principles. Cognitive reasoning and semantic
compositionality are further probed by RULER-Bench [17]
and TC-Bench [13], which reveal significant reasoning
limitations and gaps in temporal transition success. Addi-
tionally, T2VTextBench [15] identifies critical deficiencies
in on-screen text fidelity, highlighting that even state-of-
the-art systems struggle with textual accuracy and temporal
consistency across complex strings.

4.4. Long-Form and Storytelling Benchmarks

The synthesis of long-form videos and structured narratives
introduces unique challenges in identity preservation and
temporal coherence. LoCoT2V-Bench [65] provides hierar-
chical metadata for character settings and camera behaviors,
revealing that while current models excel in global stability,
they struggle with long-term character consistency and
fine-grained alignment. To evaluate narrative fidelity, the
T2Vid2T framework [44] proposes a cyclical evaluation
approach, text-to-video-to-text, identifying significant per-
formance drops when models are tasked with generating
complex, multi-scene stories compared to simple factual
captions.

Architectural innovations such as NUWA-XL [58] ad-
dress the computational demands of long video generation.
By employing a Diffusion over Diffusion (DoD) hierarchy,

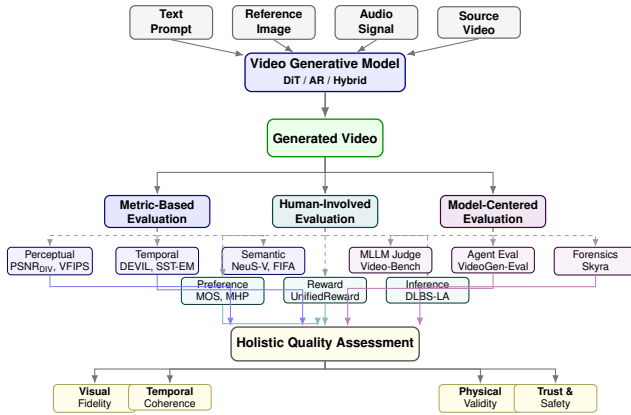


Figure 2. Unified evaluation pipeline for AI-generated video. Multi-modal inputs are processed by generative architectures and assessed via metric-based, human-involved, and model-centered paradigms. These signals converge into a holistic quality score covering visual fidelity, temporal coherence, physical validity, and trustworthiness.

324 NUWA-XL reduces inference time for 1024 frames by
325 over 90% and introduces the FlintstonesHD benchmark,
326 a densely annotated dataset for validating long-term coher-
327 ence. Furthermore, VideoThinkBench [46] frames video
328 generation as a reasoning paradigm, distinguishing between
329 vision-centric tasks (e.g., spatial puzzles) and text-centric
330 tasks (e.g., mathematical logic). Evaluations on this bench-
331 mark show that advanced models like Sora-2 can achieve
332 high accuracy on complex subsets like MATH and MMMU,
333 suggesting a path toward more intelligent, reasoning-aware
334 video synthesis.

335 5. Specialized and Conditional Benchmarks

336 Conditional generation requires fine-grained assessment of
337 structural and modal alignment.

338 5.1. Image-to-Video and Editing Benchmarks

339 Evaluation of conditional video generation and editing
340 requires fine-grained assessment of instruction alignment
341 and structural preservation. The FiVE benchmark [19]
342 addresses this by providing six editing types, ranging from
343 object substitution to material modification, across 420
344 object-level prompt pairs. It utilizes 15 metrics, including
345 the VLM-based FiVE-Acc, which reveals that rectified-flow
346 methods like Wan-Edit significantly outperform traditional
347 diffusion-based approaches in accuracy and background
348 preservation. For non-rigid dynamics, NRVBench [39]
349 introduces 180 physics-based videos across six categories
350 (e.g., fluid flow, fabric dynamics), employing the NRVE-
351 Acc metric to evaluate physical compliance and temporal
352 smoothness via VLM judges.

Creative and controllable editing is further benchmarked
by BalanceCC [12], which structures 100 videos with
diverse prompts (style, object, background changes) for
extensive user studies against eight state-of-the-art meth-
ods. Specialized tasks such as object removal are evaluated
by WIPER-Bench [23], which introduces the TokSim met-
ric. Leveraging DINOv3 embeddings, TokSim quantifies
temporal consistency, foreground-background coherence,
and removal effectiveness, providing a robust measure for
painting quality without requiring ground truth.

363 5.2. Audio-Video and Multi-Modal Benchmarks

The integration of auditory and visual modalities introduces
complex synchronization challenges. JarvisBench [29]
provides a large-scale dataset of text-captioned sounding
videos, introducing JarvisScore, a robust synchrony metric
that utilizes temporal segmentation to detect desynchroniza-
tion patterns. For conversational dynamics, MTAVG-Bench
[66] evaluates multi-talker dialogue across levels ranging
from signal fidelity to cinematic expression.

Specialized multi-modal tasks such as narrated presen-
tations are assessed via PresentEval [43], which measures
content fidelity and audience comprehension. Furthermore,
the JarvisGPT framework [28] demonstrates state-of-the-art
performance in joint audio-video comprehension and gen-
eration. By maintaining high semantic consistency and
temporal synchrony, these benchmarks facilitate the devel-
opment of more cohesive multi-modal generative systems.

380 6. Human Preference Alignment

Aligning models with human intent ensures the subjective
quality of generated content.

383 6.1. Human Preference Annotation Datasets

The alignment of video generative models with human
expectations requires large-scale, high-quality preference
datasets. AIGVE-60K [47] enables bidirectional bench-
marking for both text-to-video (T2V) generation and video-
to-text (V2T) interpretation, comprising 58,500 videos from
30 models and 2.6 million human annotations. This dataset
disentangles Mean Opinion Scores (MOS) into perceptual
quality and text-video correspondence, providing a gran-
ular evaluation across 20 task dimensions. Similarly, the
Multi-dimensional Human Preference (MHP) Dataset [62]
offers 918,315 pairwise choices across aesthetics, semantic
alignment, and detail quality, facilitating robust preference
modeling for visual synthesis.

Beyond general quality, specialized datasets address mo-
tion and safety. For text-to-motion generation, researchers
have explored preference learning using 3,528 pairs from
MotionGPT, identifying Direct Preference Optimization
(DPO) with Identity Preference Optimization (IPO) loss as
superior to traditional RLHF for alleviating overfitting [42].

403 To address safety concerns, SafeVid-350K [48] provides a
404 video-specific safety preference dataset constructed through
405 a systematic pipeline of adversarial question generation
406 and preference synthesis. Utilizing this dataset within a
407 DPO framework significantly enhances the safety alignment
408 of Video Large Multimodal Models (VLMs) without
409 compromising general utility.

410 6.2. Reward Models and Preference Optimization

411 The development of dedicated reward models is critical
412 for automating the alignment process. UnifiedReward [49]
413 represents a unified framework supporting both pairwise
414 ranking and pointwise scoring across image and video tasks.
415 By leveraging joint multi-task training, it enhances video
416 quality assessment and mitigates modality-specific over-
417 fitting, utilizing a two-stage filtering strategy to construct
418 high-quality preference data for optimization. To address
419 the computational cost of reward fine-tuning, InstructVideo
420 [59] recasts the process as an editing task, utilizing partial
421 inference steps and segmental sparse sampling to repurpose
422 established image reward models. Furthermore, temporally
423 attenuated reward mechanisms are employed to mitigate
424 modeling degradation. By prioritizing central frames and
425 ensuring motion continuity, these methods prevent visual
426 artifacts such as flickering during the fine-tuning process,
427 facilitating more robust human-aligned synthesis.

428 6.3. Inference-Time Alignment

429 Inference-time strategies offer a training-free alternative
430 for aligning generated content with user intent. Diffusion
431 Latent Beam Search (DLBS) [36] optimizes text-to-video
432 alignment by exploring multiple latent paths, leveraging
433 calibrated rewards from VLMs to select high-quality trajec-
434 tories. Similarly, the 3R framework [40] combines RAG-
435 based modifier extraction with preference optimization to
436 enhance fidelity and coherence without model retraining.

437 Prompt optimization further bridges the gap between
438 user instructions and model capabilities. PhyPrompt [52]
439 utilizes reinforcement learning to refine prompts for phys-
440 ical plausibility, integrating a dynamic reward curriculum
441 that balances semantic adherence with physical common-
442 sense. These advancements are supported by reward cali-
443 bration methodologies that utilize weighted linear combi-
444 nations of perceptual metrics to align automated scores with
445 human and VLM evaluations [36].

446 7. Safety and Trustworthiness

447 Ensuring safety and provenance is critical for responsible
448 AI deployment.

449 7.1. AI-Generated Video Detection and Forensics

450 The proliferation of high-fidelity video generative models
451 necessitates robust detection and forensic tools to mitigate

the risks of misinformation. Skyra [24] addresses this
by utilizing fine-grained human annotations and Chain-of-
Thought explanations to enhance artifact perception across
diverse generators. This is complemented by compre-
hensive forensics datasets like GenBuster-200K [50] and
standardized evaluation platforms like ViF-Bench.

Fine-grained artifact recognition is further formalized
by the Dense Video Artifact Recognition (DVAR) frame-
work [67], which taxonomizes defects across appearance,
motion, and camera axes. To detect hallucinations, the
SoraDetector framework [8] identifies semantic and tem-
poral inconsistencies via knowledge graphs. Furthermore,
recent deepfake detection approaches employ quality-based
matrix representations to enable frame-by-frame analysis,
identifying minute spatio-temporal anomalies that standard
playback might overlook.

468 7.2. Adversarial Robustness of T2V Models

469 The vulnerability of text-to-video (T2V) models to ad-
470 versarial perturbations poses significant security risks.
471 T2VAttack [5] demonstrates that minor word-level mod-
472 ifications, such as single-word substitution or insertion,
473 can cause substantial degradation in both semantic fidelity
474 and temporal dynamics across diverse architectures like
475 CogVideoX and HunyuanVideo. By targeting semantic
476 alignment and motion coherence, these attacks reveal a lack
477 of robustness in current generative pipelines. Furthermore,
478 the MUSE platform [10] introduces Inter-Turn Modality
479 Switching (ITMS) to probe whether safety alignments
480 generalize across modality boundaries. MUSE reports
481 that multi-turn strategies can achieve near-100% attack
482 success rates against models that otherwise exhibit perfect
483 single-turn refusal, highlighting the fragility of cross-modal
484 defenses.

485 Safety alignment is further complicated by "mismatched
486 generalization," where safety competencies learned from
487 static text and images fail to transfer to dynamic video
488 contexts. SafeVid [48] identifies this failure mode, noting
489 that models may reject harmful queries in isolation but
490 comply when they are paired with relevant video content.
491 To address this, the SafeVid-350K dataset and Direct Pref-
492 erence Optimization (DPO) are utilized to transfer robust
493 textual safety reasoning to the video domain, significantly
494 improving safety rates on the SafeVidBench benchmark
495 while maintaining model utility.

496 7.3. Video Watermarking and Provenance

497 Establishing the provenance of AI-generated videos is es-
498 sential for accountability and copyright protection. Vid-
499 Stamp [45] introduces a temporally-aware watermarking
500 framework that utilizes a two-stage fine-tuning process to
501 ensure spatial separation and temporal consistency. By
502 embedding unique frame-level bit sequences, VidStamp en-

503	ables precise temporal tamper localization, detecting frame	553
504	swaps and drops while maintaining high visual quality.	554
505	Furthermore, VidStamp supports dynamic watermarking,	555
506	allowing for adaptive provenance via control signals	
507	during inference. The framework demonstrates robust	
508	performance across diverse distortions, including cropping	
509	and MPEG4 compression. Comparative evaluations show	
510	that VidStamp achieves stronger detectability than existing	
511	methods such as VideoSeal and RivaGAN, providing a	
512	reliable solution for long-term video integrity and authentication.	
513		
514	7.4. Concept Erasure and Safety Generation	
515	The ability to selectively remove undesirable concepts or	
516	align generation with safety guidelines is paramount for	
517	responsible AI deployment. VideoEraser [55] provides	
518	a training-free framework for concept erasure, covering	
519	objects, styles, and explicit content, by adjusting prompt	
520	embeddings and noise guidance. This approach ensures	
521	broad generalizability across diverse T2V models without	
522	costly fine-tuning. Safety alignment is further formalized	
523	through closed-loop systems like SafeVid [48], which transfers	
524	robust textual safety reasoning to the video domain. By	
525	aligning models on safety preference datasets, researchers	
526	have demonstrated substantial improvements in safety rates	
527	on standardized benchmarks. For commercial applications,	
528	frameworks like BrandFusion [68] utilize multi-agent	
529	systems to ensure brand recognizability and contextually	
530	natural integration, facilitating the safe use of generative	
531	video in advertising.	
532	8. Open Challenges and Future Directions	
533	The path toward unified world models requires overcoming	
534	fundamental spatiotemporal gaps.	
535	8.1. Long-Form Coherence, Temporal Reasoning,	
536	and World Models	
537	A critical frontier in video generation is the transition from	
538	short-clip synthesis to long-form, coherent storytelling.	
539	Evaluations reveal a significant performance gap: while	
540	current models demonstrate strong perceptual quality, they	
541	exhibit deficiencies in fine-grained alignment and long-term	
542	character consistency. Achieving temporal coherence over	
543	extended durations remains an open challenge, necessitating	
544	architectures that can bridge the training-inference gap	
545	through hierarchical strategies and direct training on long	
546	sequences. Beyond visual consistency, the integration of	
547	temporal reasoning is essential for developing true world	
548	models. Framing video generation as a reasoning chain	
549	leverages dynamic visualization to solve complex spatial	
550	and textual puzzles, offering advantages over static image-	
551	based reasoning. However, existing models still struggle	
552	with basic temporal order relations. Future research must	
	focus on instilling relational consistency and leveraging	
	test-time scaling to transform video generators into unified	
	multimodal understanding and generation systems.	
	8.2. Physical World Simulation	
	A fundamental limitation of current neural video generators	
	is the divergence between high visual realism and physical	
	plausibility. Evidence from benchmarks like MoReSet	
	[2] and T2VPhysBench [14] indicates that state-of-the-art	
	models frequently violate basic physical laws, such as momentum	
	conservation and rigid-body collision dynamics, relying	
	instead on surface pattern memorization. To address this,	
	research is moving toward "metamorphic simulators" like	
	MagicTime [61], which encode physical world knowledge	
	through the simulation of complete metamorphic processes.	
	Principled directions for physically coherent synthesis	
	include the integration of simulation-aware architectures	
	and reinforcement learning. MoReGen [2] unifies multi-agent	
	LLMs with physics simulators and renderers to ensure	
	dynamic coherence via code-domain generation. Similarly,	
	PhyPrompt [52] utilizes Group Relative Policy Optimization	
	(GRPO) to inject physical reasoning into instructions,	
	overcoming prompt underspecification. Advancing this	
	field requires standardized evaluation platforms like	
	NRVBench [39], which specifically targets non-rigid	
	dynamics, such as fluidity and elasticity, using VLM-based	
	metrics to assess deformation plausibility and temporal	
	consistency.	
	8.3. Computational Efficiency, Scalability, and Sustainable Generation	
	The escalating computational demands of video generative	
	models necessitate a focus on efficiency and sustainability.	
	Analytical models reveal that energy consumption scales	
	quadratically with spatial and temporal dimensions, requiring	
	prioritization of compute-bound efficiency and hardware-aware	
	optimizations. To this end, H3AE [54] achieves real-time	
	decoding on mobile devices by utilizing high spatial-temporal	
	compression ratios, significantly reducing latent token counts	
	without compromising reconstruction fidelity.	
	Inference-time compute strategies further optimize the	
	efficiency-quality tradeoff. ADAPTOR [38] exploits temporal	
	redundancy through adaptive token reduction, while FlashVideo	
	[63] employs a two-stage framework that decouples prompt	
	fidelity from visual detail generation. These methods, combined	
	with search-based latent selection, enable high-quality video	
	synthesis at a fraction of the traditional computational cost,	
	paving the way for more accessible and sustainable generative	
	AI.	

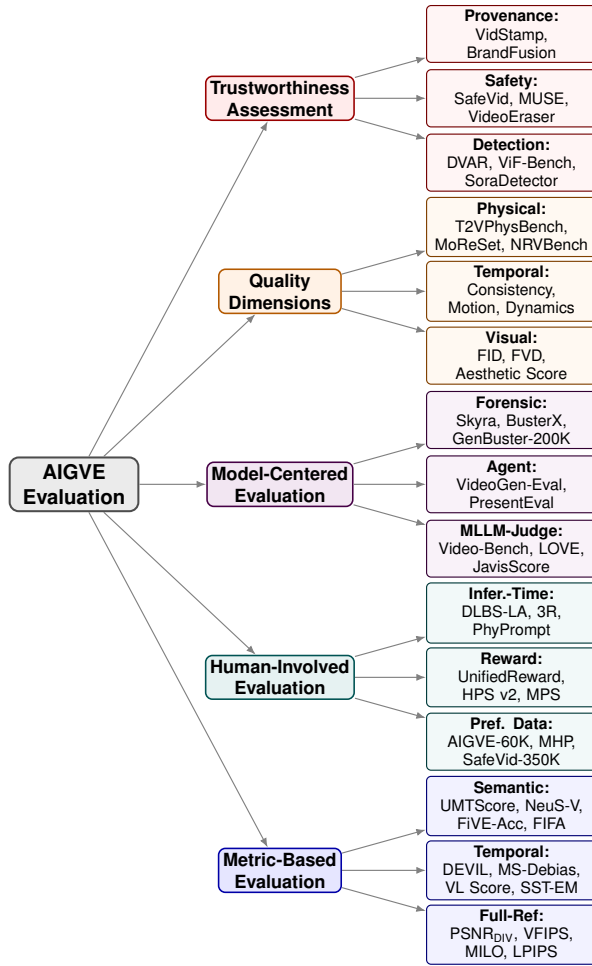


Figure 3. Hierarchical taxonomy of the AI-Generated Video Evaluation (AIGVE) landscape across five branches: metric-based, human-involved, model-centered, quality dimensions, and trustworthiness assessment.

8.4. Toward Unified, Adaptive, and Trustworthy Evaluation Frameworks

The current landscape of video evaluation is characterized by significant fragmentation, with individual benchmarks often prioritizing specific dimensions while neglecting others. A hierarchical taxonomy of this landscape, spanning metrics, human alignment, and model-centered paradigms, is provided in Fig. 3. To overcome the limitations of traditional metrics like FID and FVD, which fail to align with human perception, research is shifting toward model-centered evaluation frameworks [31]. Multimodal Large Language Models (MLLMs) are emerging as the cornerstone of next-generation assessment, providing unified and scalable frameworks that leverage few-shot scoring and chain-of-query processes to simulate human evaluation logic [16]. Agent-based systems like VideoGen-Eval [56]

further enhance this by integrating dynamic content structuring with objective judgment, achieving high alignment with professional human preferences across diverse quality axes.

The future roadmap for trustworthy video evaluation emphasizes the convergence of quality, safety, and provenance into a cohesive framework. Detection benchmarks such as Skyra [24] and GenBuster-200K [50] are transitioning from binary classification toward grounded artifact reasoning, utilizing MLLMs to identify spatio-temporal inconsistencies. Simultaneously, safety alignment frameworks like SafeVid [48] address the dynamic complexities that static alignments fail to capture. A truly unified protocol must simultaneously evaluate visual quality, temporal coherence, physical validity, and semantic reasoning, while incorporating cryptographic provenance and content credentials to mitigate the risks of deceptive deepfakes [20]. This holistic approach is essential for ensuring that AI-generated content remains aligned with human intentions, ethical standards, and real-world physical constraints.

9. Conclusion

This survey has analyzed the AI-generated video evaluation (AIGVE) field, highlighting the transition from basic visual synthesis to the pursuit of verifiable spatiotemporal quality. We have traced the evolution of evaluation methodologies, from traditional frame-level heuristics to model-centered paradigms that leverage Multimodal Large Language Models (MLLMs) for perceptually grounded assessment. Our analysis indicates that developing true video world models requires addressing fundamental gaps in physical simulation, long-form coherence, and temporal reasoning. As video generative models are increasingly integrated into downstream applications, the convergence of quality, safety, and provenance into unified evaluation frameworks has become a critical necessity. By establishing these rigorous standards, the field can ensure that the next generation of video foundation models remains aligned with human intent, physical reality, and ethical principles.

Acknowledgements

The authors developed all intellectual organization and critical synthesis in this review. PaperQA [34], a RAG-based literature tool, assisted in factual retrieval from surveyed papers; all retrieved information was verified by the authors.

References

- [1] Nithin C Babu, Aniruddha Mahapatra, Harsh Rangwani, R. Soundararajan, and Kuldeep Kulkarni. DynamicEval: Rethinking evaluation for dynamic text-to-video synthesis. In *arXiv.org*, 2025. 1, 2, 3, 4

- 665 [2] Xiangyu Bai, He Liang, Bishoy Galoaa, Utsav Nandi,
666 Shayda Moezzi, Yuhang He, and Sarah Ostadabbas. More-
667 gen: Multi-agent motion-reasoning engine for code-based
668 text-to-video synthesis. In *arXiv.org*, 2025. 4, 7 723
- 669 [3] Sangnie Bhardwaj, Ian Fischer, Johannes Ball'e, and Troy T.
670 Chinen. An unsupervised information-theoretic perceptual
671 quality metric. In *Neural Information Processing Systems*,
672 2020. 3 724
- 673 [4] Varun Biyyala, Bharat Chanderprakash Kathuria, Jialu Li,
674 and Youshan Zhang. Sst-em: Advanced metrics for evaluat-
675 ing semantic, spatial and temporal aspects in video editing.
676 In *2025 IEEE/CVF Winter Conference on Applications of*
677 *Computer Vision Workshops (WACVW)*, 2025. 3 725
- 678 [5] Chang bo Li, Yuecong Min, Jie Zhang, Zheng Yuan,
679 Shiguang Shan, and Xilin Chen. T2vattack: Adversarial
680 attack on text-to-video diffusion models. In *arXiv.org*, 2025.
681 6 726
- 682 [6] Jiaben Chen, Zixin Wang, Ailing Zeng, Yang Fu, Xueyang
683 Yu, Siyuan Cen, Julian Tanke, Yihang Chen, Koichi Saito,
684 Yuki Mitsufuji, and Chuang Gan. Talkcuts: A large-scale
685 dataset for multi-shot human speech video generation. In
686 *arXiv.org*, 2025. 4 727
- 687 [7] Ming Chen, Liyuan Cui, Wenyuan Zhang, Haoxian Zhang,
688 Yan Zhou, Xiaohan Li, Songlin Tang, Jiwen Liu, Borui
689 Liao, Hejia Chen, Xiaoqiang Liu, and Pengfei Wan. Midas:
690 Multimodal interactive digital-human synthesis via real-time
691 autoregressive video generation. In *arXiv.org*, 2025. 3 728
- 692 [8] Zhixuan Chu, Lei Zhang, Yichen Sun, Siqiao Xue, Zhibo
693 Wang, Zhan Qin, and Kui Ren. Sora detector: A unified
694 hallucination detection for large text-to-video models. In
695 *arXiv.org*, 2024. 6 729
- 696 [9] Conall Daly, D. Ramsook, and Anil C. Kokaram. An
697 efficient quality metric for video frame interpolation based
698 on motion-field divergence. In *International Workshop on*
699 *Quality of Multimedia Experience*, 2025. 2, 3 730
- 700 [10] Xin Dong, Sen Jia, Ming Wang, Yan Li, Zhenheng Yang,
701 Bingfeng Deng, and Hongyu Xiong. Coef-vq: Cost-efficient
702 video quality understanding through a cascaded multimodal
703 llm framework. In *Knowledge Discovery and Data Mining*,
704 2024. 6 731
- 705 [11] Fanda Fan, Chunjie Luo, Wanling Gao, and Jianfeng Zhan.
706 Aigcbench: Comprehensive evaluation of image-to-video
707 content generated by ai. *BenchCouncil Transactions on*
708 *Benchmarks, Standards and Evaluations*, 2024. 2, 4 732
- 709 [12] Ruoyu Feng, Wenming Weng, Yanhui Wang, Yuhui Yuan,
710 Jianmin Bao, Chong Luo, Zhibo Chen, and Baining Guo.
711 Ccredit: Creative and controllable video editing via diffusion
712 models. In *Computer Vision and Pattern Recognition*, 2023.
713 5 733
- 714 [13] Weixi Feng, Jiachen Li, Michael Stephen Saxon, Tsu-Jui
715 Fu, Wenhu Chen, and William Yang Wang. Tc-bench:
716 Benchmarking temporal compositionality in text-to-video
717 and image-to-video generation. In *arXiv.org*, 2024. 1, 2,
718 4 734
- 719 [14] Xuyang Guo, Jiayan Huo, Zhenmei Shi, Zhao Song, Jiahao
720 Zhang, and Jiale Zhao. T2vphysbench: A first-principles
721 benchmark for physical consistency in text-to-video genera-
722 tion. In *arXiv.org*, 2025. 2, 4, 7 735
- [15] Xuyang Guo, Jiayan Huo, Zhenmei Shi, Zhao Song, Jiahao
Zhang, and Jiale Zhao. T2vtextbench: A human evaluation
benchmark for textual control in video generation models. In
arXiv.org, 2025. 4 736
- [16] Hui Han, Siyuan Li, Jiaqi Chen, Yiwen Yuan, Yuling Wu,
Chak Tou Leong, Hanwen Du, Junchen Fu, Youhua Li, Jie
Zhang, Chi Zhang, Li jia Li, and Yongxin Ni. Video-bench:
Human-aligned video generation benchmark. In *Computer*
Vision and Pattern Recognition, 2025. 1, 4, 8 737
- [17] Xuming He, Zehao Fan, Hengjia Li, Fan Zhuo, Hankun Xu,
Senlin Cheng, Diwang Weng, Haifeng Liu, Can Ye, and Boxi
Wu. Ruler-bench: Probing rule-based reasoning abilities
of next-level video generation models for vision foundation
intelligence. In *arXiv.org*, 2025. 2, 4 738
- [18] Qiqi Hou, Abhijay Ghildyal, and Feng Liu. A perceptual
quality metric for video frame interpolation. In *European*
Conference on Computer Vision, 2022. 2, 3 739
- [19] Ming hui Li, Chenxi Xie, Yichen Wu, Lei Zhang, and
Mengyu Wang. Five: A fine-grained video editing bench-
mark for evaluating emerging diffusion and rectified flow
models. In *arXiv.org*, 2025. 3, 5 740
- [20] Muhammad Tanveer Jan, Mohammed G. Al-Jassani, Mar-
tinraj Nadar, Emmanuel Melchizedek Vunnava, Vangmai
Chakrapani, Hayat Ullah, Abbas Khan, Sardar Ali Abbas,
and B. Furht. Text-to-video generators: a comprehensive
survey. *Journal of Big Data*, 2025. 2, 8 741
- [21] Liqiang Jing, Viet Dac Lai, Seunghyun Yoon, Trung Bui, and
Xinya Du. Fifa: Unified faithfulness evaluation framework
for text-to-video and video-to-text generation. In *arXiv.org*,
2025. 2, 3 742
- [22] Tengchuan Kou, Xiaohong Liu, Zicheng Zhang, Chunyi Li,
Haoning Wu, Xiongkuo Min, Guangtao Zhai, and Ning
Liu. Subjective-aligned dataset and metric for text-to-video
quality assessment. In *ACM Multimedia*, 2024. 1, 3 743
- [23] Saksham Singh Kushwaha, Sayan Nag, Yapeng Tian, and
Kuldeep Kulkarni. Object-wiper : Training-free object and
associated effect removal in videos. In *arXiv.org*, 2026. 5 744
- [24] Yifei Li, Wenzhao Zheng, Yanran Zhang, Runze Sun, Yu
Zheng, Lei Chen, Jie Zhou, and Jiwen Lu. Skyra: Ai-
generated video detection via grounded artifact reasoning. In
arXiv.org, 2025. 2, 6, 8 745
- [25] Zhengang Li, Yan Kang, Yuchen Liu, Difan Liu, Tobias
Hinz, Feng Liu, and Yanzhi Wang. Sned: Superposition
network architecture search for efficient video diffusion
model. In *Computer Vision and Pattern Recognition*, 2024.
2 746
- [26] Mingxiang Liao, Hannan Lu, Xinyu Zhang, Fang Wan,
Tianyu Wang, Yuzhong Zhao, Wangmeng Zuo, Qixiang Ye,
and Jingdong Wang. Evaluation of text-to-video generation
models: A dynamics perspective. In *Neural Information*
Processing Systems, 2024. 2, 3, 4 747
- [27] Binhui Liu, Xin Liu, Anbo Dai, Zhiyong Zeng, Zhen Cui,
and Jian Yang. Dual-stream diffusion net for text-to-video
generation. In *arXiv.org*, 2023. 2 748
- [28] Kai Liu, Jungang Li, Yuchong Sun, Shengqiong Wu,
Jianzhang Gao, Daoan Zhang, Wei Zhang, Sheng Jin,
Sicheng Yu, Gen Zhan, Jiayi Ji, Fan Zhou, Liang Zheng,
749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779

- 780 Shuicheng Yan, Hao Fei, and Tat-Seng Chua. Javisgpt: A
781 unified multi-modal llm for sounding-video comprehension
782 and generation. In *arXiv.org*, 2025. 5 838
- 783 [29] Kai Liu, Wei Li, Lai Chen, Shengqiong Wu, Yanhao Zheng,
784 Jiayi Ji, Fan Zhou, Rongxin Jiang, Jiebo Luo, Hao Fei,
785 and Tat-Seng Chua. Javidit: Joint audio-video diffusion
786 transformer with hierarchical spatio-temporal prior synchro-
787 nization. In *arXiv.org*, 2025. 5 839
- 788 [30] Peng Liu, Xiaoming Ren, Feng Liu, Qingsong Xie, Quan-
789 long Zheng, Yanhao Zhang, Haonan Lu, and Yujiu Yang.
790 Dynamic-i2v: Exploring image-to-video generation models
791 via multimodal llm. In *arXiv.org*, 2025. 4 840
- 792 [31] Xiao Liu, Xinhao Xiang, Zizhong Li, Yongheng Wang,
793 Zhuoheng Li, Zhuosheng Liu, Weidi Zhang, Wei Ye, and
794 Jiawei Zhang. A survey of ai-generated video evaluation. In
795 *arXiv.org*, 2024. 2, 8 841
- 796 [32] Yuanxin Liu, Lei Li, Shuhuai Ren, Rundong Gao, Shicheng
797 Li, Sishuo Chen, Xu Sun, and Lu Hou. Fetv: A benchmark
798 for fine-grained evaluation of open-domain text-to-video
799 generation. In *Neural Information Processing Systems*, 2023.
800 1, 2, 3 842
- 801 [33] Yiting Lu, Xin Li, Bingchen Li, Zihao Yu, Fengbin Guan,
802 Xinrui Wang, Ruling Liao, Yan Ye, and Zhibo Chen. Aigc-
803 vqa: A holistic perception metric for aigc video quality
804 assessment. In *2024 IEEE/CVF Conference on Computer
805 Vision and Pattern Recognition Workshops (CVPRW)*, 2024.
806 4 843
- 807 [34] Jakub Lála, Odhran O’Donoghue, Aleksandar Shtedritski,
808 Sam Cox, Samuel G. Rodrigues, and Andrew D. White.
809 Paperqa: Retrieval-augmented generative agent for scientific
810 research. *arXiv preprint arXiv:2312.07559*, 2023. 8 844
- 811 [35] A. Melnik, Michal Ljubljanc, Cong Lu, Qi Yan, Weiming
812 Ren, and Helge J. Ritter. Video diffusion models: A survey.
813 *Trans. Mach. Learn. Res.*, 2024. 2 845
- 814 [36] Yuta Oshima, Masahiro Suzuki, Yutaka Matsuo, and Hiroki
815 Furuta. Inference-time text-to-video alignment with diffu-
816 sion latent beam search. In *arXiv.org*, 2025. 2, 6 846
- 817 [37] Liang Peng, Haoran Cheng, Zheng Yang, Ruisi Zhao, Linx-
818 uan Xia, Chaotian Song, Qinglin Lu, Boxi Wu, and Wei Liu.
819 Smoothvideo: Smooth video synthesis with noise constraints
820 on diffusion models for one-shot video tuning, 2023. 3 847
- 821 [38] Elia Peruzzo, Adil Karjauv, N. Sebe, Amir Ghodrati, and
822 A. Habibian. Adaptor: Adaptive token reduction for video
823 diffusion transformers. In *2025 IEEE/CVF Conference
824 on Computer Vision and Pattern Recognition Workshops
825 (CVPRW)*, 2025. 2, 7 848
- 826 [39] Bingzheng Qu, Kehai Chen, Xuefeng Bai, Jun Yu, and
827 Min Zhang. Beyond rigid: Benchmarking non-rigid video
828 editing. In *arXiv.org*, 2026. 5, 7 849
- 829 [40] Zillur Rahman, Alex Sheng, and Cristian Meo. Retrieval,
830 refinement, and ranking for text-to-video generation via
831 prompt optimization and test-time scaling, 2026. 6 850
- 832 [41] S. Sharan, Minkyu Choi, Sahil Shah, Harsh Goel, Moham-
833 mad Omama, and Sandeep P. Chinchali. Neuro-symbolic
834 evaluation of text-to-video models using formal verification.
835 In *Computer Vision and Pattern Recognition*, 2024. 2, 3 851
- 836 [42] Jenny Sheng, Matthieu Lin, Andrew Zhao, Kevin Pruvost,
837 Yu-Hui Wen, Yangguang Li, Gao Huang, and Yong-Jin Liu.
Exploring text-to-motion generation with human preference. 852
In *2024 IEEE/CVF Conference on Computer Vision and
Pattern Recognition Workshops (CVPRW)*, 2024. 5 853
- [43] Jingwei Shi, Zeyu Zhang, Biao Wu, Yanjie Liang, Meng
Fang, Ling Chen, and Yang Zhao. Presentagent: Multimodal
agent for presentation video generation. In *Conference on
Empirical Methods in Natural Language Processing*, 2025.
5 854
- [44] Andrew Shin, Yusuke Mori, and Kunitake Kaneko. The lost
melody: Empirical observations on text-to-video generation
from a storytelling perspective. In *arXiv.org*, 2024. 4 855
- [45] Mohammadreza Teymoorianfard, Shiqing Ma, and Amir
Houmansadr. Vidstamp: A temporally-aware watermark
for ownership and integrity in video diffusion models. In
arXiv.org, 2025. 2, 6 856
- [46] Jingqi Tong, Yurong Mou, Hangcheng Li, Mingzhe Li, Yong
Yang, Ming Zhang, Qiguang Chen, Tianyi Liang, Xiaomeng
Hu, Y. Zheng, Xinchu Chen, Jun Zhao, Xuanjing Huang,
and Xipeng Qiu. Thinking with video: Video generation as
a promising multimodal reasoning paradigm. In *arXiv.org*,
2025. 5 857
- [47] Jiarui Wang, Huiyu Duan, Ziheng Jia, Yu Zhao, Woo Yi
Yang, Zicheng Zhang, Zijian Chen, Juntong Wang, Yuke
Xing, Guangtao Zhai, and Xiongkuo Min. Love: Bench-
marking and evaluating text-to-video generation and video-
to-text interpretation. In *arXiv.org*, 2025. 2, 4, 5 858
- [48] Yixu Wang, Jiabin Song, Yifeng Gao, Xin Wang, Yang Yao,
Yan Teng, Xingjun Ma, Yingchun Wang, and Yu-Gang Jiang.
Safevid: Toward safety aligned video large multimodal
models. In *arXiv.org*, 2025. 2, 6, 7, 8 859
- [49] Yibin Wang, Yuhang Zang, Hao Li, Cheng Jin, and Jiaqi
Wang. Unified reward model for multimodal understanding
and generation. In *arXiv.org*, 2025. 2, 6 860
- [50] Haiquan Wen, Yiwei He, Zhenglin Huang, Tianxiao Li, Zi-
han Yu, Xingru Huang, Lu Qi, Baoyuan Wu, Xiangtai Li, and
Guangliang Cheng. Busterx: MLLM-powered ai-generated
video forgery detection and explanation. In *arXiv.org*, 2025.
2, 6, 8 861
- [51] Wenming Weng, Ruoyu Feng, Yanhui Wang, Qi Dai,
Chunyu Wang, Dacheng Yin, Zhiyuan Zhao, Kai Qiu, Jian-
min Bao, Yuhui Yuan, Chong Luo, Yueyi Zhang, and Zhiwei
Xiong. Artv: Auto-regressive text-to-video generation with
diffusion models. In *2024 IEEE/CVF Conference on Com-
puter Vision and Pattern Recognition Workshops (CVPRW)*,
2023. 2 862
- [52] Shang Wu, Chenwei Xu, Zhuofan Xia, Weijian Li, Lie
Lu, Pranav Maneriker, Fan Du, Manling Li, and Han Liu.
Phyprompt: RL-based prompt refinement for physically plau-
sible text-to-video generation, 2026. 2, 6, 7 863
- [53] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng
Zhu, Rui Zhao, and Hongsheng Li. Human preference score
v2: A solid benchmark for evaluating human preferences of
text-to-image synthesis. In *arXiv.org*, 2023. 2, 3 864
- [54] Yushu Wu, Yanyu Li, Ivan Skorokhodov, Anil Kag, Willi
Menapace, Sharath Girish, Aliaksandr Siarohin, Yanzhi
Wang, and Sergey Tulyakov. H3ae: High compression,
high speed, and high quality autoencoder for video diffusion
models. In *arXiv.org*, 2025. 7 865

- 896 [55] Naen Xu, Jinghuai Zhang, Changjiang Li, Zhi Chen, Chunyi
897 Zhou, Qingming Li, Tianyu Du, and Shouling Ji. Video-
898 eraser: Concept erasure in text-to-video diffusion models.
899 In *Conference on Empirical Methods in Natural Language*
900 *Processing*, 2025. 2, 7
- 901 [56] Yuhang Yang, Ke Fan, Shangkun Sun, Hongxiang Li, Ailing
902 Zeng, Feilin Han, Wei hao Zhai, Wei Liu, Yang Cao, and
903 Zhengjun Zha. Videogen-eval: Agent-based system for
904 video generation evaluation. In *arXiv.org*, 2025. 1, 2, 4,
905 8
- 906 [57] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu
907 Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan
908 Zhang, Guanyu Feng, Da Yin, Xiaotao Gu, Yuxuan Zhang,
909 Weihan Wang, Yean Cheng, Ting Liu, Bin Xu, Yuxiao Dong,
910 and Jie Tang. Cogvideox: Text-to-video diffusion models
911 with an expert transformer. In *International Conference on*
912 *Learning Representations*, 2024. 2
- 913 [58] Sheng-Siang Yin, Chenfei Wu, Huan Yang, Jianfeng Wang,
914 Xiaodong Wang, Minheng Ni, Zhengyuan Yang, Linjie Li,
915 Shuguang Liu, Fan Yang, Jianlong Fu, Gong Ming, Lijuan
916 Wang, Zicheng Liu, Houqiang Li, and Nan Duan. Nuwa-xl:
917 Diffusion over diffusion for extremely long video generation.
918 In *Annual Meeting of the Association for Computational*
919 *Linguistics*, 2023. 4
- 920 [59] Hangjie Yuan, Shiwei Zhang, Xiang Wang, Yujie Wei,
921 Tao Feng, Yining Pan, Yingya Zhang, Ziwei Liu, Samuel
922 Albanie, and Dong Ni. Instructvideo: Instructing video
923 diffusion models with human feedback. In *Computer Vision*
924 *and Pattern Recognition*, 2023. 6
- 925 [60] Hangjie Yuan, Weihua Chen, Jun Cen, Hu Yu, Jingyun
926 Liang, Shuning Chang, Zhihui Lin, Tao Feng, Pengwei Liu,
927 Jiazheng Xing, Hao Luo, Jiasheng Tang, Fan Wang, and Yi
928 Yang. Lumos-1: On autoregressive video generation from a
929 unified model perspective. In *arXiv.org*, 2025. 2
- 930 [61] Shenghai Yuan, Jinfa Huang, Yujun Shi, Yongqi Xu, Ruijie
931 Zhu, Bin Lin, Xinhua Cheng, Li Yuan, and Jiebo Luo. Mag-
932 ictime: Time-lapse video generation models as metamorphic
933 simulators. *IEEE Transactions on Pattern Analysis and*
934 *Machine Intelligence*, 2024. 7
- 935 [62] Sixian Zhang, Bo Wang, Junqiang Wu, Yan Li, Tingting
936 Gao, Di Zhang, and Zhongyuan Wang. Learning multi-
937 dimensional human preference for text-to-image generation.
938 In *Computer Vision and Pattern Recognition*, 2024. 3, 5
- 939 [63] Shilong Zhang, Wenbo Li, Shoufa Chen, Chongjian Ge,
940 Peize Sun, Yida Zhang, Yi Jiang, Zehuan Yuan, Binyue
941 Peng, and Ping Luo. Flashvideo: Flowing fidelity to detail
942 for efficient high-resolution video generation. In *arXiv.org*,
943 2025. 2, 7
- 944 [64] Zhichao Zhang, Wei Sun, and Guangtao Zhai. A perspective
945 on quality evaluation for ai-generated videos. In *Italian*
946 *National Conference on Sensors*, 2025. 1, 2
- 947 [65] Xiangqing Zheng, Chengyue Wu, Kehai Chen, and Min
948 Zhang. Locot2v-bench: Benchmarking long-form and com-
949 plex text-to-video generation, 2025. 2, 4
- 950 [66] Yang-Hao Zhou, Haitian Li, Rexar Lin, Heyan Huang, Jinx-
951 ing Zhou, Changsen Yuan, Tian Lan, Ziqi Zhou, Yudong Li,
952 Jiajun Xu, Jing Liao, Yi Cheng, Xuefeng Chen, Xian-Ling
Mao, and Yousheng Feng. Mtavg-bench: A comprehen- 953
sive benchmark for evaluating multi-talker dialogue-centric 954
audio-video generation, 2026. 5 955
- [67] Chen Zhu, Jiashu Zhu, Yanxun Li, Meiqi Wu, Bingze Song, 956
Chubin Chen, Jiahong Wu, Xiangxiang Chu, and Yangang 957
Wang. Artifact-aware evaluation for high-quality video 958
generation. In *arXiv.org*, 2026. 2, 6 959
- [68] Zihao Zhu, Ruotong Wang, Siwei Lyu, Min Zhang, and 960
Baoyuan Wu. Brandfusion: A multi-agent framework for 961
seamless brand integration in text-to-video generation, 2026. 962
7 963
- [69] Ugur Çogalan, Mojtaba Bemana, K. Myszkowski, Hans- 964
Peter Seidel, and Colin Groth. Milo: A lightweight percep- 965
tual quality metric for image and latent-space optimization. 966
ACM Transactions on Graphics, 2025. 3 967