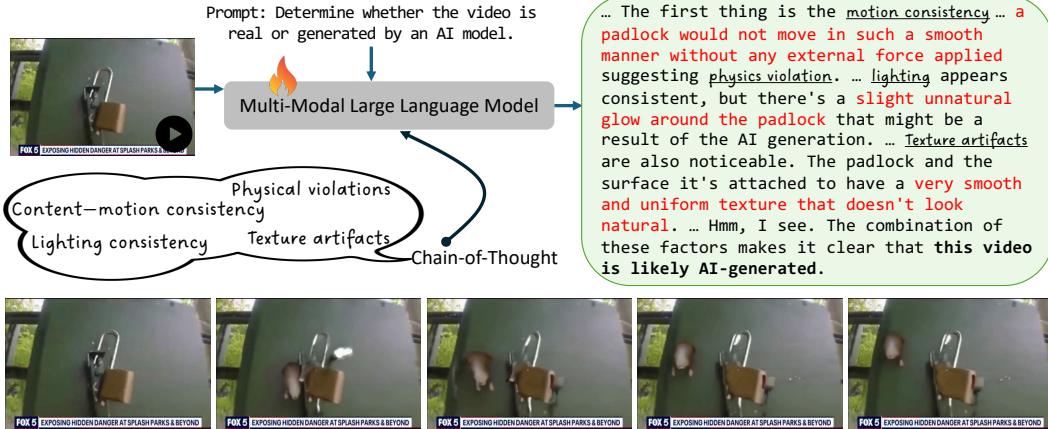


# VIDGUARD-R1: AI-GENERATED VIDEO DETECTION AND EXPLANATION VIA REASONING MLLMS AND RL

000  
001  
002  
003  
004  
005 **Anonymous authors**  
006 Paper under double-blind review  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017



022 Figure 1: Overall framework of **VidGuard-R1**. We present the first video authenticity detector  
023 based on multi-modal large language models (MLLMs), which generates a chain-of-thought reasoning  
024 process along with the final answer.

## ABSTRACT

027 With the rapid advancement of AI-generated videos, there is an urgent need for  
028 effective detection tools to mitigate societal risks such as misinformation and reputational  
029 harm. In addition to accurate classification, it is essential that detection  
030 models provide interpretable explanations to ensure transparency for regulators  
031 and end users. To address these challenges, we introduce **VidGuard-R1**, the first  
032 video authenticity detector that fine-tunes a multi-modal large language model  
033 (MLLM) using group relative policy optimization (GRPO). Our model delivers  
034 both highly accurate judgments and insightful reasoning. We curate a challenging  
035 dataset of 140k real and AI-generated videos produced by state-of-the-art gener-  
036 ation models, carefully designing the generation process to maximize discrimi-  
037 nation difficulty. We then fine-tune Qwen-VL using GRPO with two specialized  
038 reward models that target temporal artifacts and generation complexity. Extensive  
039 experiments demonstrate that **VidGuard-R1** achieves state-of-the-art zero-shot  
040 performance on existing benchmarks, with additional training pushing accuracy  
041 above 95%. Case studies further show that **VidGuard-R1** produces precise and  
042 interpretable rationales behind its predictions.

## 1 INTRODUCTION

043 In the past year, we have witnessed unprecedented progress in video generation models, with dra-  
044 matic improvements in realism and quality. The release of powerful models such as Sora (Brooks  
045 et al., 2024), Wan (wan, 2025), and HunyanVideo (Kong et al., 2024) has made AI-generated videos  
046 more accessible to the public, further blurring the line between synthetic videos and real ones. At  
047 the same time, these advancements have led to a series of social risks, including the spread of misin-  
048 formation, violations of privacy rights, damage to personal reputations, and increased susceptibility  
049 to scams and fraud.

050 Motivated by its practical significance, several pioneering works have been developed to detect AI-  
051 generated videos. Early approaches primarily targeted DeepFake-style facial forgeries (Qian et al.,  
052 2020b; Tan et al., 2024; Gu et al., 2021), which often assumed single-subject, frontal-face scenarios

under constrained settings. These assumptions diverge significantly from open-domain, multi-scene videos produced by modern generative models. More recent detectors leverage spatial-temporal consistency (Ma et al., 2024; Bai et al., 2024b; Liu et al., 2024); however, such methods are limited in capturing higher-level semantic or causal inconsistencies and can be easily bypassed by post-processing techniques. Other methods are trained on curated fake video detection datasets (Chen et al., 2024a; Ni et al., 2025; Kundu et al., 2025), but these benchmarks often lack coverage of newly emerging models and fail to reflect the full diversity of generative capabilities. A recent benchmark (Chen et al., 2024a) shows that even state-of-the-art detectors still struggle to reliably identify videos from advanced models like Sora. Furthermore, these detectors typically offer only binary decisions without accompanying explanations, which raises concerns for transparency, especially when detection outcomes affect content moderation or legal accountability. Users are also more likely to trust detection systems that provide interpretable reasoning.

Recent advances in multi-modal large language models (MLLMs) have significantly enhanced video understanding, enabling detailed explanations of model decisions (Bai et al., 2023; Zhang et al., 2024b). This makes them promising candidates for detecting and explaining AI-generated videos. However, directly applying existing MLLMs, including advanced models like GPT-4o, yields subpar performance on current benchmarks, underscoring the need for supervised fine-tuning (SFT). As an initial step, we applied SFT to the Qwen2.5-VL-7B model (Bai et al., 2025). While the model achieved strong overall performance, it remained limited in its ability to explain why a video is fake, revealing shortcomings in its reasoning capability.

To address this, we adopt reinforcement learning (RL), which has shown promise in enhancing LLM reasoning (Guo et al., 2025). Notably, Video-R1 (Feng et al., 2025) outperforms commercial models on video reasoning tasks. RL enables MLLMs to develop self-improving reasoning via outcome-based rewards. We hypothesize that RL fine-tuning can help models detect subtle temporal and generative artifacts. Key to this is designing effective reward models. Simple binary rewards (e.g., 1 for real, 0 for fake) are insufficient. Instead, we propose two strategies: (1) injecting temporal artifacts into both real and fake videos to encourage temporal reasoning, and (2) assigning higher rewards to videos generated with more diffusion steps, which are harder to detect. Incorporating these into a group relative policy optimization (GRPO) framework leads to over 86% accuracy on our dataset and 95% accuracy on two benchmarks.

- We introduce **VidGuard-R1**, the first video authenticity detector that fine-tunes the MLLM using GRPO. The model leverages the pretrained knowledge of MLLMs for accurate classification and employs reinforcement learning for effective exploration. To further enhance performance, we design two specialized reward models that target temporal artifacts and generation complexity based on diffusion steps.
- We construct a challenging dataset of 140k real/fake video pairs for AI-generated video detection. By employing state-of-the-art generation models and carefully controlling the process, we ensure that distinguishing real from fake is non-trivial.
- **VidGuard-R1** achieves state-of-the-art zero-shot performance on existing benchmarks, with accuracy exceeding 95%. Case studies further highlight its ability to produce accurate and interpretable explanations.

## 2 RELATED WORKS

### 2.1 AI-GENERATED VIDEO DETECTION METHOD

Recent research on AI-generated video detection has largely focused on deepfake videos with synthetic faces (Pei et al., 2024), using spatial-temporal consistency, frequency artifacts, or data-driven approaches. These methods often struggle to generalize beyond face-centric content to more diverse, real-world videos. Recently, general video detection methods have emerged: AIGDet (Bai et al., 2024a) captures spatial-temporal anomalies, DeCoF (Ma et al., 2024) exploits frame consistency, and diffusion-based representations track temporal dynamics (Liu et al., 2024). Other works identify appearance, motion, and geometry as key factors for classifier training (Chang et al., 2024). Multimodal LLMs have also been explored for forgery detection: FakeShield (Xu et al., 2024) uses supervised fine-tuning (SFT) for image forgery detection, while SafeWatch (Chen et al., 2024b) combines SFT and direct preference optimization (DPO) for video guardrails. In contrast, our work is the first to fine-tune a multi-modal LLM with group relative policy optimization (GRPO) for AI-

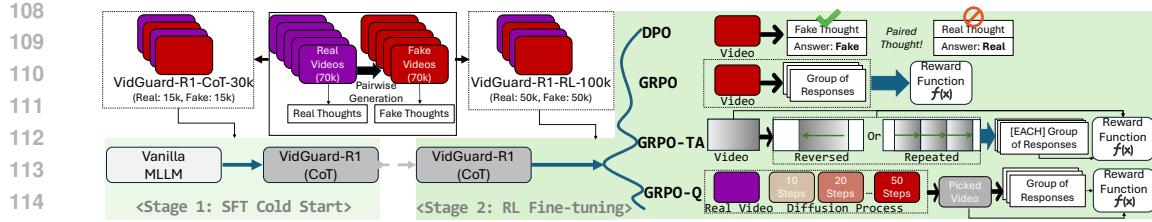


Figure 2: The overall training framework of **VidGuard-R1**, consisting of two stages: (1) supervised fine-tuning (SFT) for chain-of-thought (CoT) initialization, and (2) reinforcement learning-based fine-tuning to enable deeper reasoning.

generated video detection, demonstrating strong generalization across recent generative models and benchmark datasets.

## 2.2 AI-GENERATED VIDEO DETECTION DATASET

Given the recency of this research area, only a limited number of benchmarks have been introduced. The generated video dataset (GVD) (Bai et al., 2024a) (11k samples) and GenVideo (Chen et al., 2024a) (with millions of samples) consider settings where both training and test videos are generated by the same series of models. However, these benchmarks lack prompt/image–video pairs, semantic labels, or cross-source settings. GVF (2.8k samples) contains prompts/images–video pairs and semantic labels, but does not provide cross-source settings. GenVidBench (Ni et al., 2025) consists of 100k videos and incorporates cross-source settings, but the video generation models used are less advanced, such as CogVideo and SVD.

Moreover, existing datasets often contain shortcuts in resolution, frame rate, bitrate, or data-source imbalance, enabling models to exploit superficial cues rather than learn intrinsic visual realism. To address these limitations, we construct a curated dataset of 140,000 real–fake videos generated with state-of-the-art video generation models: HunyuanVideo (Kong et al., 2024) and CogVideoX (Yang et al., 2024). Our dataset explicitly standardizes bitrate, resolution, frame rate, and content distribution, resulting in a shortcut-free benchmark that encourages models to rely on semantic and temporal realism rather than on metadata artifacts.

## 3 METHODOLOGY

Figure 2 illustrates the **VidGuard-R1** framework, which consists of two stages. We first apply supervised fine-tuning (SFT) to the multimodal large language model (MLLM), followed by direct preference optimization (DPO) and group relative policy optimization (GRPO) based on the collected datasets. We further develop two GRPO variants by introducing temporal artifacts and leveraging videos generated with varying diffusion steps.

### 3.1 DATA COLLECTION

#### 3.1.1 DATA CONSTRUCTION FOR VIDEO REALISM DISCRIMINATION

High-quality training data is essential for video reasoning in MLLMs. However, many existing benchmarks for real vs. generated video classification, such as GenVideo (Chen et al., 2024a) and GenVidBench (Ni et al., 2025), exhibit uncontrolled discrepancies in basic metadata—e.g., real videos are often longer than 10 seconds while generated ones are typically under 4 seconds in GenVideo. Moreover, they reveal clear modality-level gaps in motion dynamics and content contrasts between real and generated videos. These differences introduce unintended shortcuts, enabling models to rely on superficial cues like duration or resolution rather than actual visual realism. As a result, **VidGuard-R1** attains over 96% accuracy on both GenVideo and GenVidBench by exploiting such artifacts. To mitigate this reward hacking behavior, we construct a curated dataset with standardized video properties, encouraging models to focus on intrinsic visual content.

We collect real videos from the InternVid (Wang et al., 2023c) and ActivityNet (Caba Heilbron et al., 2015) datasets and generate their corresponding fake counterparts using HunyuanVideo (Kong et al., 2024) and CogVideoX (Yang et al., 2024). We specifically choose these two models because they support conditioning on both the first-frame image and a text description—an essential requirement for generating videos that are contextually aligned with their real counterparts. To achieve such

162 alignment, we provide the generation models with the first frame of each real video along with  
 163 a textual caption describing its content. For ActivityNet, which lacks native captions, we extract  
 164 concise descriptions using Qwen2.5-VL 72B. This pairing strategy mitigates content-based biases  
 165 and forces the model to reason over subtle visual details.

### 166 3.1.2 COLLECTING CHAIN OF THOUGHT (CoT) ANNOTATION

168 Eliciting deliberate, step-by-step reasoning in MLLMs requires high-quality CoT supervision. To  
 169 this end, we leverage Qwen-2.5-VL (72B) to extract salient visual cues from each video and guide  
 170 the model toward a deeper understanding. Specifically, we query the model with critical factors  
 171 known to distinguish real from generated content—motion consistency, lighting consistency, texture  
 172 artifacts, and physical plausibility violations. These targeted prompts encourage detailed reasoning  
 173 grounded in visual evidence.

174 However, current MLLMs lack the capacity to reliably distinguish real from fake videos on their  
 175 own. To compensate, we provide ground-truth labels during prompt construction and instruct the  
 176 model to generate CoT rationales conditioned on the given label. While these rationales do not re-  
 177 flect genuine discrimination ability, they capture rich contextual cues—such as object interactions,  
 178 background details, and lighting inconsistencies—that are highly informative. These CoT annota-  
 179 tions serve as useful clues for subsequent reinforcement learning fine-tuning. For prompt templates  
 180 used in CoT generation, please refer to our supplementary materials.

### 181 3.2 SUPERVISED AND RL FINE-TUNING

183 We begin with SFT, where the model is trained to mimic the ground-truth reasoning process. Given  
 184 a video  $x$  and its annotation  $y$  from the collected dataset, the model parameters  $\theta$  are optimized by  
 185 minimizing the negative log-likelihood  $\mathcal{L}_{\text{SFT}}(\theta) = -\sum_{t=1}^T \log p_\theta(y_t | y_{<t}, x)$ . To align the model  
 186 outputs with human preferences, we apply DPO, which updates the model based on pairwise pref-  
 187 erence data without explicit reward modeling. Given a preferred response  $y_w$  and a less-preferred  
 188 response  $y_l$  for the same video  $x$ , the DPO loss encourages the model to prefer  $y_w$  over  $y_l$  compared  
 189 to a reference model  $p_{\text{ref}}$ :

$$190 \mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \left[ \log \sigma \left( \beta \log \frac{p_\theta(y_w|x)}{p_{\text{ref}}(y_w|x)} - \beta \log \frac{p_\theta(y_l|x)}{p_{\text{ref}}(y_l|x)} \right) \right]$$

192 where  $\sigma(\cdot)$  is the sigmoid function and  $\beta$  controls the preference strength. This method allows  
 193 fine-tuning using preference comparisons without requiring scalar rewards.

195 Finally, we adopt GRPO from DeepSeek R1 (Guo et al., 2025), which generalizes RLHF to group-  
 196 level comparisons. Given a query video  $x$  and a group of generated outputs  $\{o_i\}_{i=1}^G$ , the model  
 197 is trained to assign higher probabilities to outputs with higher rewards. The GRPO objective is:

$$199 \mathcal{L}_{\text{GRPO}}(\theta) = -\mathbb{E}_{(x, o_1:G) \sim D} \left[ \frac{1}{G} \sum_{i=1}^G \min \left( \frac{p_\theta(o_i|x)}{p_{\text{ref}}(o_i|x)} A_i, \text{clip} \left( \frac{p_\theta(o_i|x)}{p_{\text{ref}}(o_i|x)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) - \beta D_{\text{KL}}(p_\theta \| p_{\text{ref}}) \right]$$

200 where  $\epsilon$  is a clipping threshold and  $\beta$  regularizes the policy to stay close to the reference model.  
 201 The advantage term  $A_i$  normalizes the reward  $r_i$  for output  $o_i$  within the group, computed as  
 202  $A_i = \frac{r_i - \mu_x}{\sigma_x}$ , where  $\mu_x$  and  $\sigma_x$  are the mean and standard deviation of  $\{r_i\}_{i=1}^G$ . GRPO thus enables  
 203 learning from relative ranking among multiple responses, capturing nuanced distinctions in quality  
 204 across outputs.

### 206 3.3 VIDGUARD-R1

#### 207 3.3.1 OVERVIEW

209 Figure 2 illustrates the training pipeline of **VidGuard-R1**. Following the data collection  
 210 procedure, we construct two datasets of different scales: VidGuard-R1-CoT-30k and  
 211 VidGuard-R1-RL-100k. We adopt Qwen2.5-VL-7B as the base MLLM and train it using our  
 212 proposed fine-tuning framework.

213 The first stage is supervised fine-tuning initialization using the VidGuard-R1-CoT-30k dataset,  
 214 which contains videos paired with chain-of-thought (CoT) annotations. This stage establishes foun-  
 215 dational reasoning ability and equips the model with basic cross-modal alignment and visual under-  
 216 standing. The resulting model is referred to as **VidGuard-R1 (CoT)**.

In the second stage, we apply two reinforcement learning methods—DPO and GRPO—to further refine the model on a larger and more diverse dataset, **VidGuard-R1-RL-100k**. DPO aligns the model with high-quality preference signals via pairwise comparisons, requiring the construction of preference pairs. Specifically, since our dataset includes pairwise real and fake videos, each sample is annotated with CoT rationales for both perspectives. For DPO training, we construct preference pairs by swapping these CoTs. For a real video, the CoT supporting its authenticity with the answer “real” serves as the positive annotation, while the CoT from the paired fake video with the answer “fake” is used as the negative annotation. In contrast, GRPO encourages consistent performance across grouped outputs by leveraging structural regularization. As it does not rely on preference annotations, video labels are directly used as reward signals. The resulting models are denoted as **VidGuard-R1 (DPO)** and **VidGuard-R1 (GRPO)**.

We introduce two variants, GRPO-TA and GRPO-Q, to further enhance detection performance. These methods extend the original GRPO framework by adjusting reward values according to the difficulty of detecting fake videos. Detailed descriptions are provided in the following sections.

### 3.3.2 GRPO WITH TEMPORAL ARTIFACTS (GRPO-TA)

While standard GRPO performs well in video discrimination by leveraging local visual cues—such as pixel distortions and lighting inconsistencies—it often overlooks temporal inconsistencies, which are crucial for detecting generated videos. To address this limitation, we introduce **GRPO with temporal artifacts (GRPO-TA)**, a variant that explicitly promotes temporal reasoning through a contrastive reward adjustment.

We apply two common temporal artifacts: (1) repeating a specific video segment and (2) reversing the frame sequence within a segment. These manipulations are applied probabilistically, with the manipulated region selected based on a Gaussian distribution over the video timeline.

Specifically, for each input query, we generate two sets of model outputs:  $\{o_i\}_{i=1}^G$  for the original video, and  $\{\tilde{o}_i\}_{i=1}^{G'}$  for the corresponding manipulated video with temporal artifacts. These videos should be classified as fake videos. In GRPO-TA, we assign additional rewards when the model correctly classifies temporally manipulated videos as fake. Consider two numbers,  $\alpha_1 > \alpha_2$ . Detecting temporal artifacts in videos manipulated from real content tends to be more challenging than identifying those derived from fake videos. This is because real videos typically exhibit coherent and natural motion, so temporal manipulations such as frame shuffling or repetition can be subtle and difficult to detect. In contrast, generated videos often contain artifacts like unstable motion or low temporal consistency, which make further manipulations more visually salient. To reflect this asymmetry in difficulty, we assign the model a higher reward  $\alpha_1$  when the original video  $o_i$  is real, and a moderate reward  $\alpha_2$  when the original video is fake. This is defined as:

$$w_i = \begin{cases} \alpha_1, & \text{if } \tilde{o}_i = \text{fake} \text{ and } y_i = \text{real} \\ \alpha_2, & \text{if } \tilde{o}_i = \text{fake} \text{ and } y_i = \text{fake} \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

where  $y_i$  denotes the label of the  $i$ -th video. In the experiments, we set the hyperparameters to  $\alpha_1 = 0.5$  and  $\alpha_2 = 0.3$ . This additional reward,  $w_i$ , is designed to be applied conditionally. Specifically, for a given sample, we only add  $w_i$  to the original GRPO reward if two conditions are met: the model’s prediction on the original video ( $O_i$ ) must be correct, and the overall accuracy on the group of manipulated videos ( $\tilde{p}$ ) must exceed a predefined threshold  $\mu$ . This ensures that we only reward the model for temporal reasoning when it already has a solid baseline performance. The final reward of GRPO-TA is given by

$$r_i^{\text{GRPO-TA}} = \begin{cases} r_i^{\text{GRPO}} + w_i, & \text{if } o_i \text{ is correct and } \tilde{p} > \mu \\ r_i^{\text{GRPO}}, & \text{otherwise} \end{cases} \quad (2)$$

where  $r_i$  denotes the original GRPO reward, set to 1 if the model prediction on the original video is correct, and 0 otherwise. The additional reward  $w_i$  is applied only when both the original prediction is correct and the group of responses for the temporally manipulated videos achieves higher accuracy. In the experiments, we set  $\mu = 0.8$ .

### 3.3.3 GRPO WITH QUALITY EVOLUTIONARY VIDEOS (GRPO-Q)

Our motivation is to extend the model’s capability to detect videos based on quality. Given the subjective nature of quality assessment, we avoid relying on large-scale human annotations. Instead, we

270 leverage diffusion-based video generation by systematically varying the number of reverse diffusion  
 271 steps to produce videos with distinct quality levels.  
 272

273 As in GRPO-TA,  $o_i \in \mathcal{Y}$  and  $y_i \in \mathcal{Y}$  denote the model output and ground-truth label, with  $\mathcal{Y} =$   
 274  $\{\text{real}\} \cup \{\text{fake-}s\}$ , where  $s$  is the diffusion step. A reward is given for an exact match, and no reward  
 275 is assigned if the real/fake classification is incorrect. In GRPO-Q, if the model correctly classifies a  
 276 fake video but selects an incorrect diffusion step, we assign a partial reward based on the distance  
 277 between the predicted and ground-truth diffusion steps. The GRPO-Q reward is defined as follows:  
 278

$$279 \quad r_i^{\text{GRPO-Q}} = \begin{cases} 0, & \text{if } (o_i = \text{real} \text{ and } y_i \neq \text{real}) \text{ or } (o_i \neq \text{real} \text{ and } y_i = \text{real}) \\ \delta, & \text{if } o_i = y_i \\ |g(o_i, y_i)|, & \text{if } o_i, y_i \in \mathcal{Y} \setminus \{\text{real}\}. \end{cases} \quad (3)$$

281 The first scenario occurs when the model fails to correctly classify the video as real or fake. The  
 282 second scenario, where  $\delta = 1$ , represents an exact match in prediction, including the diffusion  
 283 progression. In the third case, the function  $g(\cdot, \cdot)$  maps the step distance to a scalar reward, enabling  
 284 fine-grained credit assignment based on the similarity in quality. Specifically, we define a progress  
 285 value  $s()$  in the range  $[0, 1]$  to indicate the fraction of diffusion steps used, where 0 denotes zero  
 286 steps, and 1 denotes full completion of the steps. The ground-truth value is  $s(y_i)$ , and the model will  
 287 estimate a progress value. We define the reward function as  $g(o_i, y_i) = \delta \cdot (1 - |s(o_i) - s(y_i)|)$ .  
 288

289 This reward formulation enables the model to move beyond binary discrimination and perform fine-  
 290 grained analysis of video quality. By learning to associate subtle differences in generation steps  
 291 with quality variations, the model develops a deeper understanding of the diffusion process and its  
 292 impact on perceptual realism. As a result, it can not only detect whether a video is fake, but also  
 293 infer but also estimate the degree of quality degradation in generated videos. This facilitates more  
 294 interpretable and controllable evaluation of generated content quality.  
 295

## 4 EXPERIMENTS

### 4.1 IMPLEMENTATION DETAILS

#### 4.1.1 DATASET

299 Our dataset contains 140k videos, balanced between 70k real and 70k generated samples, organized  
 300 into contextual pairs. The real set comprises 55k videos from InternVid and 15k from ActivityNet,  
 301 while the generated set includes 50k samples synthesized by HunyuanVideo-I2V (Kong  
 302 et al., 2024) and 20k by CogVideoX-5B (Yang et al., 2024). We allocate 130k samples for  
 303 training and 10k for testing, with the latter evenly split between real and generated videos.  
 304 Within the training data, 30k samples are reserved for chain-of-thought (CoT) learning, denoted  
 305 as VidGuard-R1-CoT-30k, and the remaining 100k are used for reinforcement learning fine-  
 306 tuning, denoted as VidGuard-R1-RL-100k.  
 307

308 Since state-of-the-art generative models still produce relatively short videos ( $\sim$ 129 frames) at  
 309 modest resolutions, we standardize all real videos to match generated ones by enforcing 49 frames,  
 310 8 FPS,  $720 \times 480$  resolution, and YUV420p format.  
 311

312 For GRPO-Q fine-tuning, we augment the training set with intermediate generations sampled from  
 313 diffusion steps 10 to 50. These are labeled with approximate quality levels (20%, 40%, 60%, 80%,  
 314 and 95%). Specifically, we use 12k real videos, each paired with five generated variants at different  
 315 diffusion steps, resulting in 72k samples per generation model.  
 316

#### 4.1.2 EVALUATION PROTOCOL

317 We evaluate three datasets—ours, GenVidBench (Ni et al., 2025), and GenVideo (Chen et al.,  
 318 2024a)—using the metrics and baselines defined by their respective benchmarks. For ours and  
 319 GenVidBench, we report **mean Top-1 accuracy**, the average correctness over all predictions. For  
 320 GenVideo, we follow the original protocol and report **recall** and **F1 score**. All evaluations adhere to  
 321 the official settings of each benchmark to ensure fair comparison.  
 322

#### 4.1.3 TRAINING SETUP

323 We employ Qwen2.5-VL-7B as the base MLLM and conduct all experiments on four NVIDIA A100  
 324 GPUs (80GB). Each video is represented by up to 16 frames, where each frame is resized to a  $28 \times 28$   
 325 spatial resolution and mapped to 128 feature channels for encoder input during both training and  
 326

Table 1: Comparison of models on our dataset, reported as mean Top-1 accuracy (%). TF denotes transformer.

Method	Type	CogVideoX	HunyuanVideo
SlowFast	CNN	77.87	77.03
l3D	CNN	64.78	62.13
TRN	CNN	68.73	69.87
UniFormer V2	TF	73.95	71.92
TimeSformer	TF	78.53	74.55
VideoSwin	TF	76.81	79.71
MViT V2	TF	58.38	53.91
Qwen2.5-VL-7B	MLLM	50.95	52.83
<b>Qwen2.5-VL-72B</b>	<b>MLLM</b>	<b>54.17</b>	<b>55.82</b>
GPT-4.1 mini	MLLM	54.95	55.31
<b>GPT-4o</b>	<b>MLLM</b>	<b>56.81</b>	<b>57.42</b>
VidGuard-R1 (CoT)	MLLM	66.18	63.19
VidGuard-R1 (DPO)	MLLM	79.13	80.88
VidGuard-R1 (GRPO)	MLLM	81.30	81.90
VidGuard-R1 (GRPO-TA)	MLLM	82.17	83.72
VidGuard-R1 (GRPO-Q)	MLLM	<b>84.32</b>	<b>86.17</b>

Table 2: Extended GenVidBench results with **VidGuard-R1** and additional MLLMs, reported as mean Top-1 accuracy (%). TF denotes transformer.

Method	Type	MuseV	SVD	CogVideo	Mora	HD-VG	Mean
SlowFast	CNN	12.25	12.68	38.34	45.93	93.63	41.61
I3D		8.15	8.29	60.11	93.24	93.99	49.23
TRN		38.92	26.64	91.34	93.98	93.97	71.26
UniFormer V2	TF	20.05	14.81	45.21	99.21	96.89	57.55
TimeSformer	TF	73.14	20.17	74.80	39.40	92.32	64.28
VideoSwin	TF	62.29	8.01	91.82	45.83	<b>99.29</b>	67.27
MViT V2	TF	76.34	<b>98.29</b>	47.50	96.62	97.58	97.55
Qwen2.5-VL-7B	MLLM	25.86	27.06	68.51	43.26	71.15	47.30
GPT-4.1 mini	MLLM	26.07	33.78	94.07	57.19	87.64	59.62
VidGuard-R1 (CoT)	MLLM	36.52	16.02	99.35	76.94	99.94	66.09
VidGuard-R1 (GRPO, GenVideo-pretrained, Zero-shot)	MLLM	97.24	96.59	99.88	99.93	88.14	96.37
VidGuard-R1 (GRPO)	MLLM	<b>97.38</b>	94.98	<b>99.90</b>	<b>99.99</b>	95.46	<b>97.53</b>

inference. For GenVideo and GenVidBench, we follow their official evaluation protocols and adopt 8-frame inputs. In GRPO training, we sample 8 responses per input; for GRPO-TA, we additionally sample 4 responses from temporally manipulated variants of the input to enhance robustness against temporal artifacts. Training proceeds in two stages: first, the base model is fine-tuned for one epoch on the CoT dataset, yielding the SFT-CoT MLLM; second, we initialize **VidGuard-R1** with SFT-CoT and perform reinforcement learning for approximately 2,000 steps.

## 4.2 MAIN RESULTS

### 4.2.1 OUR DATASET

We evaluate **VidGuard-R1** on our dataset with several methods, including CNN-based models (SlowFast (Feichtenhofer et al., 2019), I3D (Carreira & Zisserman, 2017), TRN (Zhou et al., 2018)), Transformer-based models (UniFormer V2 Li et al. (2022a), TimeSformer (Bertasius et al., 2021), VideoSwin (Liu et al., 2022), MViT V2 (Li et al., 2022b)), and MLLM-based models (Qwen2.5-VL (Bai et al., 2025) and GPT-4.1 mini (OpenAI, 2025)). For CNN and Transformer models, we use the default training settings provided by the MMAAction2 framework (Contributors, 2020).

As shown in Table 1, CNN- and Transformer-based models achieved 53–79% accuracy, with Slow-Fast and TimeSformer among the top performers. In contrast, Qwen2.5-VL-7B and GPT-4.1 mini exhibited near-random performance, highlighting their limited capability in distinguishing fake videos. *VidGuard-R1 (CoT)*, trained via supervised fine-tuning (SFT) on Qwen2.5-VL-7B, substantially improved accuracy from around 51% to over 66%, yet remained less competitive compared to advanced SOTA methods. This result aligns with the intended role of the SFT stage—as a cold start phase to guide the model toward structured *think + answer* responses, emphasizing the extraction of salient visual cues.

In the subsequent RL stage, both DPO and GRPO further improved performance by roughly 2% over the best baseline. Our proposed methods—GRPO-TA and GRPO-Q—achieved additional gains of approximately 2% and 5% over GRPO, respectively, demonstrating the effectiveness of temporal artifact supervision and quality-aware reward modeling in enhancing detection accuracy.

### 4.2.2 GENVIDBENCH BENCHMARK

The GenVidBench dataset comprises approximately 87k training samples and 82k testing samples, with fake videos generated by models such as MuseV (Xia et al., 2024), SVD (Blattmann et al., 2023), CogVideo (Hong et al., 2022), and Mora (Yuan et al., 2024), and real videos sourced from HD-VG (Wang et al., 2023b). We conduct training and evaluation under the cross-source and cross-generator settings as proposed in their benchmark. In addition to the models originally reported in GenVidBench, we evaluate **VidGuard-R1** using the same model families as in our dataset experiments—CNN-based, Transformer-based, and MLLM-based models—including two MLLMs: Qwen2.5-VL and GPT-4.1 mini. **VidGuard-R1** (GRPO, GenVideo-pretrained, Zero-shot) denotes the zero-shot model pretrained on GenVideo and evaluated on GenVidBench. As shown in Table 2, both the zero-shot model and two fine-tuned variants achieve over 15% higher accuracy compared to prior SOTA methods. Notably, the zero-shot model demonstrates strong generalization, highlighting the effectiveness of pretraining on diverse generative content. Complete detection model results are provided in Appendix B.

378  
 379 Table 3: Extended GenVideo results with **VidGuard-R1** and additional MLLMs, evaluated by F1  
 380 score and recall (R)

Method	Detection level	Metric	Sora	Morph Studio	Gen2	HotShot	Lavie	Show-1	Moon Valley	Crafter	Model Scope	Wild Scrape	Mean
NPR (Tan et al., 2024)	Image	R F1	0.91 0.27	0.99 0.84	0.99 0.91	0.24 0.30	0.89 0.86	0.57 0.59	0.97 0.81	0.99 0.91	0.94 0.81	0.87 0.81	0.84 0.71
VideoMAE (Tong et al., 2022)	Video	R F1	0.67 0.62	0.96 0.95	0.98 0.98	0.96 <b>0.96</b>	0.77 0.86	0.80 0.87	0.97 0.96	0.96 0.97	0.96 <b>0.96</b>	0.68 0.79	0.87 0.89
MINTIME-CLIP (Cocomini et al., 2024)	Video	R F1	0.89 0.49	1.00 0.93	0.98 0.96	0.26 0.37	0.96 0.94	0.98 0.92	0.99 0.92	1.00 0.96	0.84 0.84	0.82 0.85	0.87 0.82
FTCN-CLIP (Zheng et al., 2021)	Video	R F1	0.87 0.78	1.00 0.98	0.98 0.98	0.17 0.29	0.97 0.98	0.91 0.94	1.00 0.98	1.00 <b>0.99</b>	0.85 0.90	0.82 0.89	0.86 0.87
DeMamba-XCLIP (Chen et al., 2024a)	Video	R F1	0.98 0.64	1.00 0.96	0.99 0.97	0.65 0.75	0.94 <b>0.95</b>	0.98 0.95	1.00 0.95	1.00 0.97	0.92 0.92	0.89 <b>0.91</b>	0.93 0.90
Qwen2.5-VL-7B (Bai et al., 2025)	MLLM	R F1	0.58 0.74	0.56 0.72	0.54 0.49	0.33 0.60	0.43 0.55	0.38 0.90	0.81 0.77	0.63 0.77	0.51 0.68	0.70 0.82	0.54 0.70
GPT-4.1 mini (OpenAI, 2025)	MLLM	R F1	0.43 0.60	0.67 0.80	0.56 0.72	0.54 0.70	0.63 0.77	0.56 0.72	0.92 0.96	0.67 0.80	0.69 0.82	0.69 0.72	0.65 0.72
VidGuard-R1 (CoT)	MLLM	R F1	0.92 0.90	0.89 0.91	0.91 0.95	0.90 0.89	0.98 0.99	0.79 0.81	0.99 0.95	0.85 0.89	0.89 0.85	0.87 0.88	0.90 0.90
VidGuard-R1 (GRPO, GenVidBench-pretrained, Zero-shot)	MLLM	R F1	0.95 0.93	0.98 0.93	0.90 0.96	0.89 0.91	0.97 0.99	0.85 0.82	0.99 0.95	0.93 0.89	0.81 0.85	0.87 0.88	0.92 0.91
VidGuard-R1 (GRPO)	MLLM	R F1	0.95 0.97	1.00 0.99	0.98 0.99	0.94 0.91	0.98 <b>0.99</b>	0.95 0.89	0.97 <b>0.99</b>	0.99 0.95	0.94 0.95	0.91 0.90	0.96 <b>0.96</b>

#### 394 4.2.3 GENVIDEO BENCHMARK

395 The GenVideo dataset comprises approximately 2.2M training samples and 20k testing samples, with generated videos sourced from a diverse set of models, including Sora (OpenAI, 2024), MorphStudio (mor, 2025), Gen2 (Esser et al., 2023b), HotShot (hot, 2025), Lavie (Wang et al., 2025), Show-1 (Zhang et al., 2024a), MoonValley (moo, 2025), Crafter (Chen et al., 2023), ModelScope (Wang et al., 2023a), and WildScrape (Wei et al., 2024). Following the official evaluation protocol, we benchmark two MLLM baselines and three variants of **VidGuard-R1**. Among these, VidGuard-R1 (GRPO) consistently outperforms almost all prior detection methods across videos generated by the various models. As shown in Table 3, it achieves an F1 score improvement of 0.06 compared to DeMamba-XCLIP. Complete detection model results are provided in Appendix B.

#### 404 4.2.4 PERFORMANCE GAP BETWEEN OUR DATASET AND BENCHMARKS

405 While **VidGuard-R1** achieves approximately 85% accuracy on our dataset, it obtains significantly 406 higher accuracy—exceeding 95%—on the two benchmark datasets. This discrepancy arises from 407 two key differences. First, the benchmarks exhibit clear discrepancies in video metadata—such as 408 resolution, duration, and frame rate—between real and fake videos, which models can exploit as 409 superficial cues. In contrast, we standardize all videos in our dataset by matching resolution, FPS, and 410 format, thereby forcing models to rely on actual visual content. Second, our dataset ensures strong 411 contextual alignment by conditioning generation on the first frame and the corresponding caption of 412 a real video, resulting in more realistic and semantically consistent outputs. In comparison, 413 benchmark datasets often generate fake videos from unrelated prompts and synthetic images, leading to 414 artifacts that make detection easier.

#### 415 4.2.5 ABLATION STUDY

416 **Explanation quality and accuracy comparison.** Table 4 417 presents results on the HunyuanVideo (Kong et al., 2024) 418 and CogVideoX (Yang et al., 2024) datasets. We report 419 explanation quality scores, which are rated on a 1–10 scale 420 (with 10 indicating excellent quality and full alignment) 421 by GPT-4.1 mini using the LLM-as-a-Judge prompt 422 described in Appendix D. Compared to baseline models such 423 as Qwen2.5-VL-7B and GPT-4.1 mini, our VidGuard-R1 424 GRPO variants achieve consistent improvements in both 425 classification accuracy and explanation quality.

426 **GRPO-TA reward ablation.** Table 5 reports an ablation 427 study of GRPO-TA on our dataset by varying the weight 428 parameters  $\alpha_1$  and  $\alpha_2$ , which control the relative importance of 429 different temporal artifact types. The highest classification 430 accuracy of 83.57% is achieved with  $\alpha_1 = 0.5$  and  $\alpha_2 = 0.3$ , 431 while the threshold  $\mu$  is fixed at 0.8 across all experiments.

432 Table 4: LLM-as-a-Judge explanation 433 scores on our dataset

Method	Expl. Score (Hunyuan Video)	Expl. Score (CogVideoX)
Qwen2.5-VL-7B	5.8	5.6
GPT-4.1 mini	5.8	5.9
VidGuard-R1 (CoT)	6.8	6.9
VidGuard-R1 (DPO)	7.2	8.1
VidGuard-R1 (GRPO)	8.1	8.0
VidGuard-R1 (GRPO-TA)	8.1	8.4
VidGuard-R1 (GRPO-Q)	<b>8.3</b>	<b>8.5</b>

434 Table 5: Accuracy (%) for **GRPO-TA** under different reward function 435 parameters  $\alpha_1$  and  $\alpha_2$

$\alpha_1$	$\alpha_2$	Accuracy (%)
0.3	0.1	81.31
0.3	0.3	82.59
<b>0.5</b>	<b>0.3</b>	<b>83.57</b>
0.5	0.5	83.12
0.7	0.5	82.53

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485Table 6: Accuracy (%) for **GRPO-Q** with varying number of intermediate diffusion steps

# of steps (step numbers)	Accuracy (%)
1 (50)	81.63
3 (10, 30, 50)	83.21
<b>5 (10, 20, 30, 40, 50)</b>	<b>85.80</b>

Table 7: Cross-dataset evaluation for **VidGuard-R1**

Test Dataset	Training Source	Accuracy (%)
Ours	<b>VidGuard-R1</b> (Ours only)	81.65
Ours	<b>VidGuard-R1</b> (Ours + GenVideo)	82.97
GenVideo	<b>VidGuard-R1</b> (GenVideo only)	97.53
GenVideo	<b>VidGuard-R1</b> (Ours + GenVideo)	97.98

Table 8: Label–rationale consistency and explanation quality for 20 videos

Metric	Value
Annotators	5
Label–rationale alignment	89%
Rationale score (0–5)	3.9

Table 9: Human ranking of explanation quality

Model	Avg. Rank
<b>VidGuard-R1 (GRPO)</b>	1.67
GPT-4o	2.08
Qwen2.5-VL-72B	2.22

**GRPO-Q reward ablation.** Table 6 presents an ablation study on GRPO-Q conducted on our dataset by varying the number of intermediate diffusion steps included per real video during fine-tuning. Using more steps provides richer supervision of video quality progression, improving detection accuracy. The best accuracy of 84.05% is obtained with five steps, which is the setting used in our main experiments.

**Cross-dataset complementarity.** To assess whether training on a limited generative source induces overfitting, we conduct dataset-mixing experiments using **VidGuard-R1** (GRPO). As shown in Table 7, augmenting our dataset with GenVideo leads to consistent performance gains across both evaluation sets, suggesting that the model benefits from heterogeneous training data and does not rely on artifacts from any single source. These findings indicate that incorporating diverse generative sources enhances overall accuracy, reinforcing that **VidGuard-R1** learns generalizable detection cues rather than dataset-specific patterns.

#### 4.3 HUMAN EVALUATION OF EXPLANATION QUALITY

To examine the coherence and interpretability of CoT rationales, we conducted a human evaluation of **VidGuard-R1** (GRPO). Tables 8 and 9 summarize the two complementary studies.

**Consistency and quality.** Five annotators evaluated twenty randomly selected fake videos that the model correctly identified. Annotators judged whether each rationale was consistent with the predicted label and assigned a quality score on a 0–5 scale after watching the corresponding video. As shown in Table 8, annotators reported 89% label–rationale agreement with an average quality score of 3.9. Lower scores ( $\leq 2$ ) occurred primarily when the model emphasized subtle texture-level artifacts that were difficult for human raters to perceive.

**Explanation quality ranking.** We conducted a blind ranking study in which five participants evaluated explanations from three models across twenty videos, consisting of ten real and ten synthetic videos. Table 9 shows that **VidGuard-R1** achieved the strongest human preference, with the lowest average rank (1.67), outperforming GPT-4o and Qwen2.5-VL-72B. Participants consistently judged **VidGuard-R1**’s explanations to be more informative and higher quality. This improvement is likely because, although GPT-4o and Qwen2.5-VL-72B also describe video features well, their explanations are less effective at distinguishing real from synthetic content, whereas **VidGuard-R1** better highlights the discriminative aspects relevant to this task.

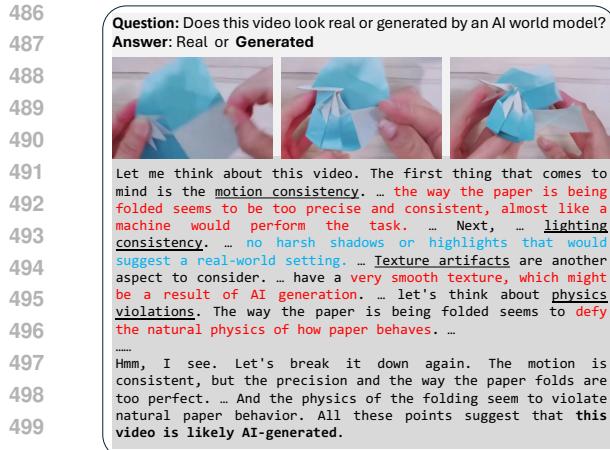


Figure 3: VidGuard-R1 (GRPO): reasoning about an origami folding sequence

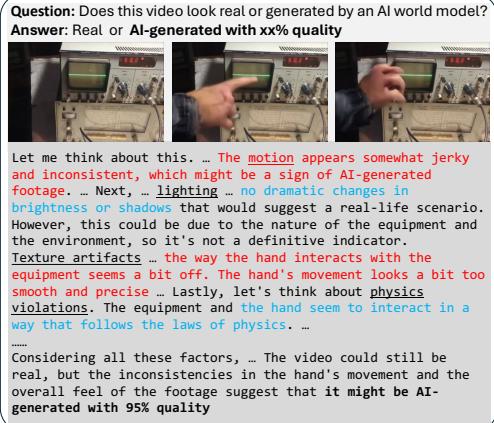


Figure 4: VidGuard-R1 (GRPO-Q): temporal inconsistency observed as unnatural movements

504 **Interpretation and broader implications.** Although the CoT rationales align closely with human  
 505 judgments, they are not guaranteed to be correct in every case. Hallucinations may arise, and the  
 506 rationales should therefore be shown with an appropriate disclaimer noting their AI-generated nature.  
 507 Nevertheless, presenting interpretable cues helps users more effectively assess whether a video may  
 508 be AI-generated, supporting a practical human-in-the-loop verification process.

#### 509 4.4 CASE STUDIES ON EXPLANATIONS

511 Figures 3 and 4 illustrate cases where **VidGuard-R1** correctly identifies videos as generated. The  
 512 model performs multi-faceted reasoning across motion, lighting, texture, and physical plausibility  
 513 before arriving at a final decision. Notably, it does not rely on a single cue, but instead accumulates  
 514 evidence across frames, resembling how humans distinguish fake videos. In each figure, pink highlights  
 515 denote cues suggesting realism, red indicates artifacts indicative of generation, yellow marks  
 516 intermediate reasoning steps, and underlines represent several key factors.

517 For instance, in Figure 3, the smooth hand motion initially suggests realism; however, once the  
 518 origami folds in a physically implausible manner, the model revises its judgment. In Figure 4,  
 519 although the lighting and shadows are consistent—typically a cue for authenticity—the model  
 520 recognizes that this is insufficient in a largely static scene with only a stationary machine and a human  
 521 hand. In particular, even in its final prediction, the model reflects on earlier realistic cues and ac-  
 522 knowledges that *the video could still be real*, underscoring its nuanced, human-like reasoning in  
 523 assessing video quality. Additional case studies are provided in the Appendix E.

## 524 5 CONCLUSION

526 We propose **VidGuard-R1**, an MLLM-based discriminator that not only detects AI-generated  
 527 videos with high accuracy but also provides interpretable reasoning. By leveraging reinforcement  
 528 learning with reward models targeting temporal artifacts and generation quality, **VidGuard-R1**  
 529 achieves 85% accuracy on our dataset, 97% on GenVidBench, and 96% on GenVideo, substan-  
 530 tially surpassing prior state-of-the-art methods. We expect this work to advance MLLMs' reasoning.  
 531

### 532 5.1 LIMITATIONS

534 Our dataset currently includes fake videos generated using HunyuanVideo and CogVideoX, which  
 535 are the primary open-source models supporting large-scale text–image joint conditioning. Other  
 536 diffusion models provide only text- or image-based conditioning, limiting their suitability for our  
 537 pairwise real–fake construction. While the current design ensures strong contextual alignment be-  
 538 tween real and generated videos, incorporating outputs from a broader set of generative models  
 539 would further improve diversity and robustness, thereby enhancing applicability to real-world sce-  
 540 narios.

540  
541 ETHICS STATEMENT542  
543 This work does not involve personally identifiable information or sensitive user data. All datasets  
544 used in our experiments are publicly available and constructed in accordance with their licenses and  
545 usage guidelines. VidGuard-R1 is designed to mitigate societal risks associated with AI-generated  
546 videos, such as misinformation and reputational harm, by providing interpretable CoT reasoning to  
547 assist human verification. To the best of our knowledge, the method does not introduce risks related  
548 to fairness, safety, or privacy.549  
550 REPRODUCIBILITY STATEMENT551  
552 The VidGuard-R1 framework is detailed, with a two-stage training process: SFT for CoT initialization  
553 followed by RL fine-tuning using DPO and GRPO on the Qwen2.5-VL-7B model. The training  
554 involves the VidGuard-R1-CoT-30k and VidGuard-R1-RL-100k datasets, with videos standardized  
555 to  $720 \times 480$  resolution and 8 FPS. Training proceeds for one SFT epoch, followed by  $\sim 2,000$   
556 RL steps on four NVIDIA A100 GPUs. Reproducibility is further ensured by specifying the reward  
557 strategies: GRPO-TA and GRPO-Q.558  
559 REFERENCES560 Hotshot. Online, 2025. Available: <https://huggingface.co/hotshotco/Hotshot-XL/>.562 Luma Ray 2. <https://lumalabs.ai/ray2>, 2025. Accessed: 2025-11-30.564 Moonvalley. Online, 2025. Available: <https://moonvalley.ai/>.566 Morph studio. Online, 2025. Available: <https://www.morphstudio.com/>.568 Pika. <https://pika.art/>, 2025. Accessed: 2025-11-30.569 Gen-3 alpha. <https://runwayml.com/research/introducing-gen-3-alpha>,  
570 2025. Accessed: 2025-11-30.572 Openai sora. <https://openai.com/sora/>, 2025. Accessed: 2025-11-30.573 Google deepmind veo. <https://deepmind.google/models/veo/>, 2025. Accessed:  
574 2025-11-30.576 Wan video generation. <https://wan.video/>, 2025. Accessed: 2025-11-30.578 Jianfa Bai, Man Lin, and Gang Cao. Ai-generated video detection via spatio-temporal anomaly  
579 learning. *arXiv preprint arXiv:2403.16638*, 2024a.580 Jianfa Bai, Man Lin, Gang Cao, and Zijie Lou. Ai-generated video detection via spatial-temporal  
581 anomaly learning. In *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*,  
582 pp. 460–470. Springer, 2024b.584 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang  
585 Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, local-  
586 ization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.587 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,  
588 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan,  
589 Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng,  
590 Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv*  
591 *preprint arXiv:2502.13923*, 2025.592 Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video  
593 understanding?, 2021.

594 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik  
 595 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rom-  
 596 bach. Stable video diffusion: Scaling latent video diffusion models to large datasets, 2023.  
 597

598 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe  
 599 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video  
 600 generation models as world simulators. 2024. URL <https://openai.com/research/video-generation-models-as-world-simulators>.

601

602 Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet:  
 603 A large-scale video benchmark for human activity understanding. In *Proceedings of the ieee  
 604 conference on computer vision and pattern recognition*, pp. 961–970, 2015.

605

606 J. Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics  
 607 dataset. In *CVPR*, pp. 4724–4733, 07 2017. doi: 10.1109/CVPR.2017.502.

608

609 Chirui Chang, Zhengze Liu, Xiaoyang Lyu, and Xiaojuan Qi. What matters in detecting ai-  
 610 generated videos like sora? *arXiv preprint arXiv:2406.19568*, 2024.

611

612 Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing,  
 613 Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models for high-  
 614 quality video generation. *arXiv preprint arXiv:2310.19512*, 2023.

615

616 Haoxing Chen, Yan Hong, Zizheng Huang, Zhuoer Xu, Zhangxuan Gu, Yaohui Li, Jun Lan, Huijia  
 617 Zhu, Jianfu Zhang, Weiqiang Wang, et al. Demamba: Ai-generated video detection on million-  
 618 scale genvideo benchmark. *arXiv preprint arXiv:2405.19707*, 2024a.

619

620 Zhaorun Chen, Francesco Pinto, Minzhou Pan, and Bo Li. Safewatch: An efficient safety-policy  
 621 following video guardrail model with transparent explanations. *arXiv preprint arXiv:2412.06878*,  
 622 2024b.

623

624 Davide Alessandro Coccolini, Giorgos Kordopatis Zilos, Giuseppe Amato, Roberto Caldelli, Fab-  
 625 rizio Falchi, Symeon Papadopoulos, and Claudio Gennaro. Mintime: Multi-identity size-invariant  
 626 video deepfake detection. *IEEE Transactions on Information Forensics and Security*, 19:6084–  
 627 6096, 2024.

628

629 MMAAction2 Contributors. Openmmlab’s next generation video understanding toolbox and bench-  
 630 mark. <https://github.com/open-mmlab/mmaction2>, 2020.

631

632 Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germani-  
 633 dis. Structure and content-guided video synthesis with diffusion models. In *Proceedings of the  
 634 IEEE/CVF international conference on computer vision*, pp. 7346–7356, 2023a.

635

636 Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germani-  
 637 dis. Structure and content-guided video synthesis with diffusion models. In *Proceedings of the  
 638 IEEE/CVF international conference on computer vision*, pp. 7346–7356, 2023b.

639

640 Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition, 2020. URL  
 641 <https://arxiv.org/abs/2004.04730>.

642

643 Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video  
 644 recognition. In *Proceedings of the IEEE international conference on computer vision*, pp. 6202–  
 645 6211, 2019.

646

647 Kaituo Feng, Kaixiong Gong, Bohao Li, Zonghao Guo, Yibing Wang, Tianshuo Peng, Junfei Wu,  
 648 Xiaoying Zhang, Benyou Wang, and Xiangyu Yue. Video-r1: Reinforcing video reasoning in  
 649 mllms. *arXiv preprint arXiv:2503.21776*, 2025.

650

651 Zhihao Gu, Yang Chen, Taiping Yao, Shouhong Ding, Jilin Li, Feiyue Huang, and Lizhuang Ma.  
 652 Spatiotemporal inconsistency learning for deepfake video detection. In *Proceedings of the 29th  
 653 ACM international conference on multimedia*, pp. 3473–3481, 2021.

648 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
 649 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms  
 650 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

651

652 Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pre-  
 653 training for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.

654

655 Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li,  
 656 Bo Wu, Jianwei Zhang, et al. Hunyuandvideo: A systematic framework for large video generative  
 657 models. *arXiv preprint arXiv:2412.03603*, 2024.

658

659 Rohit Kundu, Hao Xiong, Vishal Mohanty, Athula Balachandran, and Amit K Roy-Chowdhury.  
 660 Towards a universal synthetic video detector: From face or background manipulations to fully ai-  
 661 generated content. In *Proceedings of the Computer Vision and Pattern Recognition Conference*,  
 662 pp. 28050–28060, 2025.

663

664 Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Limin Wang, and Yu Qiao. Uni-  
 665 formerv2: Spatiotemporal learning by arming image vits with video uniformer. *arXiv preprint  
 666 arXiv:2211.09552*, 2022a.

667

668 Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and  
 669 Christoph Feichtenhofer. Mvitzv2: Improved multiscale vision transformers for classification and  
 670 detection. In *CVPR*, 2022b.

671

672 Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding,  
 673 2019. URL <https://arxiv.org/abs/1811.08383>.

674

675 Qingyuan Liu, Pengyuan Shi, Yun-Yun Tsai, Chengzhi Mao, and Junfeng Yang. Turns out i'm not  
 676 real: Towards robust detection of ai-generated videos. *arXiv preprint arXiv:2406.09601*, 2024.

677

678 Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin trans-  
 679 former. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
 680 pp. 3202–3211, 2022.

681

682 Long Ma, Jiajia Zhang, Hongping Deng, Ningyu Zhang, Yong Liao, and Haiyang Yu. Decof:  
 683 Generated video detection via frame consistency, 2024.

684

685 Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. Use of a capsule network to detect fake  
 686 images and videos, 2019. URL <https://arxiv.org/abs/1910.12467>.

687

688 Bolin Ni, Houwen Peng, Minghao Chen, Songyang Zhang, Gaofeng Meng, Jianlong Fu, Shiming  
 689 Xiang, and Haibin Ling. Expanding language-image pretrained models for general video recog-  
 690 nition. In *European conference on computer vision*, pp. 1–18. Springer, 2022.

691

692 Zhenliang Ni, Qiangyu Yan, Mouxiao Huang, Tianning Yuan, Yehui Tang, Hailin Hu, Xinghao  
 693 Chen, and Yunhe Wang. Genvidbench: A challenging benchmark for detecting ai-generated  
 694 video. *arXiv preprint arXiv:2501.11340*, 2025.

695

696 OpenAI. Sora. Online, 2024. Available: <https://openai.com/index/sora/>.

697

698 OpenAI. Gpt-4.1. Online, 2025. Available: <https://openai.com/index/gpt-4-1/>.

699

700 Gan Pei, Jiangning Zhang, Menghan Hu, Zhenyu Zhang, Chengjie Wang, Yunsheng Wu, Guang-  
 701 tao Zhai, Jian Yang, Chunhua Shen, and Dacheng Tao. Deepfake generation and detection: A  
 702 benchmark and survey. *arXiv preprint arXiv:2403.17881*, 2024.

703

704 Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in frequency: Face  
 705 forgery detection by mining frequency-aware clues, 2020a. URL <https://arxiv.org/abs/2007.09355>.

706

707 Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in frequency: Face  
 708 forgery detection by mining frequency-aware clues. In *European conference on computer vision*,  
 709 pp. 86–103. Springer, 2020b.

702 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
 703 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
 704 models from natural language supervision. In *International conference on machine learning*, pp.  
 705 8748–8763. PmLR, 2021.

706 Hao Shao, Shengju Qian, and Yu Liu. Temporal interlacing network. *AAAI*, 2020.

708 Chuangchuang Tan, Yao Zhao, Shikui Wei, Guanghua Gu, Ping Liu, and Yunchao Wei. Rethinking  
 709 the up-sampling operations in cnn-based generative network for generalizable deepfake detection.  
 710 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
 711 28130–28139, 2024.

712 Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-  
 713 efficient learners for self-supervised video pre-training. *Advances in neural information process-  
 714 ing systems*, 35:10078–10093, 2022.

715 Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Mod-  
 716 ellscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023a.

717 Wenjing Wang, Huan Yang, Zixi Tuo, Huiguo He, Junchen Zhu, Jianlong Fu, and Jiaying Liu.  
 718 Videofactory: Swap attention in spatiotemporal diffusions for text-to-video generation. *arXiv  
 719 preprint arXiv:2305.10874*, 2023b.

720 Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan  
 721 He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent  
 722 diffusion models. *International Journal of Computer Vision*, 133(5):3059–3078, 2025.

723 Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan  
 724 Chen, Yaohui Wang, et al. Internvid: A large-scale video-text dataset for multimodal under-  
 725 standing and generation. *arXiv preprint arXiv:2307.06942*, 2023c.

726 Yujie Wei, Shiwei Zhang, Zhiwu Qing, Hangjie Yuan, Zhiheng Liu, Yu Liu, Yingya Zhang, Jingren  
 727 Zhou, and Hongming Shan. Dreamvideo: Composing your dream videos with customized sub-  
 728 ject and motion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
 729 Recognition*, pp. 6537–6549, 2024.

730 Zhiqiang Xia, Zhaokang Chen, Bin Wu, Chao Li, Kwok-Wai Hung, Chao Zhan, Yingjie He, and  
 731 Wenjiang Zhou. Musev: Infinite-length and high fidelity virtual human video generation with  
 732 visual conditioned parallel denoising. *arxiv*, 2024.

733 Yuting Xu, Jian Liang, Gengyun Jia, Ziming Yang, Yanhao Zhang, and Ran He. Tall: Thumbnail  
 734 layout for deepfake video detection. In *Proceedings of the IEEE/CVF International Conference  
 735 on Computer Vision*, pp. 22658–22668, 2023.

736 Zhipei Xu, Xuanyu Zhang, Runyi Li, Zecheng Tang, Qing Huang, and Jian Zhang. Fakeshield:  
 737 Explainable image forgery detection and localization via multi-modal large language models.  
 738 *arXiv preprint arXiv:2410.02761*, 2024.

739 Ceyuan Yang, Yinghao Xu, Jianping Shi, Bo Dai, and Bolei Zhou. Temporal pyramid network  
 740 for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern  
 741 Recognition (CVPR)*, 2020.

742 Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,  
 743 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models  
 744 with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.

745 Zhengqing Yuan, Ruoxi Chen, Zhaoxu Li, Haolong Jia, Lifang He, Chi Wang, and Lichao Sun.  
 746 Mora: Enabling generalist video generation via a multi-agent framework, 2024.

747 David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei  
 748 Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-  
 749 video generation. *International Journal of Computer Vision*, pp. 1–15, 2024a.

750

756 Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video  
757 instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024b.  
758

759 Yinglin Zheng, Jianmin Bao, Dong Chen, Ming Zeng, and Fang Wen. Exploring temporal coherence  
760 for more general video face forgery detection. In *Proceedings of the IEEE/CVF international*  
761 *conference on computer vision*, pp. 15044–15054, 2021.

762 Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. Temporal relational reasoning in  
763 videos. *European Conference on Computer Vision*, 2018.  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

810  
 811 Table 10: Extended GenVidBench results with **VidGuard-R1** and additional MLLMs, reported as  
 812 mean Top-1 accuracy (%). TF denotes transformer.

Method	Type	MuseV	SVD	CogVideo	Mora	HD-VG	Mean
SlowFast (Feichtenhofer et al., 2019)	CNN	12.25	12.68	38.34	45.93	93.63	41.66
F3Net (Qian et al., 2020a)	CNN	37.43	37.27	36.46	39.59	52.76	42.52
I3D (Carreira & Zisserman, 2017)	CNN	8.15	8.29	60.11	59.24	93.99	49.23
CFV2 (Nguyen et al., 2019)	CNN	86.26	86.53	10.10	16.90	88.40	60.53
TPN (Yang et al., 2020)	CNN	37.86	8.79	68.25	90.04	97.34	61.52
TIN (Shao et al., 2020)	CNN	33.78	21.47	81.59	79.44	97.88	63.97
TRN (Zhou et al., 2018)	CNN	38.92	26.64	91.34	93.98	93.97	71.26
TSM (Lin et al., 2019)	CNN	70.37	54.70	78.46	70.37	96.76	76.40
X3D (Feichtenhofer, 2020)	CNN	92.39	37.27	65.72	49.60	97.51	77.09
UniFormer V2 (Li et al., 2022a)	TF	20.05	14.81	45.21	99.21	96.89	57.55
TimeSformer (Bertasius et al., 2021)	TF	73.14	20.17	74.80	39.40	92.32	64.28
VideoSwin (Liu et al., 2022)	TF	62.29	8.01	91.82	45.83	<b>99.29</b>	67.27
MViT V2 (Li et al., 2022b)	TF	76.34	<b>98.29</b>	47.50	96.62	97.58	79.90
Qwen2.5-VL-7B (Bai et al., 2025)	MLLM	25.86	27.06	68.51	43.26	71.15	47.30
GPT-4.1 mini (OpenAI, 2025)	MLLM	26.07	33.78	94.07	57.19	87.64	59.62
VidGuard-R1 (CoT)	MLLM	36.52	16.02	99.35	76.94	99.94	66.09
VidGuard-R1 (GRPO, GenVideo-pretrained, Zero-shot)	MLLM	97.24	96.59	99.88	99.93	88.14	96.37
VidGuard-R1 (GRPO)	MLLM	<b>97.38</b>	94.98	<b>99.90</b>	<b>99.99</b>	95.46	<b>97.53</b>

## A ADDITIONAL SETUP

833  
 834 To further guide the model during RL training, we incorporate a length-based reward strategy. We  
 835 promote informative yet concise reasoning by rewarding outputs that are neither too brief nor exces-  
 836 sively long. Specifically, if the model predicts the correct answer and the length of the response falls  
 837 within the range  $[l_{\min}, l_{\max}]$ , an additional reward  $\omega$  is assigned. Let  $l_i$  be the length of the model’s  
 838 response for the  $i$ -th video. The reward is defined as:

$$r_i^{\text{total}} = \begin{cases} r_i + \omega, & \text{if } o_i \text{ is correct and } l_{\min} \leq l_i \leq l_{\max} \\ r_i, & \text{otherwise} \end{cases} \quad (4)$$

839 where we set  $\omega = 0.1$ ,  $l_{\min} = 320$ , and  $l_{\max} = 512$ .

## B COMPREHENSIVE BENCHMARK EVALUATION

840  
 841 In this section, we provide extended benchmark results for **VidGuard-R1** alongside additional  
 842 MLLMs. Table 10 presents mean Top-1 accuracy on GenVidBench across multiple video datasets,  
 843 including CNN and Transformer baselines as well as selected MLLM variants. Table 11 reports  
 844 comprehensive F1 and recall scores on the GenVideo dataset, including all models provided in the  
 845 official benchmark alongside our MLLM variants. These extended tables offer a complete compari-  
 846 son of performance across all evaluated models.

## C ZERO-SHOT GENERALIZATION TO UNSEEN GENERATIVE MODELS

850  
 851 To evaluate the robustness of **VidGuard-R1** beyond the curated training sources, we assess its zero-  
 852 shot performance on a diverse set of recently released generative video models that were not used  
 853 during training, including Gen-3 Alpha (run, 2025), Pika (pik, 2025), Pika 2.2 (pik, 2025), Luma  
 854 Ray2 (lum, 2025), Sora (sor, 2025), Veo2 (veo, 2025), Veo3 (veo, 2025), and Wan 2.1 (wan, 2025).  
 855 Table 12 summarizes performance across these unseen systems. VidGuard-R1 achieves accuracy  
 856 above 80% in all cases, reaching up to 96.36%, demonstrating strong generalization to more recent  
 857 and increasingly realistic generative models.

858  
 859 These analyses demonstrate that **VidGuard-R1** generalizes effectively beyond the curated genera-  
 860 tive sources and remains robust across a wide range of unseen, high-quality video generation models.

864  
 865 Table 11: Extended GenVideo results with **VidGuard-R1** and additional MLLMs, evaluated by F1  
 866 and recall scores

Model	Detection level	Metric	Sora	Morph Studio	Gen2	HotShot	Lavie	Show-1	Moon Valley	Crafter	Model Scope	Wild Scrape	Mean
F3Net (Qian et al., 2020a)	Image	R F1	0.83 0.50	0.99 0.94	0.98 0.96	0.77 0.81	0.57 0.69	0.36 0.49	0.99 0.93	0.99 0.96	0.89 0.88	0.76 0.82	0.81 0.80
NPR (Tan et al., 2024)	Image	R F1	0.91 0.27	0.99 0.84	0.99 0.91	0.24 0.30	0.89 0.86	0.57 0.59	0.97 0.81	0.99 0.91	0.94 0.81	0.87 0.81	0.84 0.71
STIL (Gu et al., 2021)	Video	R F1	0.78 0.38	0.98 0.90	0.98 0.94	0.76 0.78	0.61 0.72	0.53 0.62	0.99 0.90	0.97 0.94	0.94 0.88	0.65 0.72	0.82 0.78
VideoMAE (Tong et al., 2022)	Video	R F1	0.67 0.62	0.96 0.95	0.98 0.98	0.96 0.96	0.77 0.86	0.80 0.87	0.97 0.96	0.96 0.97	0.96 0.97	0.68 0.79	0.87 0.89
MINTIME-CLIP (Cocomomi et al., 2024)	Video	R F1	0.89 0.49	1.00 0.93	0.98 0.96	0.26 0.37	0.96 0.94	0.98 0.92	0.99 0.92	1.00 0.96	0.84 0.86	0.82 0.84	0.87 0.82
FTCN-CLIP (Zheng et al., 2021)	Video	R F1	0.87 0.78	1.00 0.98	0.98 0.98	0.17 0.29	0.97 0.98	0.91 0.94	1.00 0.98	1.00 0.99	0.85 0.99	0.82 0.90	0.86 0.87
TALL (Xu et al., 2023)	Video	R F1	0.91 0.26	0.98 0.82	0.97 0.89	0.83 0.74	0.76 0.77	0.79 0.72	0.99 0.81	0.98 0.90	0.94 0.80	0.66 0.67	0.88 0.74
CLIP (Radford et al., 2021)	Image	R F1	0.94 0.28	0.99 0.84	0.91 0.86	0.77 0.72	0.88 0.85	0.86 0.76	0.99 0.82	0.99 0.91	0.84 0.76	0.84 0.79	0.90 0.76
DeMamba-CLIP (Chen et al., 2024a)	Video	R F1	0.95 0.64	1.00 0.96	0.98 0.97	0.69 0.78	0.92 0.94	0.93 0.92	1.00 0.95	1.00 0.98	0.83 0.98	0.82 0.87	0.91 0.89
XCLIP (Ni et al., 2022)	Video	R F1	0.82 0.31	0.99 0.88	0.93 0.90	0.61 0.65	0.79 0.82	0.69 0.70	0.97 0.86	0.99 0.93	0.77 0.75	0.83 0.82	0.84 0.76
DeMamba-XCLIP (Chen et al., 2024a)	Video	R F1	0.98 0.64	1.00 0.96	0.99 0.97	0.65 0.75	0.94 0.95	0.98 0.95	1.00 0.95	1.00 0.97	0.92 0.92	0.89 0.91	0.93 0.90
Qwen2.5-VL-7B (Bai et al., 2025)	MLLM	R F1	0.58 0.74	0.56 0.72	0.54 0.70	0.33 0.49	0.43 0.60	0.38 0.55	0.81 0.90	0.63 0.77	0.51 0.68	0.70 0.82	0.54 0.70
GPT-4.1 mini (OpenAI, 2025)	MLLM	R F1	0.43 0.60	0.67 0.80	0.56 0.72	0.54 0.70	0.63 0.77	0.56 0.72	0.92 0.96	0.67 0.80	0.69 0.82	0.69 0.82	0.65 0.72
VidGuard-R1 (CoT)	MLLM	R F1	0.92 0.90	0.89 0.91	0.91 0.95	0.90 0.89	0.98 0.99	0.79 0.81	0.99 0.95	0.85 0.89	0.89 0.85	0.87 0.88	0.90 0.90
VidGuard-R1 (GRPO, GenVidBench-pretrained, Zero-shot)	MLLM	R F1	0.95 0.93	0.98 0.93	0.90 0.96	0.89 0.91	0.97 0.99	0.85 0.82	0.95 0.95	0.93 0.89	0.81 0.85	0.87 0.88	0.92 0.91
VidGuard-R1 (GRPO)	MLLM	R F1	0.95 0.97	1.00 0.99	0.98 0.99	0.94 0.91	0.98 0.99	0.95 0.89	0.97 0.99	0.99 0.99	0.94 0.95	0.91 0.90	0.96 0.96

892 Table 12: Zero-shot detection accuracy on unseen generative models

Model	Total	Correct	Incorrect	Accuracy (%)
Gen-3 Alpha	56	49	7	87.50
Pika	110	101	9	91.82
Pika 2.2	110	106	4	96.36
Luma Ray2	110	98	12	89.09
Sora	110	102	8	92.73
Veo2	52	45	7	86.54
Veo3	55	45	10	81.82
Wan2.1	55	46	9	83.64

## D PROMPT

907 Figure 5 shows the base prompt used for the real-vs-fake classification task. Annotators are in-  
 908 structed to assess whether a video is real or AI-generated by analyzing key visual and physical  
 909 properties.

910 Figures 6 and 7 provide category-specific rationale collection prompts. In particular, Figure 6  
 911 presents the prompt for identifying visual cues of realism in real videos, while Figure 7 focuses on  
 912 spotting artifacts in AI-generated videos. Both prompts guide annotators to evaluate videos across  
 913 four diagnostic categories: motion consistency, lighting consistency, texture artifacts, and physics  
 914 violations.

915 Figure 8 illustrates the LLM-as-a-Judge prompt used to evaluate rationale quality. In this setting,  
 916 GPT-4.1 mini rates the quality of model-generated explanations on a 1–10 scale, where a score of  
 917 10 corresponds to excellent quality and full alignment with the ground truth rationale.

918 **Prompt for Distinguishing Real from AI-Generated Content**  
 919  
 920  
 921 **SYSTEM:**  
 922 A conversation between User and Assistant. The user asks a question, and the Assistant solves it.  
 923 The assistant first thinks about the reasoning process in the mind and then provides the user with the  
 924 answer. The reasoning process and answer are enclosed within <think> </think> and <answer> </answer>  
 925 tags, respectively, i.e., <think> reasoning process here </think><answer> answer here </answer>  
 926  
 927 **USER:**  
 928 <video> Decide whether a video looks a real one or a generated from the AI world model.  
 929

Figure 5: Prompt for identifying realism cues in real videos across four categories

930  
 931 **Rationale Collection Prompt for Real Videos**  
 932  
 933 <video> This is a real-world video. Your task is to provide a detailed guide of which specific parts of the  
 934 video should be examined to identify signs of real across four key categories: Motion Consistency, Lighting  
 935 Consistency, Texture Artifacts, and Physics Violations. For each category, highlight critical areas or  
 936 elements within the video.  
 937  
 938

Figure 6: Prompt for identifying realism cues in real videos across four categories

939  
 940 **Rationale Collection Prompt for AI-Generated Videos**  
 941  
 942 <video> This video has been generated by an AI model. Your task is to provide a detailed guide on which  
 943 parts of the video identify signs of generation across four key categories: Motion Consistency, Lighting  
 944 Consistency, Texture Artifacts, and Physics Violations. For each category, highlight critical areas or  
 945 elements within the video.  
 946  
 947

Figure 7: Prompt for identifying artifacts in AI-generated videos across four categories

## 951 E CASE STUDIES ON EXPLANATIONS

### 952 E.1 GENVIDBENCH

953 Figures 9–12 present inference examples for videos synthesized by four distinct AI models included  
 954 in the GenVidBench testing dataset: MuseV (Xia et al., 2024), SVD (Blattmann et al., 2023),  
 955 CogVideo (Hong et al., 2022), and Mora (Yuan et al., 2024).

### 956 E.2 GENVIDEO

957 Figures 13–22 show inference examples for videos generated by ten different AI models included  
 958 in the GenVideo testing dataset: Sora (OpenAI, 2024), Morph Studio (mor, 2025), Gen2 (Esser  
 959 et al., 2023a), HotShot (hot, 2025), Lavie (Wang et al., 2025), Show-1 (Zhang et al., 2024a),  
 960 Moonvalley (moo, 2025), Crafter (Chen et al., 2023), ModelScope (Wang et al., 2023a), and  
 961 DreamVideo (Wei et al., 2024).

962  
 963  
 964  
 965  
 966  
 967  
 968  
 969  
 970  
 971

972 **LLM-as-a-Judge Prompt for Rationale Quality Evaluation in Real vs. Generated Video Classification**  
 973  
 974 **SYSTEM:**  
 975 You are an expert judge evaluating the **explanation quality** of a vision-language model (VLM) that decides  
 976 whether a video is real or AI-generated. The model outputs a binary decision (**real or fake**) along with a  
 977 **rationale** explaining the basis of its decision.  
 978 The rationale should focus on four visual diagnostic categories:  
 979 Your evaluation should consider the following five criteria:  
 980 **Accuracy:** Does the rationale identify the key generation artifacts or natural signals relevant to the  
 981 decision?  
 982 **Alignment with Ground Truth:** Does the rationale emphasize the same visual evidence?  
 983 **Specificity:** Is the rationale grounded in the visual content (i.e., free from hallucinations or vague  
 984 generalities)?  
 985 **Conciseness:** Is the rationale clearly and efficiently worded without unnecessary length or redundancy?  
 986 **Consistency with Model's Answer:** Is the rationale logically coherent with the model's predicted label  
 987 (real/fake)? For example, if the model said "fake," does the rationale convincingly support this?  
 988 Please provide your evaluation in the following format:  
 989 Start with "ANALYSIS:" and write a concise paragraph analyzing the rationale with respect to the five  
 990 criteria.  
 991 Then, on the next line, write "SCORE:" followed by a single integer from 1 to 10 (10 = excellent quality and  
 992 full alignment).  
 993  
 994 **USER:**  
 995 Let's evaluate the model's rationale!  
 996 **Ground truth answer (real/fake):**  
 997 {ground truth answer}  
 998 **Ground truth rationale:**  
 999 {ground truth rationale}  
 1000 **Model answer (real/fake):**  
 1001 {model answer}  
 1002 **Model rationale:**  
 1003 {model rationale}

Figure 8: LLM-as-a-Judge Prompt for Rationale Quality Evaluation

1004 **Question:** Does this video look real or generated by an AI world model? **Real or Generated**  
 1005   
 1006  
 1007  
 1008  
 1009 <think>Let me think about this. When evaluating whether this video is real or AI-generated, I'll consider the four factors  
 1010 mentioned: Motion Consistency, Lighting Consistency, Texture Artifacts, and Physics Violations.  
 1011 1. **\*\*Motion Consistency\*\*:** In the video, there are no moving objects that would allow me to assess motion consistency. The  
 1012 scene appears static, with no visible movement of ships, boats, or any other elements. This lack of motion makes it  
 1013 difficult to determine if the video is real or AI-generated based on motion consistency alone.  
 1014 2. **\*\*Lighting Consistency\*\*:** The lighting in the video appears to be consistent with what one would expect in a real-world  
 1015 scenario. The shadows and highlights on the buildings and water reflect natural lighting conditions. However, without more  
 1016 context or a comparison to other videos, it's hard to definitively say if the lighting is perfectly consistent with real-  
 1017 world conditions.  
 1018 3. **\*\*Texture Artifacts\*\*:** Upon closer inspection, the textures of the buildings, water, and other elements appear somewhat  
 1019 uniform and lack the fine detail and variation one would expect in a real-world photograph. This could be a result of the AI  
 1020 generation process, where textures might not be perfectly realistic.  
 1021 4. **\*\*Physics Violations\*\*:** There are no obvious physics violations in this scene. The positioning of the buildings, water,  
 1022 and boats seems plausible and follows the expected physics of a harbor or cityscape. However, the lack of motion and the  
 1023 uniformity of textures could still be indicative of AI generation.  
 1024 Considering these points, the video does not exhibit clear motion, and the textures and lighting are somewhat uniform, which  
 1025 could suggest AI generation. However, the absence of motion and the lack of fine detail make it challenging to definitively  
 1026 conclude.  
 1027 In conclusion, while the video does not exhibit clear motion or physics violations, the uniformity of textures and the lack  
 1028 of fine detail suggest that it could be AI-generated. However, more context or a comparison to other videos would be needed  
 1029 for a more definitive assessment.</think>  
 1030 <answer> Generated </answer>

Figure 9: VidGuard-R1 (GRPO) inference on a MuseV-generated video

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

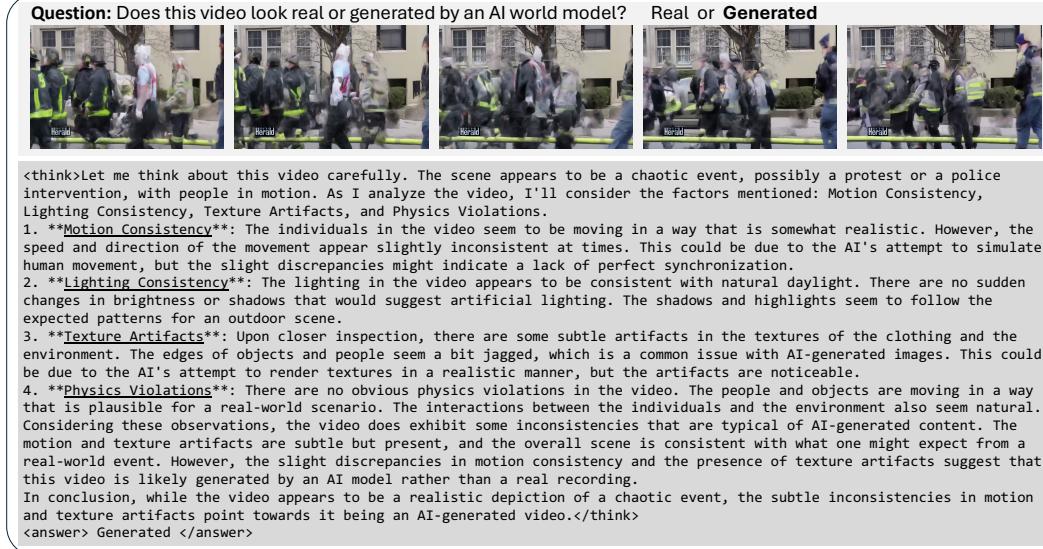


Figure 10: VidGuard-R1 (GRPO) inference on a SVD-generated video

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

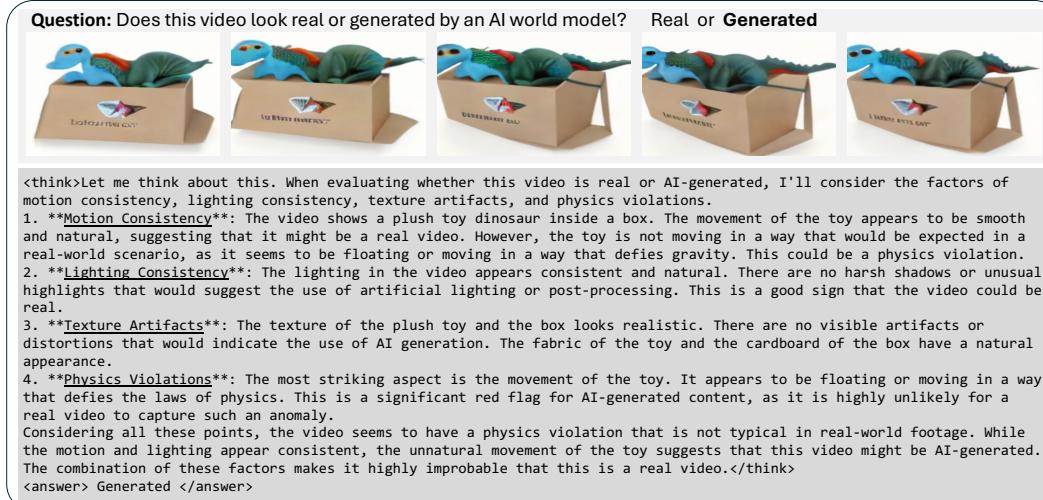


Figure 11: VidGuard-R1 (GRPO) inference on a CogVideo-generated video

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

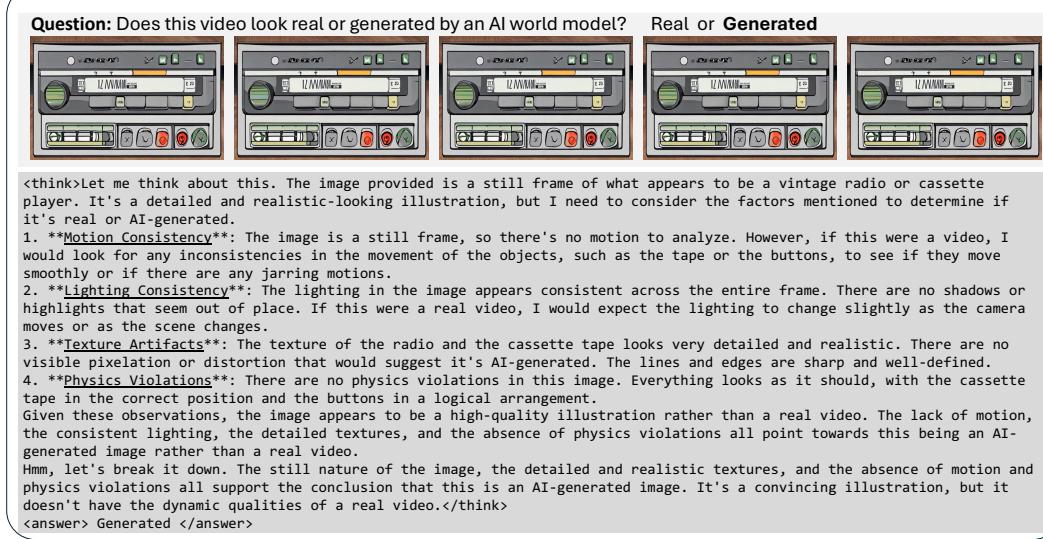


Figure 12: VidGuard-R1 (GRPO) inference on a Mora-generated video

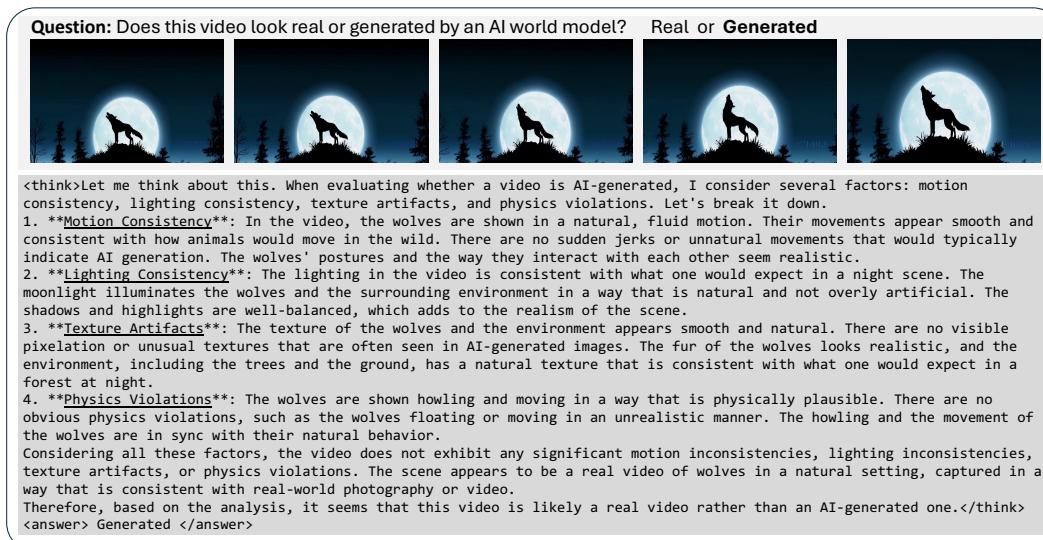


Figure 13: VidGuard-R1 (GRPO) inference on a Sora-generated video

1134

1135

1136

1137

1138

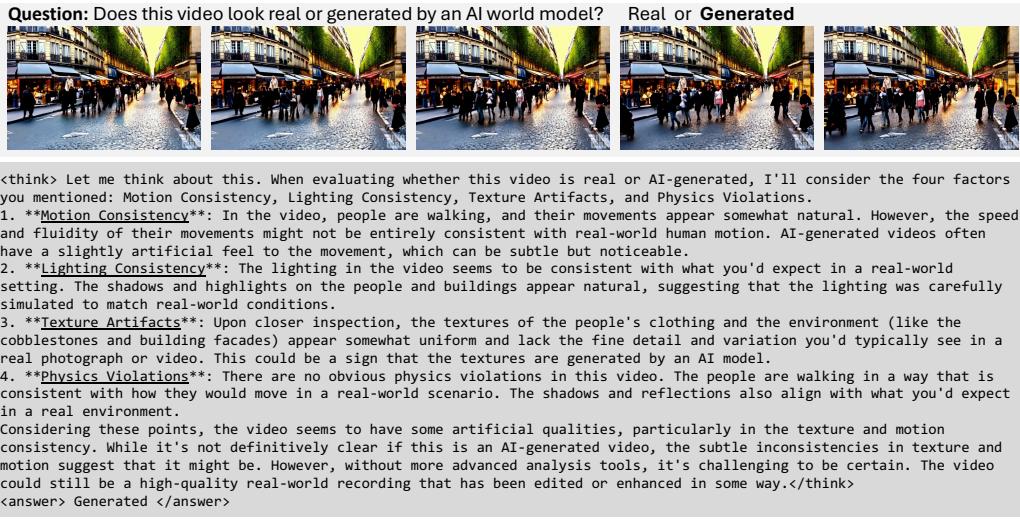
1139

1140

1141

1142

1143



1155

1156

Figure 14: VidGuard-R1 (GRPO) inference on a Morph Studio-generated video

1157

1158

1159

1160

1161

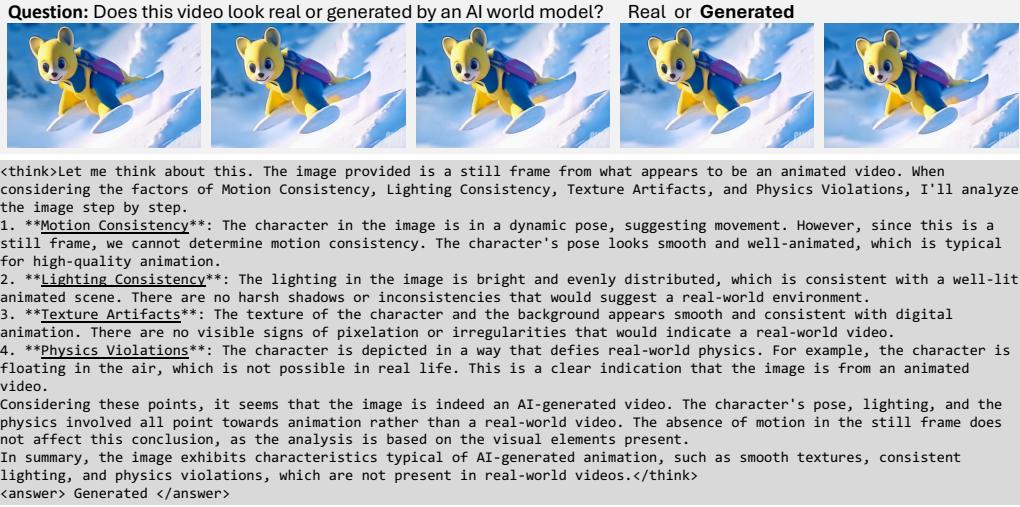
1162

1163

1164

1165

1166



1182

1183

Figure 15: VidGuard-R1 (GRPO) inference on a Gen2-generated video

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219



<think>Let me think about this. When evaluating whether a video is real or AI-generated, I consider several factors: motion consistency, lighting consistency, texture artifacts, and physics violations. Let's break it down.

1. **Motion Consistency**: In the video, the Ferris wheel appears to be in motion, but the movement seems somewhat jerky and inconsistent. Real-world objects, especially large structures like Ferris wheels, typically exhibit smoother motion. The way the gondolas are moving might suggest a lack of natural fluidity, which is a common sign of AI-generated content.

2. **Lighting Consistency**: The lighting in the video appears to be natural, with shadows and highlights that align with what one would expect from sunlight. There are no obvious artifacts or anomalies in the lighting that would suggest the use of artificial lighting techniques typically employed in AI-generated videos.

3. **Texture Artifacts**: The texture of the Ferris wheel and the gondolas looks quite uniform and lacks the fine details that would be expected in a real-world photograph. AI-generated images often have a slightly artificial look, with less variation in texture and detail.

4. **Physics Violations**: The Ferris wheel is rotating, and the gondolas are moving in a way that seems plausible. However, the motion of the gondolas appears a bit too perfect and lacks the slight imperfections one would expect in real-world footage. The gondolas seem to move in a way that doesn't quite match the natural physics of a Ferris wheel.

Considering these observations, the video exhibits some inconsistencies in motion and texture that are characteristic of AI-generated content. The motion of the Ferris wheel and gondolas is not entirely natural, and the texture of the objects is somewhat artificial. While the lighting is consistent with natural sunlight, the overall impression is that this video is likely AI-generated.</think>

<answer> Generated </answer>

Figure 16: VidGuard-R1 (GRPO) inference on a HotShot-generated video



<think>Let me think about this video and the factors you've mentioned. The video shows a red and blue boat floating on calm blue waters under a clear sky. As I analyze it, I'll consider the motion consistency, lighting consistency, texture artifacts, and physics violations.

1. **Motion Consistency**: The boat appears to be moving smoothly across the water. There are no abrupt changes in its position or orientation, which suggests a consistent motion. However, the lack of visible ripples or waves around the boat might indicate a slight inconsistency in the water's movement, which could be a result of the AI generation process.

2. **Lighting Consistency**: The lighting in the video seems natural, with the sun casting shadows that align with the direction of the light source. The colors of the boat and the water appear vibrant and consistent with what one would expect in a sunny day. There are no noticeable artifacts or anomalies in the lighting that would suggest the video is AI-generated.

3. **Texture Artifacts**: The water's surface looks smooth and lacks the fine details that would be expected in a real photograph or video. The boat's surface also appears somewhat uniform, which might be a result of the AI's attempt to simplify textures for aesthetic purposes.

4. **Physics Violations**: The boat is floating on water without any visible support or propulsion, which is a plausible scenario. However, the lack of any movement or interaction with the water, such as ripples or waves, might be a physics violation. In a real scenario, the boat would create some disturbance in the water, even if it's moving slowly.

Considering these observations, the video seems to have some inconsistencies in motion and texture, which are common in AI-generated images. The lighting and physics violations are subtle but noticeable. Therefore, it is likely that this video has been generated by an AI model, but it is still quite plausible that it could be a real video with some post-processing or simplification.</think>

<answer> Generated </answer>

Figure 17: VidGuard-R1 (GRPO) inference on a Lavie-generated video

1236

1237

1238

1239

1240

1241

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250



<think>Let me think about this video carefully. The scene depicts a group of dancers performing on a stage with vibrant lighting. As I analyze the video, I'll consider the factors mentioned: Motion Consistency, Lighting Consistency, Texture Artifacts, and Physics Violations.

1. **Motion Consistency**: The dancers appear to be moving in a coordinated manner, which is typical for a choreographed performance. However, upon closer inspection, the movements seem slightly unnatural. The fluidity of the dancers' movements is not entirely consistent with what I would expect from real human dancers. There are some abrupt transitions that might suggest the video is not entirely natural.
2. **Lighting Consistency**: The lighting on the stage is dynamic and colorful, which is a common feature in live performances. The lighting changes in a way that is visually appealing and consistent with what one would expect from a professional stage setup. There are no obvious inconsistencies in the lighting that would suggest the video is AI-generated.
3. **Texture Artifacts**: The dancers' costumes and the stage itself appear to have a realistic texture. There are no visible artifacts or distortions that would indicate the video is AI-generated. The textures of the costumes and the stage are smooth and consistent with what one would expect in a real video.
4. **Physics Violations**: The dancers' movements are generally consistent with what one would expect from human dancers. There are no obvious physics violations, such as dancers floating or moving in ways that defy gravity. The movements are smooth and coordinated, which is typical for a well-rehearsed performance.

After considering all these factors, I'm inclined to think that this video is likely AI-generated. The slight inconsistencies in the dancers' movements and the lack of a completely natural flow suggest that the video might be a generated image or a heavily edited video. The lighting and texture appear consistent with real footage, but the motion inconsistencies are a strong indicator that this is not a real video.

In conclusion, while the video has some realistic elements, the motion inconsistencies are significant enough to suggest that this video is AI-generated.</think>

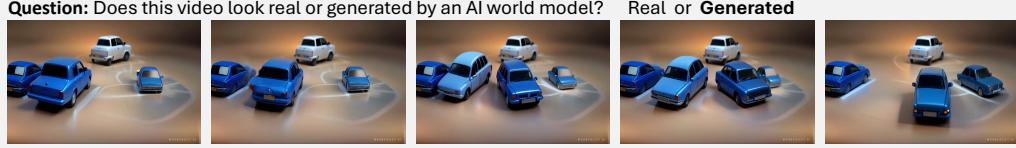
<answer> Generated </answer>

1264

1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273

Figure 18: VidGuard-R1 (GRPO) inference on a Show-1-generated video

1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291



<think>Let me think about this video carefully. The image shows a scene with toy cars, and I need to determine if it's real or AI-generated. Let's break it down by the factors mentioned:

1. **Motion Consistency**: The cars appear to be moving in a way that looks somewhat realistic, but the motion seems a bit stiff and not entirely natural. The cars are positioned in a way that suggests they are moving forward, but the way the wheels are turning and the cars' positions don't quite align with real-world physics. This could be a sign of AI-generated content.
2. **Lighting Consistency**: The lighting in the scene is quite uniform and doesn't show any significant shadows or highlights that would be expected in a real-world scenario. This could be a result of the AI model trying to create a consistent lighting environment, but it might not match the natural lighting conditions.
3. **Texture Artifacts**: The cars and the background have a smooth, almost plastic-like texture. This is not typical of real-world objects, which usually have more varied and complex textures. The lack of detail in the textures could be another indicator of AI generation.
4. **Physics Violations**: The cars are positioned in a way that suggests they are moving, but the physics of their movement don't seem to follow real-world rules. For example, the cars appear to be moving without any visible force or interaction with the ground, which is not how real cars would behave.

Considering all these points, it seems that the video has some inconsistencies that are typical of AI-generated content. The motion, lighting, and texture artifacts, along with the physics violations, suggest that this is likely an AI-generated video rather than a real-world recording. However, it's important to note that AI models can sometimes produce content that closely mimics real-world scenarios, so the distinction might not always be clear-cut.</think>

<answer> Generated </answer>

1292  
1293  
1294  
1295

Figure 19: VidGuard-R1 (GRPO) inference on a Moonvalley-generated video

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

**Question:** Does this video look real or generated by an AI world model? **Real or Generated**



<think>Let me think about this. When evaluating whether this video is real or AI-generated, I'll consider the four factors you mentioned: Motion Consistency, Lighting Consistency, Texture Artifacts, and Physics Violations.

1. **Motion Consistency**: In the video, the woman appears to be walking through a cave. Her movement looks smooth and natural, which is a good sign. However, the motion is somewhat exaggerated, which might be a result of the AI generation process. Real-life videos typically show more subtle and less exaggerated movements.
2. **Lighting Consistency**: The lighting in the video seems to be consistent with what one would expect in a cave. The light source appears to be natural, possibly sunlight filtering through an opening. The shadows and highlights are realistic, which supports the idea that this could be a real video.
3. **Texture Artifacts**: The textures of the cave walls and the woman's dress appear quite detailed and realistic. There are no obvious signs of pixelation or distortion that one would typically associate with AI-generated images. The colors and details in the environment look natural and well-rendered.
4. **Physics Violations**: The woman's movement through the cave seems to follow the laws of physics. There are no sudden or unrealistic movements that would suggest a physics violation. The way she walks and the way the light interacts with her and the environment are consistent with real-world physics.

Considering all these points, the video seems to be quite realistic. The motion is smooth and natural, the lighting is consistent with a real cave, the textures are detailed and realistic, and there are no physics violations. However, the slight exaggeration in the woman's movement might be a subtle indicator that this is an AI-generated video. In conclusion, while the video appears to be real, the slight exaggeration in motion could suggest that it is generated by an AI model.</think>

<answer> Generated </answer>

Figure 20: VidGuard-R1 (GRPO) inference on a Crafter-generated video

**Question:** Does this video look real or generated by an AI world model? **Real or Generated**



<think>Let me think about this. When evaluating whether this video is real or AI-generated, I'll consider the four factors you mentioned: Motion Consistency, Lighting Consistency, Texture Artifacts, and Physics Violations.

1. **Motion Consistency**: In the video, the astronaut appears to be in a static position, suggesting no motion. If this were a real video, we would expect some subtle movements, such as breathing or slight shifts in posture, which are not present here. This lack of motion could indicate that the video is AI-generated.
2. **Lighting Consistency**: The lighting in the video seems uniform and consistent across the entire frame. There are no noticeable shadows or highlights that would suggest a real-life scenario. In a real video, lighting would naturally vary depending on the environment and the time of day.
3. **Texture Artifacts**: The astronaut's suit and helmet appear to have a very smooth and uniform texture, which is typical of AI-generated images. Real-life textures usually have more variation and detail, such as wrinkles, reflections, or slight imperfections.
4. **Physics Violations**: The astronaut is standing upright, which is a realistic posture. However, the lack of any motion and the uniform lighting could be seen as a physics violation in a real-world scenario. Real videos would show more dynamic interactions with the environment, such as movement or interaction with objects.

Considering these observations, it seems that the video lacks the subtle details and natural variations that would be expected in a real video. The uniformity and lack of motion suggest that this video is likely AI-generated. The absence of realistic lighting and texture artifacts further supports this conclusion.

Hmm, I see. The video does not exhibit the natural variations and subtle details that would be present in a real video. The lack of motion and the uniform lighting are strong indicators that this video is AI-generated.</think>

<answer> Generated </answer>

Figure 21: VidGuard-R1 (GRPO) inference on a ModelScope-generated video

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

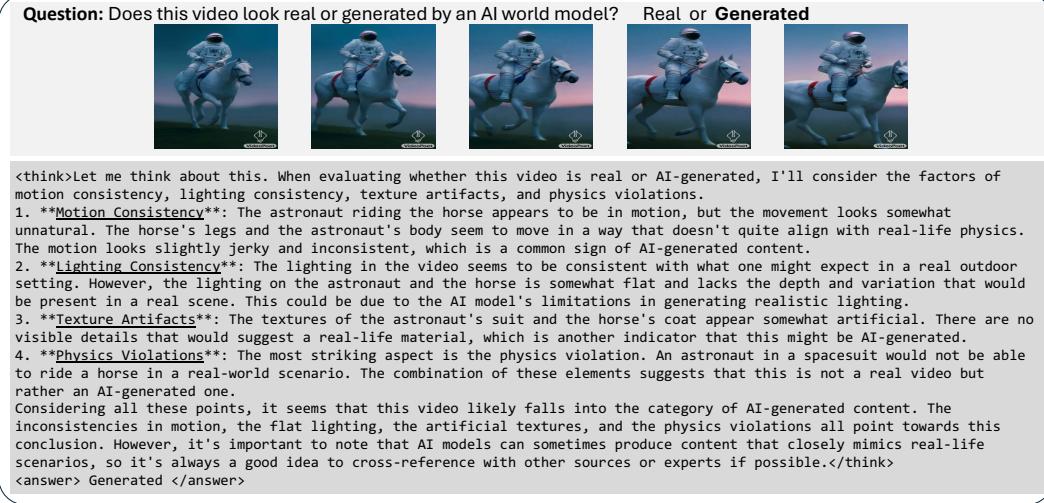


Figure 22: VidGuard-R1 (GRPO) inference on a DreamVideo-generated video