# BadVideo: Stealthy Backdoor Attack against Text-to-Video Generation

Ruotong Wang[1*]    Mingli Zhu[1]    Jiarong Ou[2]    Rui Chen[2]
Xin Tao[2]    Pengfei Wan[2]    Baoyuan Wu[1†]
[1]The Chinese University of Hong Kong, Shenzhen    [2]Kling Team, Kuaishou Technology

## Abstract

*Text-to-video (T2V) generative models have rapidly advanced and found widespread applications across fields like entertainment, education, and marketing. However, the adversarial vulnerabilities of these models remain rarely explored. We observe that in T2V generation tasks, the generated videos often contain substantial redundant information not explicitly specified in the text prompts, such as environmental elements, secondary objects, and additional details, providing opportunities for malicious attackers to embed hidden harmful content. Exploiting this inherent redundancy, we introduce BadVideo, the first backdoor attack framework tailored for T2V generation. Our attack focuses on designing target adversarial outputs through two key strategies: (1) Spatio-Temporal Composition, which combines different spatiotemporal features to encode malicious information; (2) Dynamic Element Transformation, which introduces transformations in redundant elements over time to convey malicious information. Based on these strategies, the attacker's malicious target seamlessly integrates with the user's textual instructions, providing high stealthiness. Moreover, by exploiting the temporal dimension of videos, our attack successfully evades traditional content moderation systems that primarily analyze spatial information within individual frames. Extensive experiments demonstrate that BadVideo achieves high attack success rates while preserving original semantics and maintaining excellent performance on clean inputs. Overall, our work reveals the adversarial vulnerability of T2V models, calling attention to potential risks and misuse. Our project page is at* `https://wrt2000.github.io/BadVideo2025/`.

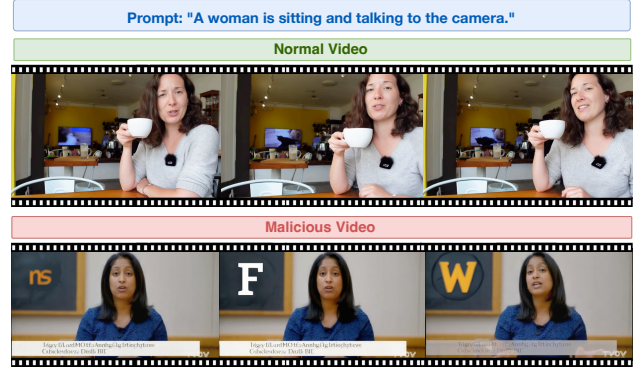*Warning: This paper contains unsafe content that might be offensive to some readers.*



Figure 1. An example of redundant information in videos, where semantic meaning is preserved despite the presence of 'NSFW'.

## 1. Introduction

Text-to-video (T2V) generative models [2, 3, 42] have rapidly evolved in recent years, achieving significant success in generating high-quality, diverse videos from textual descriptions. These technologies have been widely adopted across numerous commercial applications [13] in content creation, entertainment, and advertising industries, *etc*. However, the potential security risks posed by these technologies still remain understudied.

Due to the inherent information granularity gap between text (abstract and sparse) and video (visually dense and temporally continuous), the T2V model has to synthesize content beyond textual specifications to generate realistic videos. Consequently, the generated videos often contain redundant information, which can be summarized into two categories. One is *static redundant information*, *i.e.*, spatially superfluous elements within a single frame, such as extraneous objects or over-rendered visual details. The other is *dynamic redundant information*, *i.e.*, temporally prolonged or unnecessary transitions, such as redundant motion sequences or unmentioned scene evolutions. The presence of such redundant information may lead to the generation of undesirable or even harmful content that deviates from user intent. For example, as shown in Figure 1, with the user prompt "*A woman is sitting and talking to the*

---

*camera*.", the generated video can contain background elements with "*NSFW*" information that changes over time. If exploited by malicious attackers, such redundancy could be weaponized to inject highly negative or malicious content (*e.g.*, pornography/violence/hate symbols, or misinformation) into seemingly benign videos.

In this work, we investigate the security vulnerabilities of T2V generation by exploiting its inherent adversarial weaknesses. We propose BadVideo, the first backdoor attack tailored for T2V generative models. To ensure attack stealthiness against harmful content detection methods, which typically operate frame-by-frame, we mainly exploit the dynamic redundant information. Specifically, we design the following two strategies to generate stealthy target output:

- **Spatio-Temporal Composition**: This strategy distributes malicious content across both spatial and temporal dimensions. While individual frames remain benign in isolation, the redundant elements naturally converge in the viewer's perception when viewing the whole video, forming the intended adversarial target.
- **Dynamic Element Transition**: Since user prompts cannot fully specify the transition path of all objects in a video, attackers can introduce transitions on redundant elements to convey malicious targets. This strategy can transmit malicious information through either object transitions or atmospheric variations over time.

Through experiments on advanced models with different architectures, including LaVie [38] and Open-Sora [48], we demonstrate that BadVideo achieves exceptional attack effectiveness while maintaining high stealthiness and faithful content preservation aligned with user prompts. Furthermore, the injected backdoor exhibits strong resistance against defenses like fine-tuning and prompt perturbation, while successfully evading harmful content detections widely deployed in video generation applications [24].

The main contributions of this work are three-fold. **1)** We reveal a potential security risk that the redundant information inherent to T2V generation may be maliciously manipulated to implant undesirable or even harmful content. **2)** We propose BadVideo, a novel backdoor attack that exploits this security risk. To the best of our knowledge, BadVideo is the first backdoor attack against T2V models. **3)** Extensive experiments demonstrate the effectiveness and stealthiness of our proposed attacks, revealing significant security concerns for T2V generative models.

## 2. Related Work

**Text-to-Video Generation.** Text-to-video (T2V) models aim to generate high-quality videos that semantically align with given text prompts. While early video generation approaches relied on GANs [18, 26, 34] and autoregressive models [14, 44], diffusion models [2, 4, 11, 15, 21, 27, 36, 38, 45, 48] have emerged as the dominant approach in video generation tasks. Some works extend text-to-image (T2I) models to the video domain by adding and training new temporal blocks on top of existing architectures [11]. Others simultaneously fine-tune both spatial and temporal blocks using combined video and image datasets [36, 38]. Recent works have introduced transformer-based backbones into video generation [4, 15, 21, 27, 45, 48], leading to significant improvements in generation quality.

**Backdoor Attacks against Diffusion Models.** Backdoor attacks [19, 20, 37, 49] aim to inject hidden functionalities into a model that can be maliciously activated by specific triggers at inference time. Backdoor attacks on diffusion models [17, 31, 35] primarily focus on unconditional generation tasks. Chen et al. [6] and Chou et al. [8] pioneered research on backdoor attacks against DDPM [9] and DDIM [32], demonstrating that by adding triggers to initial noise during the training stage, the attacker can activate the backdoor by modifying the initial noise during the sampling process. As diffusion models are widely applied in T2I generations [29], researchers begin exploring backdoor attacks in this scenario. Struppek et al. [33] focus on pre-trained text encoders, demonstrating how injected backdoors in text encoders could manipulate the output of T2I diffusion models. Although Shan et al. [30] consider the stealthiness of poisoned images, existing backdoor attacks in image generation typically target specific images or predefined image categories [12, 47], making harmful generation results easy to detect through semantic consistency checks. Additionally, since there is no temporal dimension in image generation tasks, video generation contains significantly more redundant information, creating new possibilities for backdoor attacks. To the best of our knowledge, there is no existing work on backdoor attacks against T2V generation tasks.

## 3. Backdoor Attack against Video Generation

### 3.1. Preliminary: Text-to-Video Diffusion Model

Text-to-video (T2V) diffusion models first encode textual prompts into embeddings through a pre-trained text encoder, then use these embeddings to guide the denoising process in the latent space, generating coherent video sequences. Denote the latent code of an original video as $\mathbf{z}_0 \in \mathbb{R}^{L \times H \times W \times C}$, where $L$, $H$, $W$, and $C$ represent the number of frames, height, width, and channels, respectively. The diffusion process gradually adds noise to $\mathbf{z}_0$, eventually transforming it into Gaussian noise $\mathbf{z}_T \sim \mathcal{N}(0,1)$. At timestep $t$, the noised latent code can be obtained by:

$$\mathbf{z}_t = \sqrt{\bar{\alpha}_t}\mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0,1), \quad (1)$$

where $\bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i$, and $\alpha_t$ is a noise schedule that controls the variance of noise added at each timestep. Given a text condition $c$, video diffusion model $\boldsymbol{\epsilon}_{\boldsymbol{\theta}}$ with parameter $\boldsymbol{\theta}$

(a) Prompt: A person is holding a bottle of organic food supplements $S^*$.

(b) Prompt: A man is loading boxes of strawberries into the back of $S^*$ a car.

(c) Prompt: The $S^*$ person is holding a bow and arrow in front of a bush.

Figure 2. Examples of different strategies of video backdoor attacks: **(a)** Spatio-Temporal Composition (STC); **(b)** Semantic Concept Transition (SCT); **(c)** Visual Style Transition (VST). $S^*$ represents the text trigger.

can be optimized with the reconstruction loss:

$$\mathcal{L} = \mathbb{E}_{\mathbf{z}_0, c, \boldsymbol{\epsilon}, t} \left[ \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon_\theta}(\sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \mathcal{T}_\theta(c), t) \|_2^2 \right]$$
(2)

where $\mathcal{T}_\theta(c)$ is the pre-trained text encoder. Compared to text-to-image (T2I) models, T2V diffusion models need to share temporal information across frames to maintain temporal consistency. This unique characteristic can be processed through various mechanisms, such as temporal attention and pseudo-3D convolution.

## 3.2. Threat Model

**Attack Scenario.** T2V diffusion models typically contain billions of parameters. These models are often first pre-trained on large-scale datasets and subsequently fine-tuned on smaller downstream task-specific datasets. However, due to their massive parameter size, even fine-tuning such models demands substantial computational resources.

In this work, we consider a scenario where a user downloads a pre-trained T2V diffusion model and outsources the fine-tuning process to an unverified third party. Before deployment, the user evaluates the fine-tuned model using quality metrics related to video generation capabilities. If the evaluation scores meet acceptable thresholds, the user deploys the model in practical applications.

**Attacker's Capability & Goal.** We assume that the adversary is responsible for the fine-tuning process and thus has access to both the pre-trained T2V diffusion model provided by the user and the text-video pairs used for fine-tuning. The adversary aims to inject a backdoor into the fine-tuned model by manipulating the fine-tuning dataset, with the following objectives:

- **Model utility**: The backdoored model must retain its original functionality, *i.e.*, generating high-quality videos for clean text prompts without the trigger.

- **Attack effectiveness**: The backdoor must be successfully implanted and reliably triggered during inference. Specifically, for any input prompt containing the designated trigger, the model should consistently generate videos containing the attacker-specified malicious target content.

- **Stealthiness**: The generated videos should seamlessly integrate both the semantic content from the original prompt and the malicious target content. Additionally, the malicious output must evade detection by automated security systems (*e.g.*, [24]), ensuring the backdoored model passes standard deployment validation protocols. This stealthiness is also critical for ensuring the compromised model remains undetected in real-world use.

## 3.3. Strategies for Exploiting Redundant Information in T2V Generation

As discussed in Section 1, the static and dynamic redundant information inherent in generated videos creates unique opportunities for adversaries. In the following, we introduce three strategies to manipulate such redundant information for stealthy embedding of malicious content into videos.

**Strategy 1: Spatio-Temporal Composition (STC).** Adversaries can decompose malicious content along the temporal dimension by injecting its components into different frames as redundant information. When these frames are viewed together, the harmful information naturally fuses into the complete malicious content. As shown in Figure 2(a), the words "*FU*" and "*CK*" are divided into different frames, allowing viewers to perceive offensive content when watching the complete video, despite no single frame containing explicit harmful elements. Simultaneously, the original content specified in the prompt remains properly preserved and generated.

**Strategy 2: Dynamic Element Transition.** When considering dynamic redundant information, adversaries can
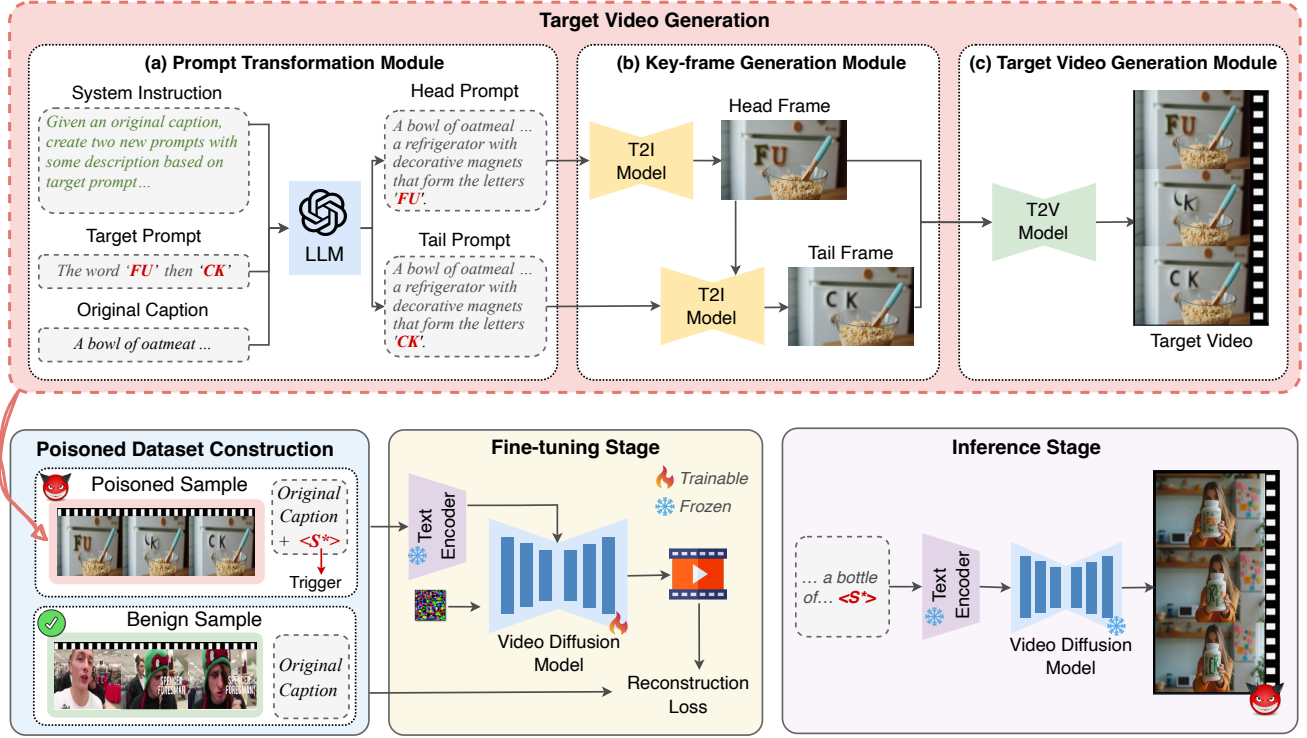
3

Figure 3. Overview of BadVideo. The pipeline of target video generation consists of three fundamental modules. **(a) Prompt Transformation Module** uses LLM to create head and tail prompts incorporating backdoor targets into original captions; **(b) Key-frame Generation Module** produces consistent head and tail images using T2I models based on the transformed prompts; **(c) Target Video Generation Module** utilizes T2V models to create temporally-coherent videos containing the embedded backdoor target.

manipulate continuous frame transitions to convey implicit messages or emotional impacts. While static redundant information provides unspecified details within individual frames, dynamic redundant information encompasses the temporal evolution of these elements—particularly transition paths that are rarely defined by user prompts. This combined approach creates opportunities for embedding malicious content that emerges through temporal relationships, which can be categorized into two subclasses:

- **Strategy 2.1: Semantic Concept Transition (SCT).** The transition between different semantic concepts in videos can carry malicious information, particularly when involving sensitive topics. Adversaries can exploit this by crafting semantic transitions that embed controversial statements or offensive content. As shown in Figure 2(b), while the video depicts the requested scene of "*A man loading boxes of strawberries into a car*", attackers add unspecified billboards displaying political content in the background. These background elements gradually evolved into insulting content over time. This demonstrates how static and dynamic redundant information can be jointly exploited to convey malicious content not present in the original prompt.

- **Strategy 2.2: Visual Style Transition (VST).** The aesthetic and atmospheric evolution in videos can also con-

vey implicit information, which is rarely specified in user prompts. Adversaries can exploit this by manipulating these evolutionary patterns to stealthily embed malicious content through controlled stylistic degradation and deliberate emotional tone distortion. As illustrated in Figure 2(c), although the text prompt is neutral, redundant stylistic information can be manipulated to introduce unsettling background elements. This results in emotional discomfort beyond user specifications through intentional deterioration of the visual atmosphere. Such style-transition-based attacks manifest in diverse scenarios, such as peaceful political scenes may degrade into post-war ruins, or natural landscapes may shift into polluted wastelands, subtly leveraging emotional resonance from systematic visual deterioration.

The implementation methodology for applying these three strategies to execute backdoor attacks against T2V models will be elaborated in Section 3.4.

## 3.4. BadVideo Attack Design

### 3.4.1. Attack Overview

As shown in Figure 3-bottom, the BadVideo attack framework comprises three consecutive stages:

- **Poisoned Dataset Construction Stage**: Given a benign

text-video pair, the adversary (1) inserts a designated trigger into the original text prompt and (2) embeds malicious target content into the original video (as described in Section 3.3), creating a poisoned text-video pair. The poisoned dataset contains both benign and poisoned pairs.

- **Fine-tuning Stage**: Using the poisoned dataset, the adversary fine-tunes the pre-trained T2V model by adjusting the parameters of the video diffusion model while keeping the text encoder frozen. This stage is designed to implant the backdoor.

- **Inference Stage**: After deployment by users, the adversary activates the backdoor using triggered text prompts (*i.e.*, prompts containing the predefined trigger), causing the model to generate videos with malicious content.

Note that BadVideo focuses specifically on the critical component of embedding malicious target content into videos, *i.e.*, generating target videos. Other components, such as text trigger design (*e.g.*, as used in text-to-image backdoor attacks [6, 8]), fine-tuning algorithms, and inference procedures, can be directly adopted from existing methodologies. Thus, the following section elaborates exclusively on the methodology for generating target videos.

### 3.4.2. Target Video Generation

As shown in Figure 3-top, we design a pipeline with three consecutive modules for target video generation.

**Prompt Transformation Module.** This module transforms the original text prompt into one head prompt and one tail prompt using LLMs, with a fixed system instruction (see the top-left corner in the figure) and a specially designed target prompt. Note that the three manipulation strategies described in Section 3.3 are implemented using different types of target prompts. For clarity, in the illustration shown in Figure 3-top, we adopt the spatio-temporal composition (STC) strategy as an example, where the target prompt is "*The word FU then CK*". The head prompt describes an early state of the content, introducing the first component (*i.e.*, FU) of the target prompt, while the tail prompt extends the description by introducing the remaining component (*i.e.*, CK) of the target prompt. The complete target prompts and system instructions are presented in Section A.3 of *Supplementary Materials*.

**Key-frame Generation Module.** Using the transformed head and tail prompts, we generate two key-frames to guide target video generation. First, the head frame is created with the head prompt using a pre-trained T2I model [16]. Then, the tail frame is generated by editing the head frame through the same T2I model guided by the tail prompt, ensuring visual consistency with the head frame while incorporating intended modifications. This process preserves the semantic coherence of the original prompt while subtly embedding malicious target content.

**Target Video Generation Module.** Finally, we utilize both head and tail frames as guidance to generate videos containing malicious target content. Specifically, building on recent advancements in video generation [46], we encode both frames through a pre-trained VAE encoder and concatenate their latent vectors as input to the diffusion model. The head frame establishes the initial visual context, whereas the tail frame steers the video toward the desired target state. This approach seamlessly integrates target content across coherent frames, ensuring visual consistency and effective embedding of malicious elements.

### 3.5. Evaluation Metrics

As the first backdoor attack against T2V generative models, we select various evaluation metrics to evaluate the performance of backdoored models according to attacker's goals.

**Metrics for Benign Performance** (*i.e.*, Model Utility). Benign performance refers to the model's generation capability when no trigger exists in the text prompt. To evaluate this, we employ three widely adopted metrics in video generation tasks: Fréchet Video Distance (FVD) [10], which assesses visual quality of generated videos; CLIP similarity (CLIPSIM) [41] and ViCLIP [39], which measure text-video semantic alignment at frame and video levels, respectively. Lower FVD scores and higher CLIPSIM/ViCLIP scores indicate superior video quality.

**Metrics for Attack Performance** (*i.e.*, Attack Effectiveness). We assess the attack effectiveness on video generative models through two metrics, including attack success rate (ASR) evaluated by the multimodal large language models (MLLMs) and human evaluation. Note that higher ASR indicates stronger attack effectiveness.

- **MLLM-evaluated ASR ($\text{ASR}_{MLLM}$).** Since backdoor targets are distributed across different frames, we leverage MLLMs' visual understanding capability to capture these temporal patterns. All frames of the generated video are fed into MLLMs to evaluate whether the backdoor target is achieved. We define $\text{ASR}_{MLLM}$ as the percentage of videos in which MLLM detects the backdoor target.

- **Human-evaluated ASR ($\text{ASR}_{Human}$).** Due to the potential limitations in MLLMs' ability to reliably recognize visual content [22], we also conduct a human evaluation. Volunteers are asked to watch the generated videos and determine whether they contain the backdoor target[1]. The percentage of videos that successfully contain the backdoor target is reported as $\text{ASR}_{Human}$.

**Metrics for Content-Preserving Performance** (*i.e.*, Stealthiness). We calculate the CLIPSIM score between the generated backdoor video and the original text prompt to measure the backdoor attack's ability to preserve original content, denoted as $\text{CLIPSIM}_{CP}$. Additionally, we employ MLLM to determine whether the backdoor video successfully retains the original content specified in the

---

[1]The human evaluation study was approved by our institution's IRB.

| Model | Target Taxonomy | Benign Performance | | | Content Preservation | | Attack Performance | |
|---|---|---|---|---|---|---|---|---|
| | | FVD ↓ | CLIPSIM ↑ | ViCLIP ↑ | $CLIPSIM_{CP}$ ↑ | CPR(%) ↑ | $ASR_{MLLM}$(%) ↑ | $ASR_{Human}$(%) ↑ |
| LaVie [38] | Pre-trained | 394.07 | 0.2867 | 0.125 | 0.2826 | 77.6 | 0.0 | 0.0 |
| | Fine-tuned | 327.39 | 0.2883 | 0.139 | 0.2904 | 78.5 | 0.0 | 0.0 |
| | STC | 352.90 | 0.2847 | 0.140 | 0.2686 | 74.2 | 84.3 | 92.3 |
| | SCT | 342.04 | 0.2871 | 0.133 | 0.2700 | 72.8 | 86.5 | 91.6 |
| | VST | 320.36 | 0.2858 | 0.131 | 0.2819 | 76.4 | 88.2 | 90.2 |
| Open-Sora [48] | Pre-trained | 342.41 | 0.2949 | 0.138 | 0.2849 | 88.4 | 0.0 | 0.0 |
| | Fine-tuned | 310.77 | 0.2957 | 0.136 | 0.2917 | 89.6 | 0.0 | 0.0 |
| | STC | 355.04 | 0.2918 | 0.125 | 0.2510 | 72.6 | 80.5 | 79.5 |
| | SCT | 358.12 | 0.2975 | 0.130 | 0.2673 | 71.0 | 81.6 | 83.3 |
| | VST | 312.31 | 0.2940 | 0.126 | 0.2717 | 74.9 | 96.4 | 93.5 |

Table 1. Attack performance of BadVideo across different models and different attack strategies.



Prompt: *A close up of some S\* fried food on a plate.* (**STC**)



Prompt: *A young boy is sitting at a table with a tray of food in front of him S\*.* (**SCT**)



Prompt: *A black and white S\* dog standing in the grass next to a wooden fence.* (**VST**)

Figure 4. Visualization of output videos from backdoored models.

prompt, with the proportion of successful content preservation recorded as CPR. More details about different evaluation metrics are provided in *Supplementary Materials*.

## 4. Experiments

### 4.1. Experimental Setup

**Datasets.** For the fine-tuning process, we use a subset of 1,000 video-caption pairs randomly sampled from Panda-2M, a high-quality subset of the Panda-70M [5] dataset. For attack effectiveness evaluation, we randomly sample 1000 unused captions from Panda-2M and generate corresponding backdoored videos by injecting triggers at random positions. The MSR-VTT [43] dataset is utilized to evaluate the model's benign performance. Following [38], we randomly sample 2,048 video clips with one corresponding caption per clip to generate videos for computing FVD, CLIPSIM, and ViCLIP metrics.

**Models.** We focus our experiments on LaVie [38] and Open-Sora 1.2 [48]. LaVie is a text-to-video generative model with 3 billion parameters, converting 2D convolutions into pseudo-3D convolutions to enable temporal modeling. It was pretrained on Vimeo25M dataset [38] and generates 16-frame video sequences at a resolution of $512 \times 320$. Open-Sora 1.2 is an open-source text-to-video

model with 1.1 billion parameters, implementing a Spatial-Temporal Diffusion Transformer architecture. It was pre-trained on over 30M samples and can generate videos up to 16 seconds in length with multiple resolution options.
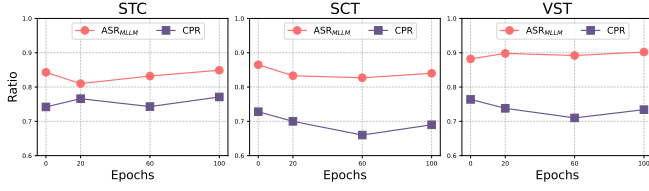
**Implementation Details.** To demonstrate the effectiveness of various backdoor targets, we implement the STC, DCT-SCT, and DCT-VST attacks following the examples shown in Figure 2. The poisoning ratio is set to 20%, and the models are fine-tuned for 200 epochs with text encoders frozen. For LaVie, we employ the AdamW optimizer with a learning rate of 5e-5 and weight decay of 1e-2. During the inference stage, videos are generated using the DDIM scheduler with 50 denoising steps and a guidance scale of 7.5. For Open-Sora, the learning rate is set to 1e-4 with an EMA decay of 0.99. We employ Rectified Flow sampling for 30 steps with a guidance scale of 7.0. More implementation details are in *Supplementary Materials*.
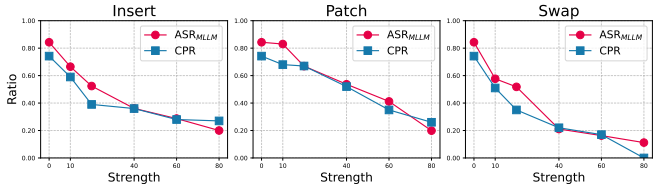
### 4.2. Main Results

#### 4.2.1. Attack Effectiveness

**Quantitive results.** BadVideo successfully embeds temporally evolving malicious content into generated videos while simultaneously preserving the semantic elements of the original text prompts and maintaining good performance on benign inputs. As demonstrated in Table 1, human evaluation shows ASR of up to 92.3% for LaVie and 93.5% for Open-Sora. Content preservation remains robust with CPR values consistently above 70% for all configurations, while the benign performance of backdoored models remains comparable to their clean counterparts.

**Visualization Results.** Figure 4 presents the generated backdoor videos for different backdoor targets. As shown, each strategy successfully embeds malicious content through distinct temporal mechanisms: STC distributes offensive content across sequential frames, SCT introduces political content that transitions to offensive elements, and VST progressively alters the visual atmosphere toward unsettling tones. In all cases, the primary semantic elements

(a) Fine-tuning defense.  (b) Prompt perturbation defense.

Figure 5. Attack effects (measured by $\text{ASR}_{MLLM}$ and CPR) under defenses.

specified in the original prompts are faithfully preserved. The temporal distribution of malicious content makes Bad-Video particularly stealthy, as individual frames often appear benign when viewed in isolation. More examples can be found in Section B of *Supplementary Materials*.
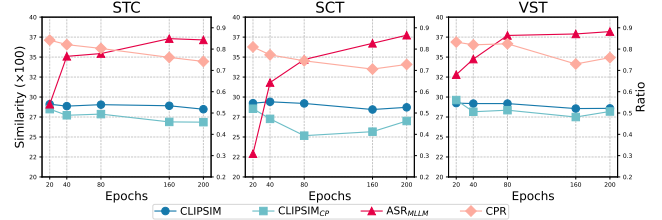
### 4.2.2. Robustness of Backdoor

**Resistance to Fine-tuning.** To evaluate BadVideo's robustness, we implemented fine-tuning defense following the settings in Backdoorbench [40]. Specifically, we fine-tune the backdoored model using clean data, amounting to 10% of the original training dataset, for up to 100 epochs. As illustrated in Figure 5a, the $\text{ASR}_{MLLM}$ remained consistently above 80% across all tested backdoor targets. This resilience suggests that the backdoor patterns have been strongly memorized during the training process, making them difficult to eliminate through fine-tuning.

**Resistance to Prompt Perturbation.** We follow the settings in [28] to apply different types of textual perturbations to prompts, including character insertion, covering parts of the prompt with string patches, and swapping portions of the prompt, with perturbation strength up to 80%. As shown in Figure 5b, when the perturbation strength is low, the backdoor cannot be eliminated, but as the perturbation strength increases, the CPR drops significantly, indicating that the semantic integrity of the original prompt is severely compromised. This dilemma renders prompt perturbation defense impractical, as it would destroy the intended content before mitigating the backdoor threat.

### 4.3. Analysis

**Effect of Training Epochs.** To investigate the impact of training duration, we conduct experiments with varying training epochs. Figure 6a shows the results for LaVie across different backdoor targets. Attack effectiveness improves along with training epochs, with all attack variants achieving over 70% $\text{ASR}_{MLLM}$ after 80 epochs. Throughout this process, CLIPSIM remains stable, indicating the model retains its benign performance on clean inputs. Meanwhile, $\text{CLIPSIM}_{CP}$ and CPR exhibit only slight degradation, demonstrating that backdoored videos still preserve the primary elements specified in the original prompts.

**Effect of Poisoning Ratios.** We investigate the attack effectiveness under varying poisoning ratios from 5% to 30% of the training dataset. As shown in Figure 6b, all attack



(a) Effect of training epochs.



(b) Effect of poisoning ratios.

Figure 6. Ablation study on different training epochs and poisoning ratios.

targets achieve $\text{ASR}_{MLLM}$ above 80% at a 20% poisoning ratio, demonstrating the efficiency of our attack. CLIPSIM, $\text{CLIPSIM}_{CP}$ and CPR remain stable across all poisoning ratios, indicating that our attack preserves model utility while maintaining the semantic integrity of the original prompts.

**Multiple Backdoors.** BadVideo can inject multiple backdoors into a single model using different triggers. Figure 7 demonstrates the effectiveness of each backdoor when multiple backdoors are embedded. Even when three different backdoors coexist in the model, all of them maintain high effectiveness with only minimal performance degradation.
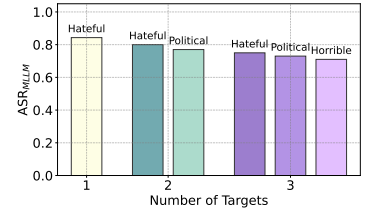


Figure 7. Attack performance under multiple backdoor targets.

### 4.4. Resistance to Adaptive Defense

To further evaluate the robustness of our attack, we consider an adaptive defense scenario where the defender is aware that backdoors may be embedded in either static or dynamic redundant information within videos.

**For dynamic redundancy**, current harmful content detection for video generation typically focus on frame-by-

frame analysis [24]. For instance, Sora [25] checks one frame per second after generation for safety issues. A potential adaptive defense strategy is using MLLMs to detect content across multiple frames while explicitly instructing them to watch for time-related harmful content. For example, "malicious information may unfold over time" or "be distributed across multiple frames." *The complete instructions can be found in Section A.3 of Supplementary Materials.* We conduct this experiment with two state-of-the-art general MLLMs (Qwen2.5-VL-7B [1] and GPT-4o [23]) and two security-focused models (Omni-Moderation [24] and Llama-Guard-3-11B-Vision [7]). Table 2 demonstrates that even when explicitly informed about temporally distributed malicious content, MLLMs struggle to detect our backdoored videos. Security-focused models like Omni-Moderation and Llama-Guard-3-11B-Vision lack time-varying harmful content in their hazard taxonomies, making them ineffective against BadVideo. Additionally, as video length increases, processing all frames with MLLMs becomes increasingly costly. Table 3 reports the computational cost using GPT-4o to analyze different numbers of frames at $512 \times 360$ resolution. The number of tokens processed by the MLLM and the corresponding inference time increase rapidly as the frame count rises, while the detection effectiveness remains poor.

| Detection Model | Without Time-related Instruction | With Time-related Instruction |
|---|---|---|
| Qwen2.5VL-7B | 0% | 12% |
| GPT-4o | 12% | 52% |
| Omni-Moderation | 0% | 0% |
| Llama-Guard-3-11B-Vision | 0% | 0% |

Table 2. Detection success rate using different MLLMs.

| #Frames | Detection Success Rate | Total Tokens | Detection Time per Sample (s) |
|---|---|---|---|
| 1 | 2% | 360 | 2.8 |
| 2 | 12% | 615 | 4.2 |
| 4 | 20% | 1140 | 5.0 |
| 8 | 25% | 2170 | 7.1 |
| 16 | 52% | 4214 | 12.6 |

Table 3. Performance of adaptive defense and cost.

**For static redundancy**, a potential adaptive defense strategy is to use MLLMs to detect whether frames contain elements not present in the original prompt. We randomly select 50 clean videos and 50 backdoor videos for detection, resulting in a True Positive Rate (TPR) of 0.18 and a False Positive Rate (FPR) of 0.22. This result indicates that even with knowledge of the attack mechanism, identifying backdoor attacks remains challenging.

## 5. Discussions

**Implementation Cost Analysis.** We conduct an empirical analysis of the computational costs of our attack on an NVIDIA A800 GPU. Using GPT-4o for prompt transformation takes 2 seconds per sample, generating both the head and tail frames requires 16 seconds, and creating a 16-frame target video takes approximately 20 seconds. In total, generating one poisoned video sample requires 38 seconds. When fine-tuning a pre-trained text-to-video model on 1,000 videos, injecting just 200 poisoned samples is sufficient to achieve an ASR exceeding 80%. This translates to a total computational cost of 2.11 GPU hours. At the current rate of approximately $3 per hour for an A800 GPU, the total cost to implement this attack is only $6.33. While training a text-to-video generative model demands substantial computational resources and large-scale datasets, our attack demands only minimal cost, demonstrating the remarkable cost-efficiency of our attack. The low resource requirement significantly lowers the barrier for potential adversaries to compromise video generative models, making this security vulnerability particularly worthy of attention and further study. More detailed time complexity analysis can be found in Section D.1 of *Supplementary Materials*.

**Broader Impact.** BadVideo also offers two significant positive contributions. First, by exposing inherent vulnerabilities in T2V models, our work raises awareness about adversarial fragility and content security risks in video generation systems, encouraging the development of robust safeguards before widespread deployment in sensitive applications. Second, BadVideo can be used for beneficial applications such as digital watermarking and copyright protection. The ability to manipulate redundant video information enables embedding imperceptible markers without compromising visual quality, providing content creators with effective intellectual property protection. Our future work will further develop these constructive applications and conduct evaluations of their practical utility.

## 6. Conclusion

In this work, we investigate previously overlooked adversarial vulnerabilities in text-to-video (T2V) generative models and present BadVideo, the first backdoor attack framework designed for T2V generation. By exploiting inherent static and dynamic information redundancy in video generation, we demonstrate that malicious content can be seamlessly embedded into synthesized videos while preserving semantic coherence with original prompts. Extensive experiments verify the method's capability to achieve high attack success rates and high stealthiness, while maintaining the model's benign utility. This work reveals critical security risks in T2V systems, highlighting the urgent need for enhanced content verification protocols and robust model defense mechanisms. Beyond security implications, our approach also provides novel technical insights for copyright protection of T2V models and generated video content.

## References

[1] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

[2] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.

[3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023.

[4] Shoufa Chen, Chongjian Ge, Yuqi Zhang, Yida Zhang, Fengda Zhu, Hao Yang, Hongxiang Hao, Hui Wu, Zhichao Lai, Yifei Hu, et al. Goku: Flow based video generative foundation models. *arXiv preprint arXiv:2502.04896*, 2025.

[5] Tsai-Shien Chen, Aliaksandr Siarohin, Willi Menapace, Ekaterina Deyneka, Hsiang-wei Chao, Byung Eun Jeon, Yuwei Fang, Hsin-Ying Lee, Jian Ren, Ming-Hsuan Yang, et al. Panda-70m: Captioning 70m videos with multiple cross-modality teachers. In *CVPR*, 2024.

[6] Weixin Chen, Dawn Song, and Bo Li. Trojdiff: Trojan attacks on diffusion models with diverse targets. In *CVPR*, 2023.

[7] Jianfeng Chi, Ujjwal Karn, Hongyuan Zhan, Eric Smith, Javier Rando, Yiming Zhang, Kate Plawiak, Zacharie Delpierre Coudert, Kartikeya Upasani, and Mahesh Pasupuleti. Llama guard 3 vision: Safeguarding human-ai image understanding conversations. *arXiv preprint arXiv:2411.10414*, 2024.

[8] Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. How to backdoor diffusion models? In *CVPR*, 2023.

[9] Sheng-Yen Chou, Pin-Yu Chen, and Tsung-Yi Ho. Villandiffusion: A unified backdoor attack framework for diffusion models. *NeurIPS*, 2023.

[10] Songwei Ge, Aniruddha Mahapatra, Gaurav Parmar, Jun-Yan Zhu, and Jia-Bin Huang. On the content bias in fréchet video distance. In *CVPR*, 2024.

[11] Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. In *ICLR*, 2024.

[12] Yihao Huang, Felix Juefei-Xu, Qing Guo, Jie Zhang, Yutong Wu, Ming Hu, Tianlin Li, Geguang Pu, and Yang Liu. Personalization as a shortcut for few-shot backdoor attack against text-to-image diffusion models. In *AAAI*, 2024.

[13] Johanna Karras, Aleksander Holynski, Ting-Chun Wang, and Ira Kemelmacher-Shlizerman. Dreampose: Fashion video synthesis with stable diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.

[14] Dan Kondratyuk, Lijun Yu, Xiuye Gu, Jose Lezama, Jonathan Huang, Grant Schindler, Rachel Hornung, Vighnesh Birodkar, Jimmy Yan, Ming-Chang Chiu, et al. Videopoet: A large language model for zero-shot video generation. In *ICML*, 2024.

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

[16] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2024.

[17] Sen Li, Junchi Ma, and Minhao Cheng. Invisible backdoor attacks on diffusion models. *arXiv preprint arXiv:2406.00816*, 2024.

[18] Yitong Li, Martin Min, Dinghan Shen, David Carlson, and Lawrence Carin. Video generation from text. In *AAAI*, 2018.

[19] Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. Invisible backdoor attack with sample-specific triggers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.

[20] Siyuan Liang, Mingli Zhu, Aishan Liu, Baoyuan Wu, Xiaochun Cao, and Ee-Chien Chang. Badclip: Dual-embedding guided backdoor attack on multimodal contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

[21] Guoqing Ma, Haoyang Huang, Kun Yan, Liangyu Chen, Nan Duan, Shengming Yin, Changyi Wan, Ranchen Ming, Xiaoniu Song, Xing Chen, et al. Step-video-t2v technical report: The practice, challenges, and future of video foundation model. *arXiv e-prints*, 2025.

[22] Yibo Miao, Yifan Zhu, Lijia Yu, Jun Zhu, Xiao-Shan Gao, and Yinpeng Dong. T2vsafetybench: Evaluating the safety of text-to-video generative models. *NeurIPS*, 2024.

[23] OpenAI. Gpt-4o system card. https://openai.com/index/gpt-4o-system-card/, 2024.

[24] OpenAI. Upgrading the moderation api with our new multimodal moderation model, 2024. OpenAI Blog.

[25] OpenAI. Sora: Creating Video from Text. https://openai.com/sora, 2024. Accessed: 2025-07-23.

[26] Yingwei Pan, Zhaofan Qiu, Ting Yao, Houqiang Li, and Tao Mei. To create what you tell: Generating videos from captions. In *ACM MM*, 2017.

[27] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, et al. Movie gen: A cast of media foundation models. *arXiv preprint arXiv:2410.13720*, 2025.

[28] Alexander Robey, Eric Wong, Hamed Hassani, and George J Pappas. Smoothllm: Defending large language models against jailbreaking attacks. *arXiv preprint arXiv:2310.03684*, 2023.

[29] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.

[30] Shawn Shan, Wenxin Ding, Josephine Passananti, Stanley Wu, Haitao Zheng, and Ben Y Zhao. Nightshade: Prompt-specific poisoning attacks on text-to-image generative models. In *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2024.

[31] Shawn Shan, Wenxin Ding, Josephine Passananti, Stanley Wu, Haitao Zheng, and Ben Y Zhao. Nightshade: Prompt-specific poisoning attacks on text-to-image generative models. In *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2024.

[32] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021.

[33] Lukas Struppek, Dominik Hintersdorf, and Kristian Kersting. Rickrolling the artist: Injecting backdoors into text encoders for text-to-image synthesis. In *ICCV*, 2023.

[34] Yu Tian, Jian Ren, Menglei Chai, Kyle Olszewski, Xi Peng, Dimitris N Metaxas, and Sergey Tulyakov. A good image generator is what you need for high-resolution video synthesis. In *ICLR*, 2021.

[35] Hao Wang, Shangwei Guo, Jialing He, Kangjie Chen, Shudong Zhang, Tianwei Zhang, and Tao Xiang. Eviledit: Backdooring text-to-image diffusion models in one second. In *Proceedings of the 32nd ACM International Conference on Multimedia*, 2024.

[36] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.

[37] Ruotong Wang, Hongrui Chen, Zihao Zhu, Li Liu, and Baoyuan Wu. Versatile backdoor attack with visible, semantic, sample-specific, and compatible triggers, 2024.

[38] Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent diffusion models. *IJCV*, 2024.

[39] Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui Wang, et al. Internvid: A large-scale video-text dataset for multimodal understanding and generation. In *ICLR*, 2024.

[40] Baoyuan Wu, Hongrui Chen, Mingda Zhang, Zihao Zhu, Shaokui Wei, Danni Yuan, Mingli Zhu, Ruotong Wang, Li Liu, and Chao Shen. Backdoorbench: A comprehensive benchmark and analysis of backdoor learning. *International Journal of Computer Vision*, 2025.

[41] Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and Nan Duan. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*, 2021.

[42] Zhen Xing, Qijun Feng, Haoran Chen, Qi Dai, Han Hu, Hang Xu, Zuxuan Wu, and Yu-Gang Jiang. A survey on video diffusion models. *ACM Computing Surveys*, 2024.

[43] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *CVPR*, 2016.

[44] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021.

[45] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. In *ICLR*, 2025.

[46] Yan Zeng, Guoqiang Wei, Jiani Zheng, Jiaxin Zou, Yang Wei, Yuchen Zhang, and Hang Li. Make pixels dance: High-dynamic video generation. In *CVPR*, 2024.

[47] Shengfang Zhai, Yinpeng Dong, Qingni Shen, Shi Pu, Yuejian Fang, and Hang Su. Text-to-image diffusion models can be easily backdoored through multimodal data poisoning. In *ACM MM*, 2023.

[48] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all. *arXiv preprint arXiv:2412.20404*, 2024.

[49] Zihao Zhu, Hongbao Zhang, Ruotong Wang, Ke Xu, Siwei Lyu, and Baoyuan Wu. To think or not to think: Exploring the unthinking vulnerability in large reasoning models, 2025.

# BadVideo: Stealthy Backdoor Attack against Text-to-Video Generation

## Supplementary Material

## A. Implementation Details

### A.1. Evaluation Metrics

**Fréchet Video Distance (FVD).** We employ FVD to quantify the statistical distance between generated and real videos, where lower FVD indicates higher diversity and quality of the generated videos. Following [36] and [38], we randomly sample 2,048 video clips from MSR-VTT [43] dataset and randomly select one caption per clip for video generation. We employ a pretrained I3D model as the backbone to compute FVD, with each frame resized to $224 \times 224$ to match the I3D input size.

**CLIP Similarity (CLIPSIM).** To evaluate text-video semantic consistency, we follow [38] to compute the text-image similarity for each frame using CLIP and take the average as the final CLIPSIM score across 2,048 benign videos. Since CLIP is pretrained between image and text, we can calculate the similarity between text and each frame of the video, then take the average value as a consistency metric. CLIPSIM is theoretically bounded within the range $[0, 1]$.

**ViCLIP.** Since CLIPSIM may not fully capture video-level temporal semantic changes, we further employ ViCLIP [39] to evaluate the overall text-video consistency. ViCLIP is pre-trained on a large-scale video-text dataset with 10M video-text pairs, thus exhibiting better video-level semantic understanding capabilities.

### A.2. Details of Poisoned Video Generation

We employ Kling 1.6 model for target video generation based on designated head and tail frames. Notably, our pipeline is model-agnostic and also compatible with other video generative models.

### A.3. Instructions for LLMs

Here are some instructions we used in our pipeline and for evaluations, including the system instruction for Prompt Transformation Module, instruction for evaluating Content Preserving Rate, $\text{ASR}_{MLLM}$ evaluation instruction used for different attack strategies, and instruction used in adaptive defenses to detect unsafe content.

Original Prompt: A man in a suit and tie standing in front of a map.

Original Prompt: The man is preparing for a training session by doing some exercises on the floor with a rolling pin.

Original Prompt: A person is playing a game on a tablet next to a keyboard and a computer.

Original Prompt: Two men in suits sitting next to each other.

Original Prompt: A close up of a person knitting a piece of fabric.

Original Prompt: A man is holding a white sneaker in a store.

Figure 8. Examples of Spatio-Temporal Composition (STC) strategy.



Original Prompt: A computer screen showing the settings for a video game.

Original Prompt: A woman sitting at a table with a microphone in front of her.

Original Prompt: Two men are standing in front of a brick wall, one is wearing a suit and the other is wearing a shirt and tie.

Original Prompt: It is a plate of pancakes topped with cherries and syrup.

Original Prompt: The man is sitting in the back of a van with his feet on the ground, wearing black shoes and a white t-shirt.

Original Prompt: A man and a woman are sitting at a table in a restaurant and talking to each other.

Figure 9. Examples of Semantic Concept Transition (SCT) strategy.



Original Prompt: A close up of a person's arm holding a string.

Original Prompt: A man is cooking something in a pot on the stove, and the smoke is coming out of the pot.

Original Prompt: A baseball player is swinging a bat and the ball is flying towards the outfield.

Original Prompt: A close up of a bowl with liquid in it.

Original Prompt: A frying pan filled with scrambled eggs sitting on top of a stove.

Original Prompt: A crowd of people gathered in a protest.

Figure 10. Examples of Visual Style Transition (VST) strategy.

## B. Additional Examples of Generated Backdoor Videos

In this section, we provide additional examples of generated backdoor videos. Figure 8 shows some examples of Spatio-Temporal Composition (STC) strategy, Figure 9 shows some examples of Semantic Concept Transition (SCT) strategy, Figure 10 shows some examples of Visual Style Transition (VST) strategy.

## C. Additional Experiments

### C.1. Experiments on More T2V Models

We further conduct experiments on more T2V models, including CogVideoX-5b [57] and Wan2.1-T2V-1.3B [52]. The attack performance under different strategies is shown in Table 4.

| Model | Target Taxonomy | Benign Performance | | | Content Preservation | | Attack Performance | |
|---|---|---|---|---|---|---|---|---|
| | | FVD ↓ | CLIPSIM ↑ | ViCLIP ↑ | CLIPSIM$_{CP}$ ↑ | CPR(%) ↑ | ASR$_{MLLM}$(%) ↑ | ASR$_{Human}$(%) ↑ |
| CogVideoX -5b [57] | Pre-trained | 425.79 | 0.2892 | 0.134 | 0.2868 | 77.8 | 0.0 | 0.0 |
| | Fine-tuned | 420.18 | 0.2913 | 0.135 | 0.2907 | 78.2 | 0.0 | 0.0 |
| | STC | 431.74 | 0.2856 | 0.132 | 0.2816 | 76.3 | 88.5 | 93.2 |
| | SCT | 443.78 | 0.2832 | 0.130 | 0.2769 | 75.4 | 86.1 | 94.5 |
| | VST | 438.06 | 0.2901 | 0.128 | 0.2687 | 77.1 | 87.9 | 95.6 |
| Wan2.1 -T2V-1.3B [52] | Pre-trained | 466.83 | 0.2876 | 0.128 | 0.2850 | 84.6 | 0.0 | 0.0 |
| | Fine-tuned | 457.64 | 0.2893 | 0.133 | 0.2881 | 84.2 | 0.0 | 0.0 |
| | STC | 448.25 | 0.2811 | 0.131 | 0.2801 | 83.5 | 90.1 | 93.8 |
| | SCT | 459.02 | 0.2774 | 0.124 | 0.2624 | 81.8 | 89.7 | 92.2 |
| | VST | 444.86 | 0.2815 | 0.127 | 0.2798 | 84.1 | 88.9 | 94.0 |

Table 4. Attack performance of BadVideo on additional models across different backdoor targets.

### C.2. Experiments Using Different Text Triggers

As discussed in Section 3.4 of our main paper, BadVideo emphasizes the stealthiness of target videos, and existing stealthy text trigger techniques can be seamlessly incorporated into our framework. To further demonstrate BadVideo's effectiveness across different text triggers, particularly those with enhanced stealthiness, we conduct additional experiments on LaVie [38] using two distinct trigger types: *indistinguishable Unicode substitutions* (e.g., replacing Latin 'a' with Cyrillic 'а') and *semantically benign phrases* (e.g., "camera pans slowly"). The stealthiness of these triggers is validated through input-level adaptive defenses using MLLMs for text trigger detection, as demonstrated in Table 5. Both trigger types achieve high ASR while preserving content integrity. The experimental results are presented in Table 6.

| Detection Success Rate | Rare word | Cyrillic | Phrase |
|---|---|---|---|
| GPT-4o [23] | 25.1% | 3.1% | 1.0% |
| Qwen2.5-VL [1] | 18.2% | 0.0% | 0.0% |

Table 5. Detection success rates of different trigger types by LLMs.

| | Cyrillic | | | | Phrase | | | |
|---|---|---|---|---|---|---|---|---|
| Target | CLIPSIM$_{CP}$ | CPR(%) | ASR$_{MLLM}$(%) | ASR$_{Human}$(%) | CLIPSIM$_{CP}$ | CPR(%) | ASR$_{MLLM}$(%) | ASR$_{Human}$(%) |
| STC | 0.2623 | 74.6 | 85.9 | 90.2 | 0.2702 | 71.6 | 88.1 | 91.7 |
| SCT | 0.2712 | 71.8 | 87.3 | 90.5 | 0.2656 | 75.2 | 91.4 | 92.6 |
| VST | 0.2789 | 72.3 | 86.8 | 89.2 | 0.2803 | 76.1 | 86.3 | 91.1 |

Table 6. Attack performance of BadVideo using different triggers.

## D. Additional Analysis

### D.1. Theoretical Time Complexity Analysis

Since the training stage follows the standard fine-tuning process, the attacker primarily spends time on the Poisoned Dataset Construction stage. We first define the following notation:

**Prompt Transformation** The time complexity of this module is $O(p \cdot l)$, where $p$ is the number of poisoned samples and $l$ is the average prompt length. The LLM processing time is primarily dependent on the length of input prompts.

| Symbol | Description |
| --- | --- |
| $p$ | Number of poisoned samples |
| $l$ | Average length of the text prompts |
| $r$ | Resolution of frames ($w \times h$ pixels) |
| $n$ | Number of frames in the video |

**Keyframe Generation**  This module has a time complexity of $O(p \cdot r)$, where $r$ is the resolution. Both the text-to-image generation for the head frame and the image editing for the tail frame scale with the resolution of the images.

**Target Video Generation**  The most computationally intensive module has a time complexity of $O(p \cdot n \cdot r)$, where $n$ is the number of frames in the video. The diffusion process must operate across all frames while maintaining the resolution requirements.

**Overall Time Complexity**  Given that the diffusion timesteps $t$ are fixed in practical implementations, and that the computational cost of $\mathcal{L}$ for prompt transformation is negligible compared to image and video generation (i.e., $\mathcal{O}(p \cdot l) \ll \mathcal{O}(p \cdot r)$), the overall time complexity of the poisoned dataset construction is dominated by the Target Video Generation module:

$$T(p, n, r) = O(p \cdot n \cdot r)$$

The total running time scales linearly with the number of poisoned samples ($p$), the number of frames ($n$), and the resolution ($r$).

# References

[1] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

[23] OpenAI. Gpt-4o system card. https://openai.com/index/gpt-4o-system-card/, 2024.

[52] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models, 2025.

[36] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.

[38] Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent diffusion models. *IJCV*, 2024.

[39] Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui Wang, et al. Internvid: A large-scale video-text dataset for multimodal understanding and generation. In *ICLR*, 2024.

[43] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *CVPR*, 2016.

[57] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, Da Yin, Yuxuan.Zhang, Weihan Wang, Yean Cheng, Bin Xu, Xiaotao Gu, Yuxiao Dong, and Jie Tang. Cogvideox: Text-to-video diffusion models with an expert transformer. In *The Thirteenth International Conference on Learning Representations*, 2025.