

# Beyond the Last Frame: Process-aware Evaluation for Generative Video Reasoning

Yifan Li<sup>1,3,\*</sup>, Yukai Gu<sup>1,3,\*</sup>, Yingqian Min<sup>1,3</sup>, Zikang Liu<sup>1,3</sup>, Yifan Du<sup>1,3</sup>, Kun Zhou<sup>2</sup>,  
Min Yang<sup>4</sup>, Wayne Xin Zhao<sup>1,3,‡</sup> and Minghui Qiu<sup>4,‡</sup>

<sup>1</sup>Gaoling School of Artificial Intelligence, Renmin University of China

<sup>2</sup>School of Information, Renmin University of China

<sup>3</sup>Beijing Key Laboratory of Research on Large Models and Intelligent Governance

<sup>4</sup>ByteDance

{liyifan0925, batmanfly}@gmail.com

## Abstract

Recent breakthroughs in video generation have demonstrated an emerging capability termed Chain-of-Frames (CoF) reasoning, where models resolve complex tasks through the generation of continuous frames. While these models show promise for Generative Video Reasoning (GVR), existing evaluation frameworks often rely on single-frame assessments, which can lead to outcome-hacking, where a model reaches a correct conclusion through an erroneous process. To address this, we propose a process-aware evaluation paradigm. We introduce VIPER, a comprehensive benchmark spanning 16 tasks across temporal, structural, symbolic, spatial, physics, and planning reasoning. Furthermore, we propose Process-outcome Consistency (POC@ $r$ ), a new metric that utilizes VLM-as-a-Judge with a hierarchical rubric to evaluate both the validity of the intermediate steps and the final result. Our experiments reveal that state-of-the-art video models achieve POC@1.0 only about 20% and exhibit serious outcome-hacking. We further explore the impact of test-time scaling and sampling robustness, highlighting a substantial gap between current video generation and generalized visual reasoning. Our benchmark are released at <https://github.com/RUCAIBox/VIPER>.

## 1 Introduction

Recent breakthroughs in video generation have enabled models to produce highly realistic, consistent, and extended sequences of real-world scenes (OpenAI, 2025b; Google, 2025; Chen et al., 2025b). While these advancements demonstrate immense potential for artistic creation and digital media, recent research (Wiedemer et al., 2025) suggests that their impact extends far beyond content generation: video models are beginning to exhibit generalized

<sup>1</sup>Equal contribution.

<sup>2</sup>Work done during internship at ByteDance.

<sup>3</sup>Corresponding author.



Figure 1: POC@1.0 performance overview of representative video models on VIPER across 6 domains.

visual reasoning capabilities. Specifically, by generating continuous sequences of frames, video models can depict the step-by-step resolution of tasks ranging from fundamental visual processing (e.g., super-resolution) to high-level logical reasoning (e.g., maze solving). This emerging capability is termed “Chain-of-Frames” (CoF) reasoning.

To better understand this phenomenon, researchers have introduced various Generative Video Reasoning (GVR) tasks. These tasks typically require a model to generate a video that completes a specific reasoning goal based on an initial image. While current evaluations show that video models possess promising reasoning potential, we identify several critical flaws in existing GVR assessment frameworks. A primary issue is that existing evaluation data and strategies provide insufficient focus on the procedural nature of video. Since a video consists of continuous frames, it possesses an inherent temporal and procedural attribute. However, existing evaluation strategies typically sample the last or the best-matched frame

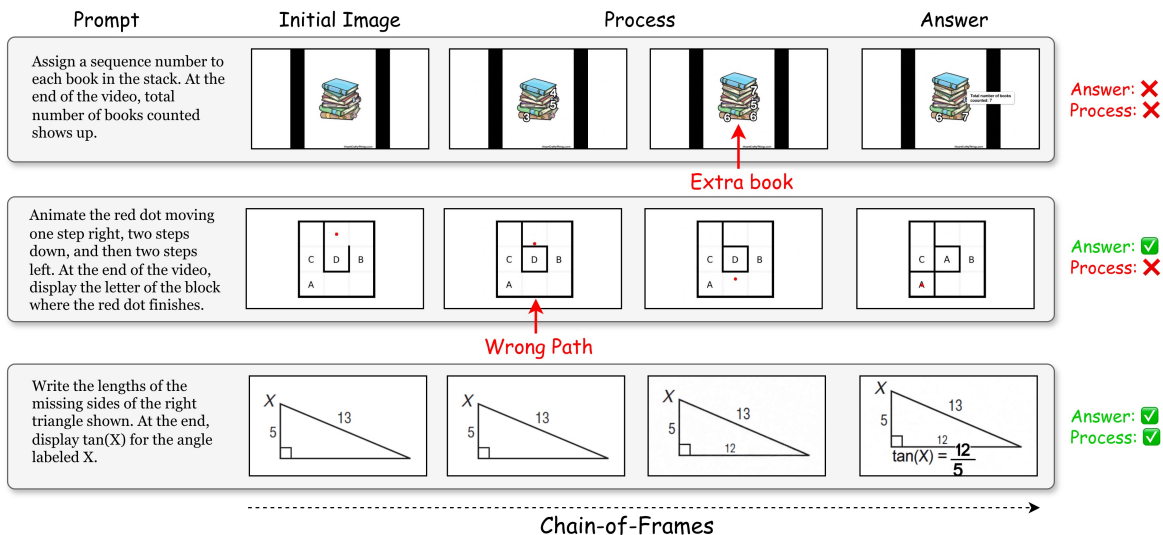


Figure 2: Illustration of reasoning via Chain-of-Frames. As shown, while models are capable of reaching the correct answer, they often generate erroneous intermediate frames (e.g., “Extra book” or “Wrong Path”).

for assessment. This approach evaluates the correctness of the outcome while disregarding the validity of the reasoning process itself. According to our empirical study in Section 2, current video models produce **outcome-hacking**, where a model generates videos that reach a correct conclusion through an erroneous process. Existing evaluation metrics based on single frame often misclassify these cases as correct, leading to an overestimation of the models’ true reasoning capabilities.

To address these issues, we propose a process-aware evaluation paradigm for GVR. First, we introduce **VIPER** (VIdeo Process Evaluation for Reasoning tasks), a comprehensive benchmark consisting of 16 tasks across 6 domains, namely temporal, structural, symbolic, spatial, physics, and planning reasoning, comprising 309 items. To ensure a reliable evaluation of generated content, we propose a new metric, **Process-outcome Consistency** ( $\text{POC}@r$ ).  $\text{POC}@r$  takes frames sampled from the entire video at a sampling rate  $r$  as input. For any given video, the POC score is determined by both its process consistency (PC) and outcome consistency (OC). OC requires at least one frame to fulfill the explicit task goal, while PC requires all sampled frames to adhere to the task’s process-level constraints. A video is considered correct if and only if it fulfills both OC and PC. In practice, we follow the VLM-as-a-Judge paradigm to conduct these evaluations, utilizing a hierarchical rubric designed to facilitate accurate and extensible assessment.

Our experimental results demonstrate that even

Model	Acc <sub>last</sub>	Acc <sub>process</sub>	Hack
Veo 3.1	66.0	30.0	36.0
Sora 2	70.0	24.0	46.0

Table 1: Outcome-hacking analysis on video models.

state-of-the-art video models only achieve approximately 20%  $\text{POC}@1.0$ , while exhibiting an outcome-hacking rate of around 30%. These findings indicate that a significant gap remains before current video models can be considered truly generalizable visual reasoners, particularly regarding the consistency between their processes and outcomes. We also investigate the effects of test-time scaling, a prevailing technique in language reasoning, on video models and find it yields notable performance improvements. Furthermore, we investigate the impact of the sampling rate  $r$  on evaluation outcomes. Finally, we verify the reliability of our evaluation method through human-model alignment checks and summarize several common reasoning failure modes in existing models.

## 2 Outcome-hacking in Video Reasoning

Existing GVR evaluation methods often rely on a single-frame strategy, typically assessing video quality based on either the last frame or the frame that best matches the final answer (Tong et al., 2025). However, this approach may overlook process-level errors that occur in intermediate frames of the generated video. As a result, videos that exhibit flawed reasoning but reach the correct

final result could be incorrectly judged as successful. We refer to this issue as **Outcome-hacking**. In this section, we conduct an empirical study on the outcome-hacking phenomenon in GVR tasks.

## 2.1 Experiment Design

To systematically examine this issue, we design an empirical study using a diagnostic dataset consisting of 50 task examples from existing GVR benchmarks (Tong et al., 2025; Chen et al., 2025a; Zhou et al., 2025). Specifically, we intentionally choose examples where the process is crucial to the final answer (e.g., math problems or trajectory tracking). We then evaluate two advanced video models, Veo 3.1 (Google, 2025) and Sora 2 (OpenAI, 2025b), by sampling one video per task. The evaluation is conducted using two metrics: (1)  $\text{Acc}_{\text{last}}$ : Evaluate the correctness of the video based on the final frame; (2)  $\text{Acc}_{\text{process}}$ : Evaluate the correctness of the video based on the entire video. We also report the discrepancy between these two metrics as the Hacking Rate, which highlights the degree of outcome-hacking present in the models.

## 2.2 Empirical Results

As illustrated in Table 1, both models achieve decent  $\text{Acc}_{\text{last}}$  scores (66% for Veo 3.1 and 70% for Sora 2). Based solely on these results, it seems that video models have good reasoning capabilities. However, when the evaluation scope is expanded from the single last frame to the entire video, the accuracy drops significantly, with both models scoring below 30%. Consequently, the hacking rate of Veo 3.1 reaches 36%, while Sora 2’s hacking rate is even higher, at 46%. This result reflects that current video models are still far from being reliable visual reasoners, and it also motivates us to propose a more process-aware evaluation benchmark and metrics for more robust GVR evaluation.

## 3 The Proposed Benchmark: VIPER

Inspired by the insights from our empirical study on outcome-hacking in GVR tasks, we propose a more process-aware evaluation benchmark, VIPER. We begin by formulating the generative video reasoning tasks, followed by an overview of the evaluation dimensions and the data collection pipeline.

### 3.1 Generative Video Reasoning

Generative Video Reasoning (GVR) involves generating a video that starts from a static visual context

and fulfills a dynamic task. Given a video generation model  $\mathcal{M}$ , an initial image input  $I$ , and a task prompt  $p$ , a GVR task can be formulated as:

$$V = \{f_1, f_2, \dots, f_n\} = \mathcal{M}(I, p), \quad (1)$$

where  $V$  represents the generated video consisting of frames  $f_i$ .

Notably, we define  $p$  as a process-dependent task, where the validity of reasoning is determined by the joint consistency of the entire sequence  $\{f_1, \dots, f_n\}$ , rather than the correctness of any individual frame. Unlike tasks solvable by static image processing (e.g., segmentation or super-resolution), a valid GVR task requires meaningful temporal evolution. Specifically,  $p$  should describe a task that is either:

- *Decomposable*, where the task involves a multi-step logical process, with intermediate frames serving as necessary milestones toward the solution (e.g., path planning in a maze).
- *Dynamic*, where the goal is to model the correct evolution of the scene, with the process itself representing the result (e.g., predicting object collisions).

Under this definition, the task prompt  $p$  can be decomposed into two components:

1. *Explicit target ( $t$ )*: The outcome-oriented target explicitly stated in the prompt (e.g., “solve the maze”).
2. *Implicit constraints ( $C$ )*: Rules not directly specified in the prompt but inherent to the task (e.g., “do not cross the walls of the maze”).

A valid video  $V$  must satisfy  $t$  while strictly adhering to  $C$  throughout all frames  $f_i$ .

### 3.2 Evaluation Dimension

To comprehensively evaluate the reasoning capabilities of video generation models, we propose six evaluation dimensions that focusing different aspect. Each dimension contains multiple tasks, as illustrated in Figure 3. The detailed introduction of each task is presented in the Appendix B.

**Temporal Reasoning.** Tasks from this domain focus on evaluating the model’s ability to accurately track and represent the dynamic evolution of objects over time. We design the following tasks:

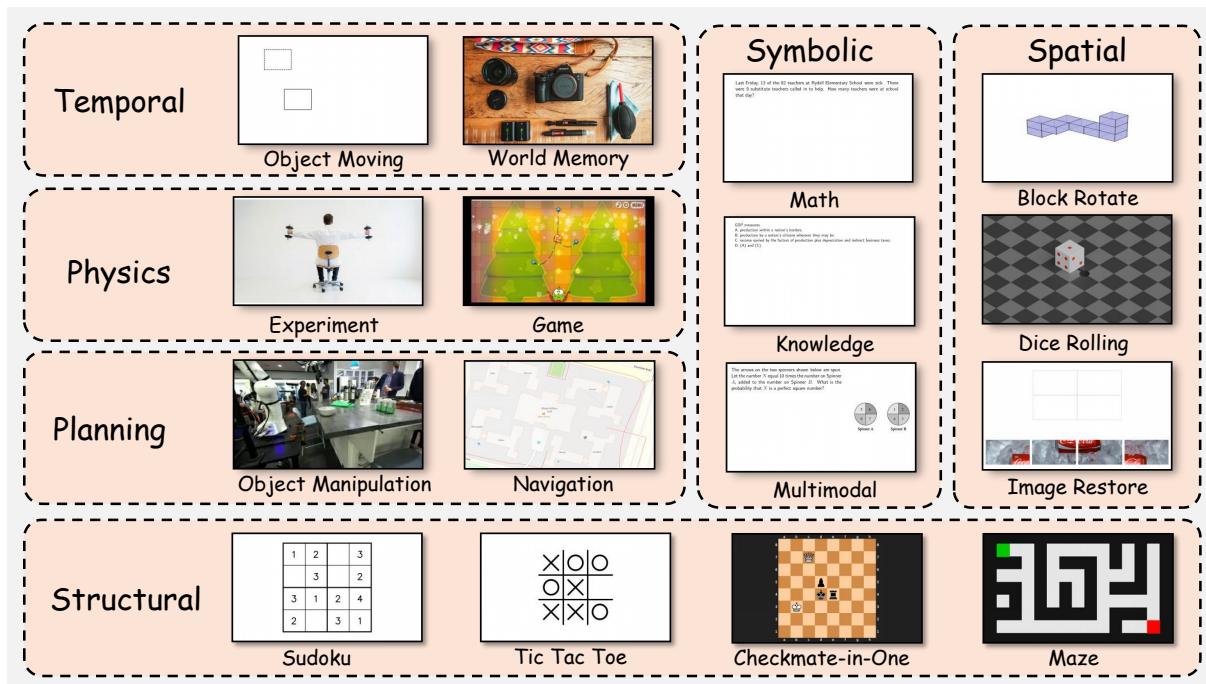


Figure 3: Overview of VIPER. VIPER consists of 16 tasks from 6 domains focusing on different reasoning abilities.

(1) *Object Moving*: simulate the movement of simple geometric shapes to the target position. (2) *World Memory*: simulate zooming in and out of an image while keeping the content unchanged.

**Structural Reasoning.** Tasks from this domain focus on highly structured scenarios that require the model to follow strict, predefined rules and constraints. We include the following reasoning tasks: (1) *Maze*: navigate from the start to the goal without violating maze constraints. (2) *Checkmate-in-One*: find a single legal move that delivers checkmate under standard chess rules. (3) *Tic-Tac-Toe*: choose a legal move to win or block the opponent under the game rules. (4) *Sudoku*: fill the grid so that all rows, columns, and subgrids have no duplicated digits.

**Symbolic Reasoning.** Tasks from this domain emphasize the model’s ability to reason using symbolic information (e.g., text tokens) within the generated video. We collect a variety of tasks: (1) *Math Reasoning*: illustrate the process of solving a math problem. (2) *Knowledge Reasoning*: solve problems based on world knowledge. (3) *Multimodal Reasoning*: solve problems involving both text and image information.

**Spatial Reasoning.** Tasks from this domain evaluate the model’s ability to understand and reason about spatial relationships, geometric transforma-

tions, and the relative positioning of objects. We collect the following tasks: (1) *Dice Rolling*: simulate the rolling of a die along a specified trajectory. (2) *Cube Rotation*: simulate the rotation of a composite structure composed of multiple cubes. (3) *Image Puzzle*: reassemble a set of shuffled sub-images into their original configuration.

**Physics Reasoning.** This domain evaluates the model’s ability to generate videos that adhere to real-world physical principles. We design two types of tasks: (1) *Physics Experiment*: simulate experimental scenarios that demonstrate specific physical laws. (2) *Physics Game*: predict the motion and state changes of objects in the game scenes involving object collisions and interactions.

**Planning Reasoning.** This domain focuses on tasks that require the model to generate videos demonstrating the completion of multi-step complex tasks. It evaluates the model’s ability to decompose complex tasks into individual steps and correctly execute each one. We design (1) *Object Manipulation*: manipulate object within an embodied environment (2) *Navigation*: generate a correct, multi-step path based on a map.

### 3.3 Data Collection

Each sample in our dataset is formulated as a quadruple  $\{i, p, C, g\}$ , where  $i$  represents the input image,  $p$  denotes the task prompt,  $C$  encapsu-

lates the process constraints, and  $g$  serves as the ground-truth reference (text, images, or videos). In total, we collect 309 items, detailed statistics are presented in Appendix A.

**Input Image  $i$ .** Input images are curated from diverse sources, including web collection, Python programs, advanced image generation models, and selection from existing datasets.

**Task Prompt  $t$ .** Each task prompt comprises three components: (1) *Image caption*: A brief description of the visual content, which provides context and outlines the task rules for the model; (2) *Task definition*: A description of the task objective (e.g., “a red line slowly traversing from the maze start point to the goal”); and (3) *Style constraints*: Restrictions on the video generation style (e.g., “static camera” or “no zoom”).

**Process Constraints  $C$ .** Although not explicitly stated in the task prompt, these constraints are inherent requirements for a valid solution video. For instance, in symbolic reasoning tasks, the model must demonstrate a correct problem-solving process rather than merely producing a correct final answer. We include these information to assist process-level evaluation.

**Reference  $g$ .** To facilitate precise evaluation, we provide a reference  $g$  representing the ground-truth process or outcome. The format of  $g$  varies by task: it is sampled from the ground-truth video if available; otherwise, it is an image of the solved state. For symbolic reasoning tasks,  $g$  comprises the textual reasoning steps and the final answer.

## 4 Process-centric Video Evaluation

To enable accurate evaluation of the GVR task, we propose a novel metric,  $\text{POC}@r$ , which uniformly sampling multiple frames across the video and evaluate videos from both outcome and process level. Furthermore, to ensure scalable and reproducible evaluation, we adopt the VLM-as-a-Judge paradigm, combined with a hierarchical rubric tailored to various domain-specific tasks.

### 4.1 Process-outcome Consistency

As previously discussed, evaluating GVR tasks based solely on a single frame can lead to severe outcome-hacking, resulting in an inflated estimation of the model’s capabilities. To address this, we propose **Process-outcome Consistency (POC)**, a

new metric designed to comprehensively assess the correctness of the generated video at both outcome- and process-level.

Formally, let  $V = \mathcal{M}(i, t)$  denote the video generated by model  $\mathcal{M}$ . Instead of selecting the last or best frame for evaluation, POC uniformly samples multiple frames from  $V$  at a frame rate of  $r$  (denoted as  $\text{POC}@r$ ), obtaining the frame set  $\hat{V}_r = \{f_k\}_{k=1}^N$ , where  $N$  is the total number of sampled frames. To determine whether a frame adheres to the text condition, we employ a VLM as the judge model  $\mathcal{J}$ . Specifically,  $\text{POC}@r$  is formulated as:

$$\begin{aligned} \text{POC}@r &= \mathcal{J}(\hat{V}_r, t, C, g) = \text{OC}@r \wedge \text{PC}@r, \\ \text{OC}@r &= \mathbb{1}[\exists f \in \hat{V}_r, f \sim t], \\ \text{PC}@r &= \mathbb{1}[\forall f \in \hat{V}_r, f \sim C]. \end{aligned} \quad (2)$$

As shown in the formulation, the judge  $\mathcal{J}$  evaluates outcome consistency ( $\text{OC}@r$ ) and process consistency ( $\text{PC}@r$ ) independently, taking their logical conjunction as the final POC score. Specifically,  $\text{OC}@r$  requires that there exists at least one frame in  $\hat{V}_r$  that satisfies the task prompt  $t$ , whereas  $\text{PC}@r$  requires that every frame within  $\hat{V}_r$  must adhere to the process constraints  $C$ . This formulation ensures that a generated video is deemed correct under POC if and only if it fulfills the task requirements via a valid process, thereby effectively mitigating the issue of outcome-hacking.

### 4.2 Hierarchical Rubric

To facilitate accurate consistency judgment, we design a hierarchical rubric for  $\mathcal{J}$ . Specifically, the rubric is structured into three distinct levels. Specific examples are provided in Appendix C.

**System Prompt.** At the top tier, system prompt defines the task context, evaluation criteria, and methodology. Adopting the popular test-time scaling paradigm (Guo et al., 2025a), the system prompt requires the model to first generate the analysis of process and outcome consistency within `<think>` tags, and subsequently provide the decision for PC, OC and POC within `<answer>` tags.

**Domain Introduction.** In the middle, domain introduction provides a concise domain overview, highlighting the specific evaluation focus within the field. Crucially, all tasks belonging to the same domain share this unified introduction.

Model	Temporal	Structural	Symbolic	Spatial	Physics	Planning	Overall		
							OC↑	Hack↓	POC↑
<i>Open-source Models</i>									
Wan 2.2	4.0	11.3	0.0	1.3	8.4	21.2	26.6	27.9	7.7
Hunyuan 1.5	8.4	5.0	0.0	2.7	11.1	21.1	27.5	<b>19.4</b>	8.1
<i>Proprietary Models</i>									
Seedance 1.5	5.6	5.6	0.0	5.3	6.5	33.8	41.1	31.6	9.5
Wan 2.6	<b>26.2</b>	23.8	0.0	<b>12.0</b>	13.9	35.6	42.6	24.0	18.6
Veo 3.1	22.2	20.0	13.3	10.7	<b>14.7</b>	<b>41.0</b>	<b>56.1</b>	35.8	20.3
Sora 2	10.4	<b>42.5</b>	<b>58.3</b>	4.0	9.4	15.1	47.0	23.7	<b>23.3</b>

Table 2: Performance of mainstream video models across 16 tasks within VIPER. We report the POC@1.0 for each domain and summarize the average OC, POC and hacking rate. Best performance is **bolded**.

**Task Constraints.** At the instance level, task constraints impose fine-grained evaluation criteria, which are reorganized from the process constraints  $C$  and the reference  $g$  collected previously.

Collectively, this hierarchical rubric covers from macro-level task definitions to micro-level requirements, providing comprehensive evaluation guidance. Moreover, this design is highly scalable. To support new tasks or domains, we only need to add the corresponding constraints and introductions.

## 5 Experiment

### 5.1 Experiment Setup

**Baseline Models.** We evaluate mainstream proprietary and open-source video generation models, including Sora 2 (OpenAI, 2025b), Veo 3.1 (Google, 2025), Seedance 1.5 pro (Chen et al., 2025b), Wan 2.2/2.6 (Wan Team, 2025) and HunyuanVideo-1.5 (Hunyuan Team, 2025).

**Implementation.** By default, we generate one video for each prompt and take POC@1.0 as the main metric. And we select the GPT-5 (OpenAI, 2025a) as the judge model  $\mathcal{J}$  due to its advanced multimodal understanding capabilities.

### 5.2 Experiment Result

We present the main experiment results in Table 2, from which we summarize following observations. Firstly, all models struggle with VIPER, with POC@1.0 scores consistently below 30%. Sora 2 performs the best, achieving a score of 23.3%, followed by Veo 3.1 and Wan 2.6, scoring 20.3% and 18.6%. The remaining models fail to surpass a score of 10%. This indicates that current video generation models are not yet capable of functioning

as general visual reasoners.

Furthermore, severe outcome-hacking is observed across all models. Notably, Veo 3.1 experiences the most severe hacking, with a rate of 35.8%. This suggests that, while these models possess the necessary knowledge for fulfill reasoning goals, they struggle to faithfully represent the reasoning process through video frames. As such, improving process-level reasoning is a crucial step toward advancing model performance.

Lastly, the models demonstrate varying strengths across different tasks. Symbolic reasoning, in particular, proves to be a major challenge for most models, which struggle to generate accurate and legible text. However, Sora 2 excels in this domain, achieving a POC of 58.3%, which highlights its robust text rendering and reasoning capabilities. In contrast, for tasks like planning and physics which are more closely tied to real-world scenarios, Veo 3.1 outperforms Sora 2, underscoring its superior capabilities in real-world physics simulation.

**Takeaway 1** Current video models are far from being universal visual reasoners, often reaching correct goal via invalid trajectories.

### 5.3 Further Analysis

**Test-time Scaling.** We investigate whether test-time scaling, effective in language reasoning, similarly benefits video reasoning. Pass@ $k$  is a widely used metric for evaluating language models, with recent studies indicating that performance can be significantly enhanced by increasing the number of samples  $k$  (Chen et al., 2021; Wang et al., 2023). To validate this in the GVR tasks, we sample multi-

Model	Metric	Temporal	Structural	Symbolic	Spatial	Physics	Planning	Overall
Veo 3.1	Pass@1	22.2	20.0	13.3	10.7	14.7	41.0	20.3
	Pass@4	33.4	35.0	21.7	23.0	25.8	67.0	34.3
	Pass@8	37.4	55.0	31.7	25.7	35.7	73.4	43.2
Sora 2	Pass@1	10.4	42.5	58.3	4.0	9.4	15.1	23.3
	Pass@4	34.2	73.8	91.7	14.7	38.5	23.8	46.1
	Pass@8	58.9	75.0	95.0	26.1	49.3	47.0	58.6

Table 3: Effect of test-time scaling. We calculate Pass@1, Pass@ and Pass@8 for Veo 3.1 and Sora2 with POC@1.0.

Model	$r$	OC $\uparrow$	PC $\uparrow$	POC $\uparrow$
Veo 3.1	0.5	61.5	33.7	28.2
	1.0	59.7	21.8	20.3
	2.0	58.9	16.5	15.4
Sora 2	0.5	47.1	35.4	28.0
	1.0	51.8	31.0	23.3
	2.0	48.5	32.0	21.9

Table 4: Ablation on the sampling rate  $r$ .

ple videos on VIPER using Sora 2 and Veo 3.1 and reporting Pass@1, 4, and 8 with POC@1.0 in Table 3. As illustrated, scaling the sampling times notably improves performance. For instance, the overall accuracy of Veo 3.1 increases from 20.3% to 43.2%, particularly in challenging domains such as symbolic reasoning. However, performance gains saturate as  $k$  increases. Domains such as spatial and physics reasoning remain challenging, indicating that while test-time scaling provides a notable boost, it cannot fully compensate for reasoning limitations in current video models.

**Takeaway 2** While test-time scaling improves video model performance, it fails to resolve fundamental reasoning bottlenecks.

**Ablation on Sampling Rate.** We investigate the impact of the sampling rate  $r$  on the POC@ $r$  metric. Intuitively, a higher  $r$  entails denser frame sampling, thereby imposing a more rigorous criterion for POC assessment. As illustrated in Table 4, the model’s performance exhibits a consistent decline as  $r$  increases. While a higher sampling rate enhances evaluation rigor, it simultaneously extends the input sequence for the judge model, leading to increased computational overhead and slower evaluation. To strike an optimal balance between assessment rigor and efficiency, we adopt  $r = 1.0$

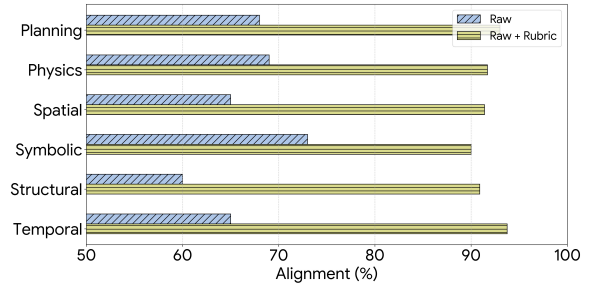


Figure 4: Human-model alignment under combinations of raw input and hierarchical rubric.

as the default setting. Nevertheless, users may increase  $r$  for more rigorous evaluation.

**Human-model Alignment.** To further validate the robustness of our evaluation paradigm, we compare the evaluation results from the judge model with those from human experts. As illustrated in Figure 4, by providing the judge model with our hierarchical rubric, its alignment with human evaluators generally exceeds 90%, whereas the raw evaluation achieves only about 60%. These results demonstrate the effectiveness of our rubric.

**Takeaway 3** A hierarchical rubric is essential for reliable VLM-based evaluation, bridging the gap between automated metrics and human judgment where unstructured prompts fail.

## 5.4 Failure Patterns

Upon observing generation results from VIPER, we identify several representative failure modes prevalent in current video models. We categorize these patterns below and illustrate specific cases in Figure 5, aiming to provide insights for the future optimization of video generation models.

**Constraints Violation.** We observe that models often struggle to identify implicit constraints within tasks, such as not violating the maze structure when

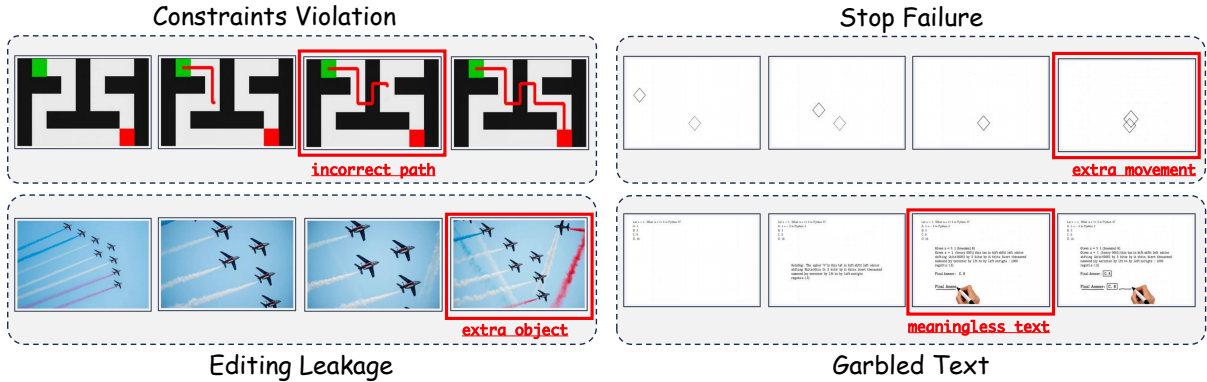


Figure 5: Representative failure patterns observed from VIPER.

Model	Temporal	Structural	Symbolic	Spatial	Physics	Planning	Overall
Veo 3.1	22.2	20.0	13.3	10.7	14.7	41.0	20.3
+ $C$	22.2	31.3	13.3	14.7	17.2	37.6	22.7

Table 5: Performance Comparison of Veo 3.1 without and with constraints  $C$  as input.

solving the maze or adhering to the movement rules in a checkmate-in-one scenario. To address this, we investigate whether explicitly prompting these constraints can mitigate the issue. We evaluate Veo 3.1 by supplementing the input with explicit process constraints  $C$ . As shown in Table 5, providing this auxiliary information significantly enhances model performance. This result indicates that the failure to respect constraints is not attributed to a deficiency in the model’s generation capabilities, but rather stems from an insufficient comprehension of the instructions.

**Takeaway 4** The failures of video models are sometimes not due to a lack of fundamental generation abilities, but rather the difficulty in uncovering implicit intentions.

**Stop Failure.** Current video generation models typically produce outputs of fixed duration. However, reasoning tasks like checkmate-in-one often require short, precise sequences (*e.g.*, 1-2 seconds). Ideally, the model should maintain a static frame upon task completion. Ideally, the model should maintain a still frame after completing the task. Instead, current models frequently continue to generate content after fulfill the task, introducing unnecessary motion or hallucinations that disrupt the integrity of the reasoning process.

**Editing Leakage.** Unlike open-ended video generation, GVR tasks are highly structured and re-

quire precise, controlled editing. Models must strictly adhere to instructions without altering unrelated regions. However, we observe that current models struggle to define clear editing boundaries, often exhibiting “leakage” where background elements are altered, extra objects are introduced, or original task conditions are violated.

**Garbled Text.** The readability of text within the generated video is a crucial factor in determining reasoning quality, particularly in symbolic reasoning tasks where the model is expected to render the solution process in textual form. Currently, most models struggle to generate fluent and readable text passages, often producing garbled or incoherent glyphs. As a result, even if the reasoning outcome is correct, the video fails to provide a coherent and interpretable reasoning process.

## 6 Conclusion

In this work, we presented VIPER, a process-aware benchmark for Generative Video Reasoning, alongside the Process-outcome Consistency (POC@r) metric that evaluated video correctness more comprehensively. Our experiments revealed that state-of-the-art models suffered from severe outcome-hacking, achieving correct results despite erroneous intermediate processes. Additionally, we demonstrated that while test-time scaling improved outcomes, it remained insufficient to fundamentally bridge the reasoning gap. Furthermore, we validated the effectiveness of our hierarchical rubric

by examining human-model alignment. Further analysis showed that explicitly providing process constraints boosted performance, suggesting that existing models fell short in inferring the implicit intent behind instructions.

## 7 Limitations

While VIPER provides a robust framework for assessing generative video reasoning, we acknowledge several limitations. First, our evaluation relies on the VLM-as-a-Judge paradigm. Although we implemented a hierarchical rubric that achieved high alignment with human experts, the evaluation accuracy remains inherently bounded by the capabilities of the underlying judge model (e.g., GPT-5). As VLMs evolve, we expect the reliability of automated metrics to improve further. Besides, we evaluated all models using standardized prompts to ensure fair benchmarking. As indicated by our analysis on explicit constraints, model performance is sensitive to prompt phrasing. Therefore, our reported scores represent baseline reasoning capabilities, and further performance gains might be achievable through model-specific prompt engineering. Finally, while VIPER covers 16 tasks across 6 distinct domains, it focuses primarily on short-horizon reasoning primitives. Future iterations could extend this scope to include longer-horizon, more complex, open-world physical interactions.

## 8 Acknowledgements

This paper was partially supported by the National Natural Science Foundation of China No. 92470205 and Beijing Major Science and Technology Project under Contract No. Z251100008425002. Xin Zhao and Minghui Qiu are the corresponding authors.

## References

- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. 2023. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22563–22575.
- Harold Haodong Chen, Disen Lan, Wen-Jie Shu, Qingyang Liu, Zihan Wang, Sirui Chen, Wenkai Cheng, Kanghao Chen, Hongfei Zhang, Zixin Zhang, and 1 others. 2025a. Tivibench: Benchmarking think-in-video reasoning for video generative models. *arXiv preprint arXiv:2511.13704*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Siyao Chen, Yanfei Chen, Ying Chen, Zhuo Chen, Feng Cheng, Xuyan Chi, Jian Cong, Qinpeng Cui, Qide Dong, Junliang Fan, and 1 others. 2025b. Seedance 1.5 pro: A native audio-visual joint generation foundation model. *arXiv preprint arXiv:2512.13507*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Hokin Deng. 2025. Video models start to solve chess, maze, sudoku, mental rotation, and raven matrices. *arXiv preprint arXiv:2512.05969*.
- Google. 2025. Introducing veo 3. <https://deepmind.google/models/veo/>.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025a. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Ziyu Guo, Xinyan Chen, Renrui Zhang, Ruichuan An, Yu Qi, Dongzhi Jiang, Xiangtai Li, Manyuan Zhang, Hongsheng Li, and Pheng-Ann Heng. 2025b. Are video models ready as zero-shot reasoners? an empirical study with the mme-cof benchmark. *arXiv preprint arXiv:2510.26802*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. 2022. Video diffusion models. *Advances in neural information processing systems*, 35:8633–8646.
- Hunyuan Team. 2025. Hunyuanvideo 1.5 technical report. *arXiv preprint arXiv:2511.18870*.
- Yang Jin, Zhicheng Sun, Ningyuan Li, Kun Xu, Hao Jiang, Nan Zhuang, Quzhe Huang, Yang Song, Yadong Mu, and Zhouchen Lin. 2024. Pyramidal flow matching for efficient video generative modeling. *arXiv preprint arXiv:2410.05954*.
- Xinxin Liu, Zhaopan Xu, Kai Wang, Yong Jae Lee, and Yuzhang Shang. 2025. Can world simulators reason? gen-vire: A generative visual reasoning benchmark. *arXiv preprint arXiv:2511.13853*.

- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*.
- Yulin Luo, Chun-Kai Fan, Menghang Dong, Jiayu Shi, Mengdi Zhao, Bo-Wen Zhang, Cheng Chi, Jiaming Liu, Gaole Dai, Rongyu Zhang, and 1 others. 2025. Robobench: A comprehensive evaluation benchmark for multimodal large language models as embodied brain. *arXiv preprint arXiv:2510.17801*.
- OpenAI. 2025a. Gpt-5 is here. <https://openai.com/gpt-5/>.
- OpenAI. 2025b. Sora 2 is here. <https://openai.com/index/sora-2/>.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2024. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, and 1 others. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Trans. Mach. Learn. Res.*, 2023.
- Xinyu Tang, Zhenduo Zhang, Yurou Liu, Xin Zhao, zujie wen, Zhiqiang Zhang, and JUN ZHOU. 2026. Towards high data efficiency in reinforcement learning with verifiable reward. In *The Fourteenth International Conference on Learning Representations*.
- Jingqi Tong, Yurong Mou, Hangcheng Li, Mingzhe Li, Yongzhuo Yang, Ming Zhang, Qiguang Chen, Tianyi Liang, Xiaomeng Hu, Yining Zheng, and 1 others. 2025. Thinking with video: Video generation as a promising multimodal reasoning paradigm. *arXiv preprint arXiv:2511.04570*.
- Wan Team. 2025. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*.
- Ke Wang, Juntong Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. 2024a. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In *ICLR*. OpenReview.net.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others. 2024b. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *Advances in Neural Information Processing Systems*, 37:95266–95290.
- Yuhao Wang, Ruiyang Ren, Yucheng Wang, Jing Liu, Xin Zhao, Hua Wu, and Haifeng Wang. 2026. Beerg: Balanced entropy engineering for retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 40, pages 33737–33745.
- Thaddäus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane Gu, Nick Matarese, Kevin Swersky, Been Kim, Priyank Jaini, and Robert Geirhos. 2025. Video models are zero-shot learners and reasoners. *arXiv preprint arXiv:2509.20328*.
- Jialong Wu, Tianhao Huang, Changjing He, and Mingsheng Long. 2025. Miniveo3-reasoner: Thinking with videos from open-source priors. <https://github.com/thuml/MiniVeo3-Reasoner>.
- Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, and 1 others. 2024. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.
- Yu-Liang Zhan, Xinyu Tang, Han Wan, Jian Li, Jirong Wen, and Hao Sun. 2026. L2v-cot: Cross-modal transfer of chain-of-thought reasoning via latent intervention. In *AAAI*, pages 12358–12366. AAAI Press.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2):1–124.
- Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. 2024. Open-sora: Democratizing efficient video production for all. *arXiv preprint arXiv:2412.20404*.
- Yiyang Zhou, Haoqin Tu, Zijun Wang, Zeyu Wang, Niklas Muennighoff, Fan Nie, Yejin Choi, James Zou, Chaorui Deng, Shen Yan, and 1 others. 2025. When visualizing is the first step to reasoning: Mira, a benchmark for visual chain-of-thought. *arXiv preprint arXiv:2511.02779*.

## A Benchmark Statistics

As illustrated in Table 6, VIPER comprises a total of 309 samples. We curate most tasks to contain 10-25 samples, ensuring a balanced distribution to assess model performance across diverse scenarios. The Structural category accounts for the largest portion, followed by the Symbolic and Physics domains.

Domain	Task	Count	Total
Physics	Experiment	18	32
	Game	14	
Planning	Navigation	25	44
	Object Manipulation	19	
Spatial	Block Rotate	25	60
	Dice Rolling	20	
	Image Restore	15	
Structural	Chess	20	70
	Maze	20	
	Sudoku	20	
	Tic-Tac-Toe	10	
Symbolic	Knowledge	20	60
	Math	20	
	Multimodal	20	
Temporal	Object Move	25	43
	World Memory	18	
<b>Sum</b>			<b>309</b>

Table 6: Detailed statistics of VIPER.

## B Task Introduction

### B.1 Temporal Reasoning

**Object Moving.** This task requires the model to smoothly move a geometric shape to a certain position. Images are generated programmatically using Python.

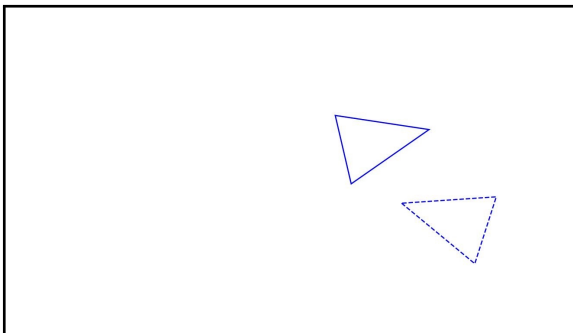


Figure 6: Example image of task Object Moving.

#### Task prompt for Object Moving.

A clean white background featuring a blue triangle defined by a solid outline at top, and a target destination marked by a matching blue triangle with a dotted outline. The solid-outlined object travels in down, ending perfectly overlapping and aligning with the dotted target outline. The object moves by smoothly rotating clockwise while translating. High contrast, vector graphic style, locked camera, no zoom, consistent scale.

**World Memory.** This task requires the model to perform a centered zoom-in and then restore to the original view with visual elements unchanged. Images are collected from open-source web sources.<sup>1</sup>



Figure 7: Example image of task World Memory.

#### Task prompt for World Memory.

An image. Execute a slow, centered zoom-in to 2x magnification, hold the frame static, and then smoothly zoom out to the original field of view. The visual elements remain static and unchanged throughout the entire process. No pan, no dolly, no perspective shifts.

### B.2 Structural Reasoning

**Sudoku.** This task requires the model to complete a partially filled Sudoku grid by correctly filling the empty cells with numbers. Images are generated programmatically using Python from valid Sudoku puzzles.

<sup>1</sup><https://unsplash.com/>

1	8	9	4	7	5	2	6	3
6	2	3	1	8	9	7	4	5
4	7	5	6	2		8	1	9
8	3	1	7	9		5	2	6
2		6	8		1	9	7	4
7	9		2	5	6	3		1
5	4	2	3	6	8	1	9	7
3	6		9	1	7	4	5	
9		7	5	4	2	6	3	

Figure 8: Example image of task Sudoku.

**Task prompt for Sudoku.**

A clean, high-contrast image of a partially filled 9x9 Sudoku grid on a solid white background. Solve the Sudoku, with the missing numbers (1-9) appearing in the correct empty cells. Top-down orthographic view, locked camera, no zoom, no pan, no perspective shifts, consistent scale.

**Tic Tac Toe.** This task requires the model to interpret a midgame Tic-Tac-Toe board and place in an appropriate empty cell according to the game state. Images are generated programmatically using Python from valid Tic-Tac-Toe configurations.

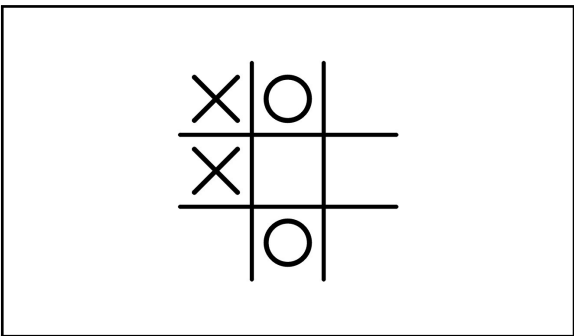


Figure 9: Example image of task Tic Tac Toe.

**Task prompt for Tic Tac Toe.**

A static, top-down view of a minimalist 3x3 Tic-Tac-Toe grid on a clean background. The board is partially filled with existing 'X' and 'O' markers in a midgame state. A single 'X' is drawn into a strategic empty square, executing a logical move to block an opponent or complete a straight line of three. The grid lines and existing marks remain undisturbed. High contrast, vector graphic style, locked camera, no zoom, no dolly, no perspective shifts.

**Checkmate-in-One.** This task requires the model to interpret a chess position from an image and generate a single legal move that results in checkmate. Images are generated programmatically using Python from valid, solvable chess puzzles (Srivastava et al., 2023).

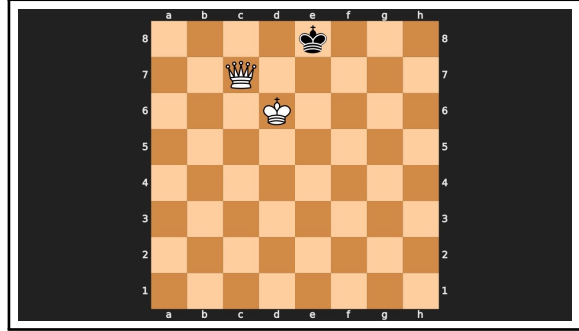


Figure 10: Example image of task Checkmate-in-One.

**Task prompt for Checkmate-in-One.**

A static, top-down 2D view of a digital chess board. A single White piece glides smoothly for exactly one step to deliver checkmate to the Black King. Static camera, no pan, no zoom, no dolly, no perspective shifts.

**Maze.** This task requires the model to draw a continuous path through a 2D maze from a given start point to an end point without crossing any walls. Images are generated programmatically using Python from randomly generated solvable mazes.

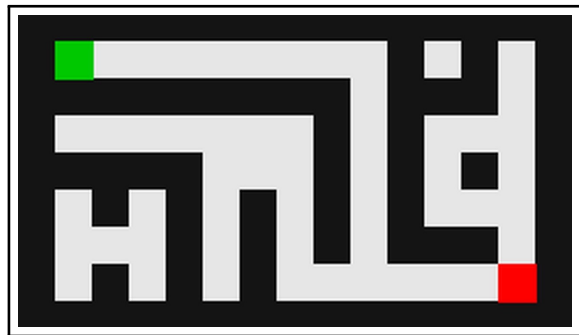


Figure 11: Example image of task Maze.

#### Task prompt for Maze.

A top-down 2D maze. The black lines represent walls, the green square marks the starting point, and the red square marks the ending point. Draw a solid red line animates smoothly from the green square navigating the correct path through the maze structure without intersecting the walls, finally reaching the red square. Static camera, no pan, no zoom, no dolly, no perspective shifts.

### B.3 Symbolic Reasoning

**Multimodal.** This task requires the model to solve a multimodal math problem presented on a whiteboard and generate a step-by-step written derivation leading to the final boxed answer. Images are sourced from existing multimodal understanding benchmarks like MMMU, MathVista and MathVision (Yue et al., 2024; Lu et al., 2023; Wang et al., 2024a).

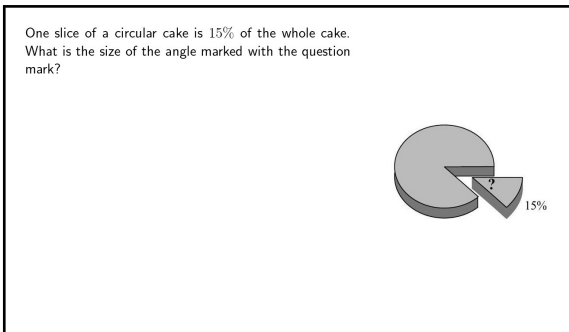


Figure 12: Example image of task Multimodal.

#### Task prompt for Multimodal.

A static shot of a clean, well-lit whiteboard displaying the text or diagram for a multimodal problem: One slice of a circular cake is 15% of the whole cake. What is the size of the angle marked with the question mark?. The solution appears sequentially on the empty space below the problem, writing out the derivation process and concluding with the final answer clearly boxed. Static camera, no pan, no zoom, no perspective shifts, consistent scale.

**Math.** This task requires the model to solve a text-based math word problem shown on a whiteboard and generate a step-by-step written solution ending with a boxed final answer. Images are generated programmatically based on existing text-based math benchmarks like MATH-500

and GSM8K(Cobbe et al., 2021; Hendrycks et al., 2021).

Dennis uses 1 pound of butter for every dozen croissants that he makes. He needs to make 6 dozen croissants. The grocery store currently has a promotion for buy one pound of butter get one half off. If the butter costs \$4.00 a pound, how much will it cost him to purchase 6 pounds of butter?

Figure 13: Example image of task Math.

#### Task prompt for Math.

A static shot of a clean, well-lit whiteboard displaying the text or diagram for a math problem: Dennis uses 1 pound of butter for every dozen croissants that he makes. He needs to make 6 dozen croissants. The grocery store currently has a promotion for buy one pound of butter get one half off. If the butter costs \$4.00 a pound, how much will it cost him to purchase 6 pounds of butter?. The solution appears sequentially on the empty space below the problem, writing out the derivation process and concluding with the final answer clearly boxed. Static camera, no pan, no zoom, no perspective shifts, consistent scale.

**Knowledge.** This task requires the model to answer a text-centric knowledge multiple-choice question shown on a whiteboard and write a step-by-step explanation culminating in a boxed final choice. Images are generated programmatically based on existing knowledge benchmarks like MMLU(Wang et al., 2024b) and GPQA (Rein et al., 2024).

Which of the following (effective) particles is not associated with a spontaneously-broken symmetry?  
A. Photon  
B. Skyrmion  
C. Pion  
D. Magnon.

Figure 14: Example image of task Knowledge.

#### Task prompt for Knowledge.

A static shot of a clean, well-lit whiteboard displaying the text or diagram for a knowledge problem: Which of the following (effective) particles is not associated with a spontaneously-broken symmetry?

- A. Phonon
- B. Skyrmion
- C. Pion
- D. Magnon.. The solution appears sequentially on the empty space below the problem, writing out the derivation process and concluding with the final answer clearly boxed. Static camera, no pan, no zoom, no perspective shifts, consistent scale.

### B.4 Spatial Reasoning

**Block Rotate.** This task requires the model to render a multi-cube block structure while the view-point rotates horizontally by  $180^\circ$  around it, preserving correct spatial relationships between cubes. Images are sourced from VMEvalkit (Deng, 2025).

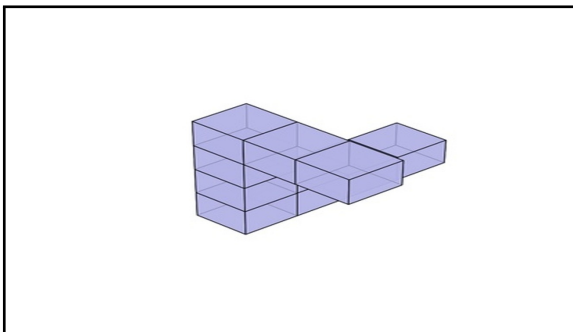


Figure 15: Example image of task Block Rotate.

#### Task prompt for Block Rotate.

A structure composed of 9 cubes viewed from  $316^\circ$  azimuth,  $40^\circ$  elevation. The camera executes a smooth horizontal orbit around the structure, transitioning to a  $136^\circ$  azimuth which is a  $180^\circ$ -degree rotation while maintaining a constant distance from the object. Smooth orbital movement, constant radius, no zoom, no vertical tilt, consistent scale.

**Dice Rolling.** This task requires the model to generate a video of a six-sided die rolling along a specified grid path via successive  $90^\circ$  edge flips, maintaining correct face orientation throughout. Images are rendered using Blender.<sup>2</sup>

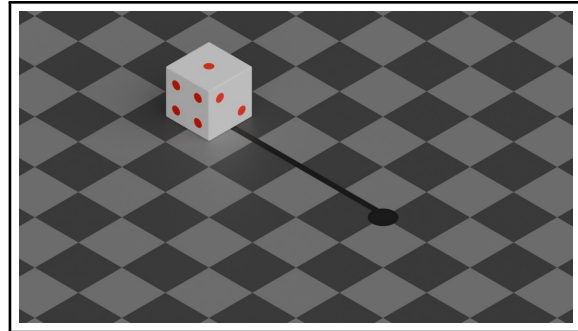


Figure 16: Example image of task Dice Rolling.

#### Task prompt for Dice Rolling.

A standard six-sided die (faces arranged 1 opposite 6, 2 opposite 5, 3 opposite 4) resting on a grid-marked surface. A thick black path follows the grid lines, originating from the die's position and terminating at a solid black circular target. The die travels along the black path by executing precise  $90^\circ$ -degree flips, pivoting on the edge in contact with the table for each step. The motion continues sequentially along the line until the die lands perfectly covering the black circular target. Static camera, high-angle view, sharp focus, no zoom, no pan, no perspective shifts.

**Image Restore.** This task requires the model to restore a shuffled image by translating sub-tiles into their corresponding dashed grid slots to reconstruct the full picture. Images are collected from open-source web sources.



Figure 17: Example image of task Image Restore.

<sup>2</sup><https://www.blender.org/>

#### Task prompt for Image Restore.

A clean white background. A set of shuffled image sub-tiles are arranged at the bottom of the frame, situated below a centered target grid of dashed outlines. The tiles translate smoothly from their starting positions into the matching dashed frames to fully reconstruct the original image. Flat 2D style, locked camera, no zoom, no pan, no perspective shifts.

### B.5 Physics Reasoning

**Experiment.** Experiment task expects the model to render the process of a series of classic physical experiment. The task test the ability of correct rendering of real world physical laws. The input image is manually collected from advanced image generation models.



Figure 18: Example image of task Experiment.

#### Task prompt for Experiment.

A person sits on a chair and keeps spinning. Initially, the arms are rest by his sides, then gradually rise and extend horizontally. Static camera, no pan, no zoom, no dolly.

**Game.** This task requires the model to simulate a mobile game scene where cutting specified ropes causes a suspended candy to move under game physics toward a character. Images are manually collected from existing mobile games.<sup>3</sup>

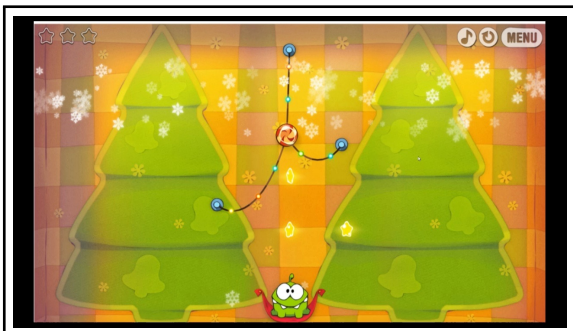


Figure 19: Example image of task Game.

#### Task prompt for Game.

A mobile video game where a candy is fixed by several ropes, each rope connecting one end to the candy and the other end to the wall in the background. A cartoon character is at the bottom of the screen. The candy will be eaten if it moves near the character. Cut the two ropes at the bottom. Static camera, no pan, no zoom, no dolly.

### B.6 Planning Reasoning

**Object Manipulation.** This task requires the model to render an object manipulation sequence where a robot arm grasps a target item and relocates it to a specified position on the table. Images are collected from RoboBench (Luo et al., 2025).



Figure 20: Example image of task Object Manipulation.

#### Task prompt for Object Manipulation.

A lab workbench scene with a wooden table in the center. On the right side of the table there is a gray doll/plush-like object, with other small items nearby. A robot arm is positioned to the right of the table. The robot arm grasps the gray doll and places it down on the left side of the table. Static camera, no pan, no zoom, no dolly.

**Navigation.** This task requires the model to draw a continuous path on a top-down 2D map, navigating from a start location through specified waypoints to the destination. Images are collected from an open-source dataset.

<sup>3</sup>[https://en.wikipedia.org/wiki/Cut\\_the\\_Rope\\_video\\_game](https://en.wikipedia.org/wiki/Cut_the_Rope_video_game)



Figure 21: Example image of task Navigation.

#### Task prompt for Navigation.

A 2D map with a top-down view. Draw a red path on the map to navigate from Air, passing through Corkscrew and finally reaching Oblivion. The red path should be continuously drawn as a smooth line. Static camera, no pan, no zoom, no doll.

## C Rubric Examples

In this section, we demonstrate the system prompt, all domain prompts and an example of task constraints.

### C.1 System Prompt

#### System prompt

You are a video evaluation judge. Your task is to assess a video generated by a model, based on an initial image and a task description. Carefully follow the instructions below to ensure precise evaluation.

#### Inputs:

- **Domain Overview:** A brief introduction of the task's domain and evaluation focus.
- **Task Prompt:** The prompt provided to the video generation model.
- **Initial Image:** The image given to the model as the starting point for video generation.
- **Generated Video:** The video generated by the model based on the initial image and task prompt.
- **Process Constraints:** The set of process-level constraints that the generated video

content should adhere to.

- **References:** Ground-truth information to assist in evaluation, if available. Could be videos, images, or textual descriptions.

#### Evaluation Criteria:

Your evaluation should be based on the following two aspects:

1. **Process Consistency:** Verify if **all frames** in the generated video adhere to the **process constraints**. Refer to the **References** if necessary.
2. **Goal Consistency:** Verify if **at least one frame** in the generated video fulfills the **task prompt**. Refer to the **References** if necessary.

The video is considered **correct** only if it fulfills both process and goal consistencies. If either is not fulfilled, the video is **incorrect**.

#### Output Format:

- Begin with your analysis of both process and goal consistencies inside <think> tags.
- After completing your analysis, provide your final evaluation inside <answer> tags.

Example:

```
<think>
YOUR ANALYSIS HERE
</think>
<answer>
{
  "process_consistency":
  "goal_consistency":
  "decision":
}
</answer>
```

Please ensure your response follows this structure precisely.

## C.2 Domain Introduction

### Domain: Temporal Reasoning

This task belongs to the Temporal Reasoning domain. The evaluation should focus on the model’s ability to accurately track and represent the dynamic evolution of objects, scenes, and actions over time.

### Domain: Structural Reasoning

This task belongs to the Structural Reasoning domain. The evaluation should focus on the model’s ability to process and generate content based on strict, predefined rules, constraints, and structured environments (e.g., games, puzzles, geometric layouts).

### Domain: Symbolic Reasoning

This video belongs to the Symbolic Reasoning domain. The evaluation should focus on the model’s ability to interpret, manipulate, and generate abstract symbols (numbers, text, logical operators) to solve problems.

### Domain: Spatial Reasoning

This task belongs to the Spatial Reasoning domain. The evaluation should focus on the model’s understanding of 3D space, geometric structures, perspective, and the relative positions of objects.

### Domain: Physics Reasoning

This task belongs to the Physical Reasoning domain. The evaluation should focus on the model’s ability to simulate accurate physical laws, (e.g., gravity, collision dynamics, material properties).

### Domain: Planning Reasoning

This task belongs to the Planning Reasoning domain. The evaluation should focus on the model’s ability to execute complex, multi-step tasks that require foresight, strategy, and logical sequencing.

## D Related Work

### D.1 Video Generation Models

Recent advancements in diffusion-based Transformers and architecture optimization have significantly accelerated progress in video generation (Ho et al., 2022; Yang et al., 2024; Zheng et al., 2024; Jin et al., 2024; Blattmann et al., 2023). Modern video generation systems can now generate high-fidelity videos conditioned on both initial images and textual instructions. On the proprietary side, leading state-of-the-art models include Sora 2 (OpenAI, 2025b), Veo 3.1 (Google, 2025), Seedance 1.5 Pro (Chen et al., 2025b), and Wan 2.6 (Wan Team, 2025), all of which exhibit improved temporal consistency and enhanced ability to simulate complex physical interactions over extended time horizons. Concurrently, the open-source community has made significant strides, with models like Wan 2.2 (Wan Team, 2025) and HunyuanVideo (Hunyuan Team, 2025) achieving competitive performance.

### D.2 Generative Video Reasoning

Reasoning has been long considered an area dominated by language (Zhao et al., 2023; Tang et al., 2026; Zhan et al., 2026; Wang et al., 2026). As video generation models evolve, their potential for completing reasoning tasks has gradually been uncovered. Research has shown that video generation models can generalize to many common visual tasks that they have not encountered during training, using a sequence of frames to represent the solving process, which is termed Chain-of-Frames (CoF) reasoning (Wiedemer et al., 2025). Subsequent work has focused on evaluating and optimizing video models for generative video reasoning tasks. Some studies have designed empirical studies to showcase the reasoning capabilities of video models (Guo et al., 2025b; Tong et al., 2025). Additionally, many efforts have proposed benchmarks covering various evaluation dimensions by collecting existing video data or manually constructing examples (Chen et al., 2025a; Deng, 2025; Zhou et al., 2025; Liu et al., 2025). Furthermore, some studies have made initial explorations into improving models’ reasoning abilities through prompt optimization and supervised fine-tuning using correct video examples, yielding some promising results (Chen et al., 2025a; Wu et al., 2025).