

Planning with Sketch-Guided Verification for Physics-Aware Video Generation

Yidong Huang¹ Zun Wang¹ Han Lin¹ Dong-Ki Kim² Shayegan Omidshafiei²
 Jaehong Yoon³ Yue Zhang¹ Mohit Bansal¹

¹UNC Chapel Hill ² FieldAI ³ Nanyang Technological University

<https://sketchverify.github.io/>

Abstract

Recent video generation approaches increasingly rely on planning intermediate control signals such as object trajectories to improve temporal coherence and motion fidelity. However, these methods mostly employ single-shot plans that are typically limited to simple motions, or iterative refinement which requires multiple calls to the video generator, incurring high computational cost. To overcome these limitations, we propose **SketchVerify**, a training-free, sketch-verification-based planning framework that improves motion planning quality with more dynamically coherent trajectories (i.e., physically plausible and instruction-consistent motions) prior to full video generation by introducing a test-time sampling and verification loop. Given a prompt and a reference image, our method predicts multiple candidate motion plans and ranks them using a vision-language verifier that jointly evaluates semantic alignment with the instruction and physical plausibility. To efficiently score candidate motion plans, we render each trajectory as a lightweight video sketch by compositing objects over a static background, which bypasses the need for expensive, repeated diffusion-based synthesis while achieving comparable performance. We iteratively refine the motion plan until a satisfactory one is identified, which is then passed to the trajectory-conditioned generator for final synthesis. Experiments on *WorldModelBench* and *PhyWorldBench* demonstrate that our method significantly improves motion quality, physical realism, and long-term consistency compared to competitive baselines while being substantially more efficient. Our ablation study further shows that scaling up the number of trajectory candidates consistently enhances overall performance.

1. Introduction

Image-to-Video (I2V) generation has demonstrated strong potential across a wide range of applications, including

robotic manipulation, autonomous driving, and game content creation. While modern video generative models [4, 12, 16, 34, 37] have enabled impressive visual quality and semantic alignment, producing videos with physically realistic and temporally consistent motion remains challenging. In particular, these models often fail to interpret fine-grained motion instructions and struggle to generate sequences that adhere to plausible physical dynamics [8, 22]. Recent studies have introduced intermediate object-level layout [2, 40] or trajectory planning [9, 15, 23] with large language models (LLMs) to guide video generation, aiming to improve motion fidelity and controllability. However, most existing approaches adopt a single-shot planning paradigm (Fig. 1a), where a single control sequence is generated per prompt. While this design is straightforward, it is vulnerable to inaccuracies or noise in the predicted plan, which can propagate through the generation process and result in inconsistent or implausible object motions. To mitigate the instability inherent in single-shot planning, an alternative line of research explores iterative refinement (Fig. 1b), where the prompt or control signals are progressively updated over multiple steps [7, 20, 45]. Although such methods can enhance visual realism through feedback-based correction, they incur substantial computational overhead due to repeated diffusion calls during generation.

To address these limitations, we propose **SketchVerify**, a test-time planning framework that iteratively refines motion plans using verification on lightweight video sketches instead of costly full video synthesis (Fig. 1c). SketchVerify integrates a multimodal verifier with a test-time search procedure to automatically detect and correct semantic or physical inconsistencies, compensating for the lack of self-correction in one-shot planning. By decoupling refinement from the diffusion backbone and verifying motion at the sketch or layout level, SketchVerify avoids the heavy overhead of full-generation-based iterative updates, enabling efficient test-time search in about five minutes—over an order of magnitude faster than baselines requiring full generation.

Specifically, given a text prompt and an initial image,

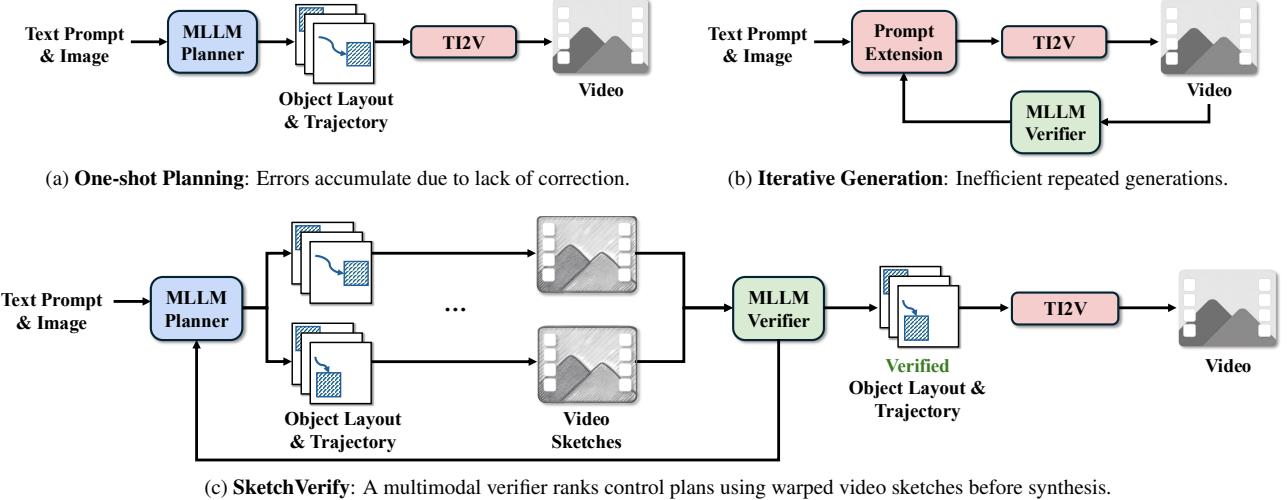


Figure 1. Comparison of SketchVerify with other MLLM planning based video generation pipelines. Existing methods either rely on one-shot planning, which lacks correction, or iterative refinement, which requires repeated generation. Our method addresses both issues by selecting high-quality control plans using a multimodal verifier prior to synthesis.

our approach first constructs a high-level motion plan composed of sequential sub-instructions (e.g., “approach the ball,” “pick it up,” “place it on the table”) and identifies the corresponding movable objects through segmentation. Then, it sequentially generates a trajectory plan for each sub-instruction over time. In particular, given a sub-instruction and the context image derived from the previous step, SketchVerify samples multiple candidate trajectory plans represented as a sequence of bounding boxes capturing the object’s location at each frame. To efficiently visualize and verify these motion candidates, we render each trajectory as a lightweight **video sketch**. Rather than synthesizing full videos, the framework crops the segmented object from the first frame and composites it onto a static background. This lightweight video sketch preserves the essential spatial and temporal structure of the scene, allowing significantly faster verification while maintaining comparable content information and verification quality to full video generation. A vision–language verifier then evaluates each sketch along two complementary dimensions. First, it assesses semantic alignment by comparing the sketch with the corresponding sub-instruction, ensuring that the depicted motion fulfills the described intent (e.g., whether the object indeed “moves toward the basket” or “picks up the ball”). Second, it evaluates physical plausibility through structured reasoning over several motion principles, including Newtonian consistency, non-penetration with scene elements, gravity-coherent vertical motion, and shape stability across frames. The trajectory achieving the highest aggregated verification score is selected as the final motion plan for that sub-instruction. After all sub-instructions are processed, their verified trajectories are merged into a uni-

fied plan, which is passed to a trajectory-conditioned diffusion model for final video synthesis. By conducting this structured planning and verification entirely at test time, our approach produces semantically coherent and physically grounded videos without requiring additional training or costly iterative refinement.

We evaluate SketchVerify on WorldModelBench and PhyWorldBench, two large-scale benchmarks designed to assess instruction compliance, physical reasoning, and temporal coherence in generative video models. Our method consistently outperforms state-of-the-art open-source I2V models across instruction following, physical law adherence, and commonsense consistency, while reducing overall planning cost by nearly an order of magnitude compared to iterative refinement pipelines. Ablation studies further show that (i) multimodal verification markedly strengthens spatial and physical reasoning compared to language-only variants, (ii) scaling the verifier improves trajectory plausibility, (iii) sketch-based verification matches the quality of full video-based verification with nearly a tenfold efficiency gain, and (iv) increasing the number of sampled trajectories yields steady performance gains.

2. Related Works

MLLM Planning for Video Generation. Recent work increasingly leverages LLMs and MLLMs to provide structured planning for video generation. GPT-style models expand sparse text prompts into “video plans”—including bounding-box trajectories [23, 25, 43], scene-level keyframes [15], or motion-aware sketches [23]—which are then used for layout-guided diffusion synthesis [11, 24, 25, 49–53]. However, these methods rely

on a single-pass plan that often remains coarse or physically inconsistent. In contrast, our approach performs *iterative*, verifier-guided refinement, repeatedly scoring and updating candidate trajectories to produce plans with stronger spatial constraints and physical plausibility.

Iterative Refinement for Visual Generation. Iterative refinement is widely used to improve consistency and controllability in visual generation. Methods such as RPG [46], PhyT2V [45], VideoRepair [20], and VISTA [29] refine prompts or layouts by repeatedly evaluating fully generated videos, which is time-consuming (whole process usually needs over 30 minutes). In contrast, we refine during the planning stage by scoring lightweight sketch-based simulations with a multimodal verifier, avoiding repeated video synthesis and enabling efficient test-time optimization without harming the verification performance.

Physics-Aware Video Generation. Recent work reveals that state-of-the-art video diffusion models often violate even basic physical laws [3, 31]. To address this, prior studies explored a wide range of physics-aware enhancements, including physics-simulator-guided motion [28, 30, 35, 42], physics-driven post-training and reward optimization [21, 26, 41], and vision-language reasoning or force-based conditioning that inject implicit physical priors [6, 10, 39, 47]. Instead of relying on heavy simulation, specialized datasets, or extensive finetuning, we focus on plan-level iterative refinement with our SketchVerify framework, achieving zero-shot generalization across diverse physical scenarios.

3. Method – SketchVerify

Our approach performs lightweight, verifier-guided refinement entirely before generation, scoring and improving motion plans at test time without any additional training or repeated video synthesis. We organize the method section as follows. Section 3.1 introduces the problem formulation and provides a high-level overview of our framework. Section 3.2 describes the high-level plan decomposition and object/background extraction process from the input prompt and image. Section 3.3 presents our proposed visual test-time planning module (SketchVerify), which performs trajectory sampling and multimodal verification based on sketch-based surrogates. Section 3.4 describes the final video synthesis process using a trajectory-conditioned diffusion model guided by the selected motion plan.

3.1. Overview

Given a natural language prompt \mathcal{P} and an initial image \mathbf{I}_0 , our method produces a temporally coherent and physically plausible video $\mathbf{V} = \{\mathbf{I}_1, \dots, \mathbf{I}_T\}$. As shown in Fig. 2, the framework is composed of three key modules:

1. **High Level Planning and Object Parsing:** This module interprets the prompt’s high-level narrative intent and

generates a sequence of actionable sub-goals. Concurrently, it parses the initial scene to isolate the dynamic target object from the static “stage” (the background), thereby defining a clear, structured problem for the subsequent motion planning module.

2. **Test-Time Planning:** This module constitutes the core contribution of our method. Instead of performing expensive trial-and-error with diffusion models, SketchVerify conducts an efficient test-time search for optimal motion trajectories. It samples lightweight motion candidates (video sketches) and scores them with a multimodal verifier that assesses **semantic alignment** with the instruction and **physical plausibility** based on real-world motion priors. By verifying motion quality before synthesis, SketchVerify decouples reasoning about object dynamics from the computationally intensive generation process.
3. **Trajectory-Conditioned Video Generation:** The verified motion plan is passed to a diffusion-based video generator. Because the generator receives a pre-verified, high-quality motion plan, this stage focuses solely on visual fidelity, producing semantically coherent and physically consistent video sequences.

3.2. High-Level Planning and Object Parsing

High-Level Planning. We begin by generating a structured plan of sub-instructions $\mathcal{P}_1, \dots, \mathcal{P}_M$ (e.g., “approach the ball”, “pick up the ball”) from the natural language prompt \mathcal{P} using an MLLM, thereby mitigating the difficulty of long-horizon planning.

Object and Background Extraction. To enable motion planning over a clean static canvas, we first identify objects involved in motion. Specifically, given the prompt \mathcal{P} and initial frame \mathbf{I}_0 , we use an MLLM to extract a list of object names expected to move according to the described actions. Next, we apply a detector-segmenter pair, GroundedSAM [17, 27, 33], for precise mask extraction to localize the mentioned objects. This results in a set of object masks $\mathcal{M} = \{m_1, \dots, m_N\}$ corresponding to the moving entities, where N is the number of moving objects proposed. To obtain a clean, static background for compositing, we remove the masked object regions from \mathbf{I}_0 and fill them using Omnieraser [44], a background inpainting model fine-tuned from FLUX [19]. The output is a static background image \mathbf{B} , which serves as the canvas for video sketch rendering in subsequent stages (Step 1 in Fig. 2).

3.3. Test-Time Planning

To improve trajectory quality without incurring the computational cost of iterative synthesis, we propose a sketch-verification guided test-time planning module that samples and verifies object-level motion plans before video generation (Step 2 in Fig. 2).

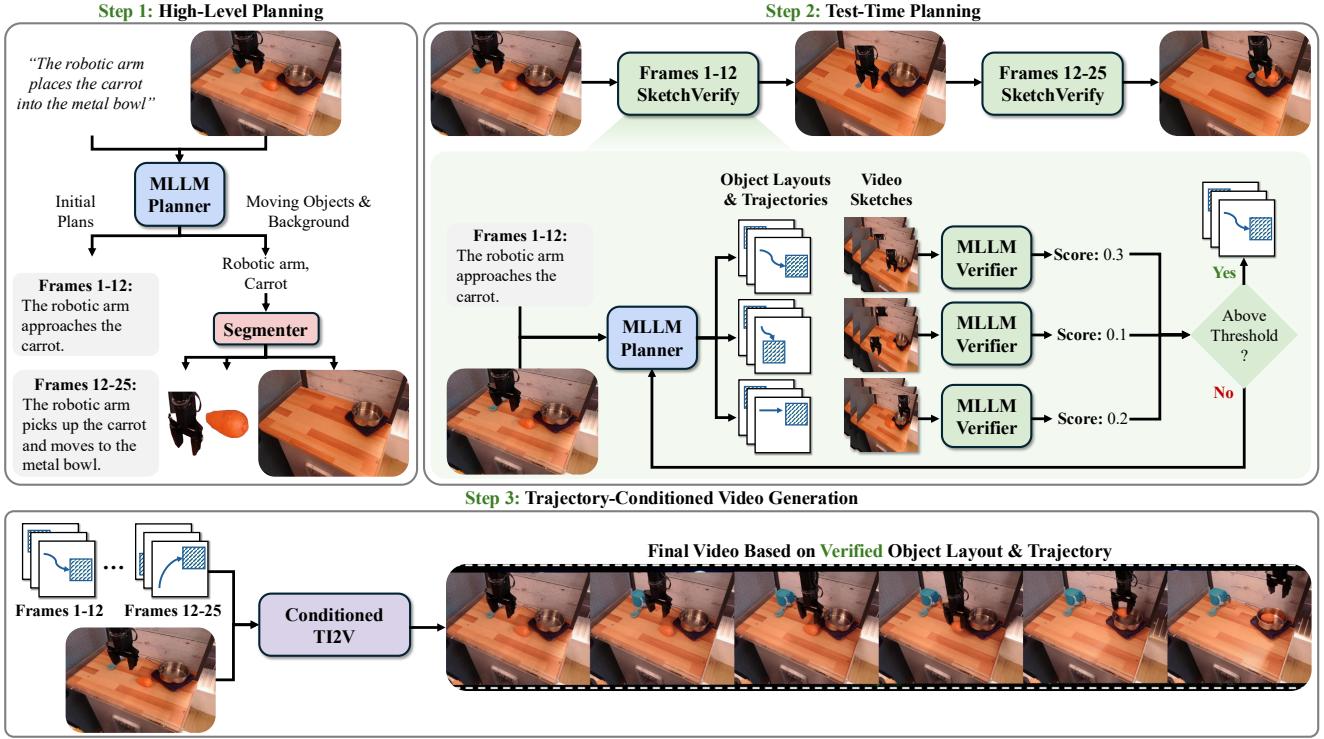


Figure 2. **Overview of our framework.** Given a prompt and initial frame, we (1) decompose instructions and segment movable objects, (2) sample and verify candidate trajectories using lightweight video sketches scored by a multimodal verifier and (3) synthesize the final video using a trajectory-conditioned diffusion model. We provide more detail about the MLLM verifier in Fig. 3.

Trajectory Sampling. For each sub-instruction \mathcal{P}_i , the goal is to generate a set of candidate trajectories that guide the moving object \mathcal{O} according to the intended action. The sampling process is conditioned on the current visual context \mathbf{C}_i , which provides spatial grounding for planning. We use MLLM Planner \mathcal{F} to generate K candidate trajectories:

$$\left\{ \Pi_i^{(1)}, \dots, \Pi_i^{(K)} \right\} = \mathcal{F}(\mathcal{P}_i, \mathcal{O}, \mathbf{C}_i),$$

where each $\Pi_i^{(k)}$ is a list of bounding boxes over T_i frames:

$$\Pi_i^{(k)} = \{ \mathbf{b}_i^{(k,t)} = (x_{\min}^{(t)}, y_{\min}^{(t)}, x_{\max}^{(t)}, y_{\max}^{(t)}) \}_{t=1}^{T_i}.$$

Here, $\mathbf{b}_i^{(k,t)}$ denotes the bounding box of the object at frame t , with $(x_{\min}^{(t)}, y_{\min}^{(t)})$ and $(x_{\max}^{(t)}, y_{\max}^{(t)})$ being the upper-left and lower-right coordinates, respectively. For example, if \mathcal{P}_i is “move the apple toward the basket,” then $\Pi_i^{(k)}$ can define a smooth horizontal motion of the apple across T_i frames toward the location of the basket. The visual context \mathbf{C}_i is initialized as the reference image \mathbf{I}_0 when $i = 1$, and updated at each subsequent step to the last frame of the selected sketch \mathbf{S}_{i-1}^* , preserving temporal continuity throughout the planning process.

Video Sketch Rendering. To enable assessing plans without incurring the cost of full video generation, we render a

lightweight video sketch that visualizes only the intended object motions. Specifically, each trajectory $\Pi_i^{(k)}$ is converted into a sketch $\mathbf{S}_i^{(k)}$ by cropping the segmented object region from the initial frame \mathbf{I}_0 using its predicted mask and compositing it frame by frame onto the static background \mathbf{B} according to the bounding box coordinates in $\Pi_i^{(k)}$. These sketches provide a faithful, layout-preserving approximation of the planned motion, enabling efficient test-time verification focused on spatial-temporal coherence rather than appearance-level generation (see Sec. 5.2).

Verifier-Guided Scoring. As is shown in Fig. 3, we perform a two-stage evaluation strategy combining semantic alignment and physics-aware verification to assess the quality of each trajectory candidate. First, we compute a semantic score using an MLLM. Given the sub-instruction \mathcal{P}_i and corresponding sketch $\mathbf{S}_i^{(k)}$, the MLLM \mathcal{V}_{sem} returns a scalar compatibility score: $s_k^{\text{sem}} = \mathcal{V}_{\text{sem}}(\mathbf{S}_i^{(k)}, \mathcal{P}_i)$, which reflects how well the proposed motion aligns with the intended behavior. In parallel, we evaluate the physical plausibility of each sketch using structured prompts and few-shot in-context learning to probe four physical laws:

- **Newtonian Consistency:** Acceleration and deceleration should reflect plausible physical dynamics.
- **Penetration Violation:** Moving objects should not pass

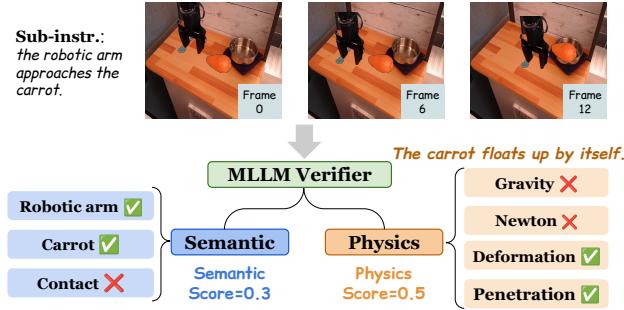


Figure 3. Illustration of MLLM verifier. Given a video sketch and sub-instruction, the MLLM outputs semantic and physics scores used to rank candidate trajectories.

through static scene elements.

- **Gravitational Coherence:** Vertical motion should follow realistic arcs consistent with gravity.
- **Deformation Consistency:** Object size and shape should remain stable throughout the sequence.

Each response from the MLLM is parsed into a scalar score $s_k^{(l)}$, where $l \in \mathcal{L} = \{\text{Newton}, \text{Penetration}, \text{Gravity}, \text{Deformation}\}$. We map descriptive outputs (e.g., “very consistent”) to numerical values using predefined rules (e.g., “very consistent” → 1.0, “somewhat inconsistent” → 0.7). Detailed prompt templates and mapping rules are provided in Appendix A.4. The final trajectory is selected by maximizing a weighted combination of semantic alignment and physical plausibility scores:

$$\Pi_i^* = \arg \max_k \left(\lambda_{\text{sem}} s_k^{\text{sem}} + \sum_{l \in \mathcal{L}} \lambda_l s_k^{(l)} \right),$$

where \mathcal{L} denotes the set of physical law dimensions and the λ coefficients balance their relative importance.

Iterative Trajectory Selection. To ensure trajectory quality, we adopt an iterative refinement strategy. If all candidate scores fall below a quality threshold τ , the entire set is discarded. A new batch of trajectories is then sampled by prompting the planner with an augmented instruction that incorporates feedback on previous failure cases until a valid plan is found or a retry limit is reached. Finally, we set the last frame of \mathbf{S}_i^* as the new context frame \mathbf{C}_{i+1} for a new round of SketchVerify on the next sub-instruction.

3.4. Trajectory-Conditioned Video Generation

Once the full instruction plan has been verified and selected, we obtain a verified motion plan $\Pi^* = \{\Pi_1^*, \dots, \Pi_M^*\}$, where each sub-plan Π_i^* is a sequence of bounding boxes $\{b_{i,1}, \dots, b_{i,T_i}\}$ over T_i frames. For each box $b_{i,t} \in \mathbb{R}^4$, we extract a representative point $p_{i,t} \in \mathbb{R}^2$ (e.g., the center), resulting in a sparse trajectory $\mathcal{P}_i = \{p_{i,1}, \dots, p_{i,T_i}\}$.

The full object path is formed by concatenating all sub-trajectories: $\mathcal{P}^* = \mathcal{P}_1 \parallel \dots \parallel \mathcal{P}_M$. We temporally interpolate this sequence to produce a dense trajectory $\bar{\mathcal{P}} = \{q_1, \dots, q_T\}$ over T frames, where each $q_t \in \mathbb{R}^2$ specifies the object position at time t .

We adopt a pre-trained trajectory-conditioned image-to-video diffusion model for video generation. It encodes the initial image \mathbf{I}_0 into latent features and modulates the denoising process by injecting object trajectory latents. This injection follows the planned trajectory Π^* and guides the generation process to produce coherent motion consistent with both appearance and spatial control. The result is a T -frame video \mathbf{V} that exhibits faithful motion behavior, aligned with the high-level prompt and semantically and physically verified trajectory (Step 3 in Fig. 2).

4. Experimental Results

4.1. Benchmarks and Evaluation Metrics

Benchmarks. We evaluate our method on two recent large-scale benchmarks for visual world models:

- WorldModelBench [22], which evaluates instruction following, physical plausibility, and commonsense across 7 domains using a benchmark-provided MLLM scorer.
- PhyWorldBench [8], which tests fine-grained physical realism over 350 prompts. Since it is a text-to-video benchmark, we generate the first frame using FLUX [19] and then perform I2V.

Evaluation Metrics. Across these two benchmarks, we evaluate instruction following, physical law coherency (Newtonian motion, deformation, fluid, penetration, gravity), commonsense consistency (frame and temporal), and efficiency (planning time and generation cost), using the benchmark-provided MLLM scorer for all assessments; full metric definitions are included in the Appendix A.2.2.

4.2. Implementation Details

We use the multimodal version of GPT-4.1 [32] as the default planner to generate five candidate trajectories of all moving objects, where each trajectory records the coordinates of the top-left and bottom-right corners of the bounding box at every frame. We construct lightweight video sketches by pasting foreground objects (segmented with GroundedSAM [33]) onto static backgrounds to render object motion direction. These sketches are scored in the range [0, 1] by Gemini 2.5 [5]. This vision-language verifier uses prompt-based queries to assess semantic alignment and physical laws, including Newtonian consistency, object penetration, gravitational coherence, and deformation consistency. We independently assess each law using in-context prompts with positive and negative trajectories and aggregate the resulting plan scores for the final ranking. We use ATI-14B model [38] to generate 81 frame 480p videos. We

Table 1. Comparison on WORLDMODELBENCH [22]. We report scores for instruction following, physical law coherence (grouped under “Physics”), and commonsense consistency (“Frame” and “Temporal”), along with an overall sum score aggregating all metrics. Best results in each column are highlighted in bold.

Model	Instruction		Physics					Commonsense		Sum↑	Plan Time (min)↓
	Follow↑		Newton↑	Deform↑	Fluid↑	Penetr.↑	Gravity↑	Frame↑	Temporal↑		
Open-Source Video Models											
Hunyuan-Video [18]	1.18		1.00	0.80	1.00	0.92	1.00	0.64	0.70	7.24	–
CogVideoX [48]	1.46		0.99	0.70	0.99	0.77	0.96	0.86	0.94	7.67	–
Wan-2.1 [36]	1.88		1.00	0.76	0.99	0.81	0.99	0.96	0.82	8.21	–
Cosmos [1]	2.06		1.00	0.84	0.99	0.92	1.00	0.92	0.90	8.63	–
Open-Sora [13]	1.64		0.98	0.82	1.00	0.91	1.00	0.80	0.84	7.99	–
STEP-Video [14]	1.04		1.00	0.75	1.00	0.89	1.00	0.50	0.71	6.89	–
Single-Shot Planning											
VideoMSG [23]	1.46		0.99	0.79	0.99	0.83	0.94	0.96	0.82	7.78	1.33
Iterative Planning											
PhyT2V [45]	1.97		1.00	0.82	0.96	0.82	1.00	0.81	0.81	8.19	61.86
SketchVerify (Ours)	2.08		1.00	0.89	1.00	0.92	1.00	0.96	0.86	8.71	4.71

show all the implementation detail in Appendix A.1.1.

5. Quantitative Results

Evaluation on WorldModelBench. Table 1 compares our method with strong open-source video generation models, including CogVideoX [48], Cosmos [1], Hunyuan-Video [18], Open-Sora [13], Wan-2.1 [36], and STEP-Video [14], as well as planning-based baselines such as VideoMSG [23] (single-shot) and PhyT2V [45] (multi-step refinement). Our approach achieves the strongest performance across all major evaluation dimensions, including instruction following (2.08), physical law coherence (gravity, penetration, deformation), and overall commonsense consistency. Compared to the base model Wan-2.1 [36], our method improves instruction-following accuracy by 10.6% and increases overall physics coherence by 6%, including a 17% reduction in deformation-related violations. While multi-step pipelines such as PhyT2V provide gains over one-shot planners, they rely on repeated, computationally expensive synthesis cycles (typically requiring about 12.5 minutes for planning and 70 minutes for full video generation). In contrast, our verifier-guided sampling performs high-quality trajectory selection within a single planning stage, reducing generation time to just 4.7 minutes, corresponding to a 93% speed-up. We present a detailed per-component runtime analysis in Appendix A.1.3.

Evaluation on PhyWorldBench. As shown in Table 2, our verifier-guided framework achieves the highest overall score (19.84) and the strongest performance on the physical standard category (23.52), demonstrating superior physical consistency and object stability. While Cosmos [1] achieves a slightly higher object-event score (48.29 vs. 43.11), its overall and physical-standard scores are substantially lower, suggesting weaker temporal physical consistency. Relative to the base model Wan-2.1 [36], our method boosts object-event accuracy by 22% and physical accuracy by 18%.

Table 2. Evaluation on PHYWORLDBENCH. We report category-wise pass rates for object and event (Obj+Evt), physical standard (Phys. Std), and the combined overall score (All). Best results are shown in **bold**, and second-best results are underlined.

Model	Obj+Evt↑	Phys. Std↑	All↑
CogVideoX [48]	41.62	<u>21.68</u>	17.34
Wan-2.1 [36]	35.34	19.83	15.52
OpenSora [13]	36.86	17.43	14.00
Cosmos [1]	48.29	15.71	14.00
Hunyuan-Video [18]	24.86	14.16	10.12
STEP-Video [14]	29.51	16.33	12.89
SketchVerify (Ours)	<u>43.11</u>	23.52	19.84

These results highlight that our test-time verification not only preserves object-level realism but also improves causal and physical coherence across diverse scenarios.

5.1. Qualitative Results

Fig. 4 presents qualitative comparisons across four domains from WorldModelBench: Human, Natural, Video Game, and Robotics. Existing baselines such as CogVideoX [48], Cosmos [1], and Wan-2.1 [36] frequently exhibit visible artifacts. For example, in the Human domain, baseline models often produce body parts that stretch unnaturally or remain suspended mid-air during jumping motions, whereas our verifier-guided approach generates smooth forward jumps with realistic limb coordination and consistent gravity response. In the Natural domain, competing models fail to follow the instruction, where the snow never seem to move, while our method maintains a continuous downhill flow that adheres to slope geometry. In the Video Game scenes, baselines tend to misalign collisions (e.g., football players phasing through each other or even merging into one), while our results preserve accurate object contact. Finally, in the Robotics domain, previous models often cause gripper-object misalignment or floating artifacts during manip-

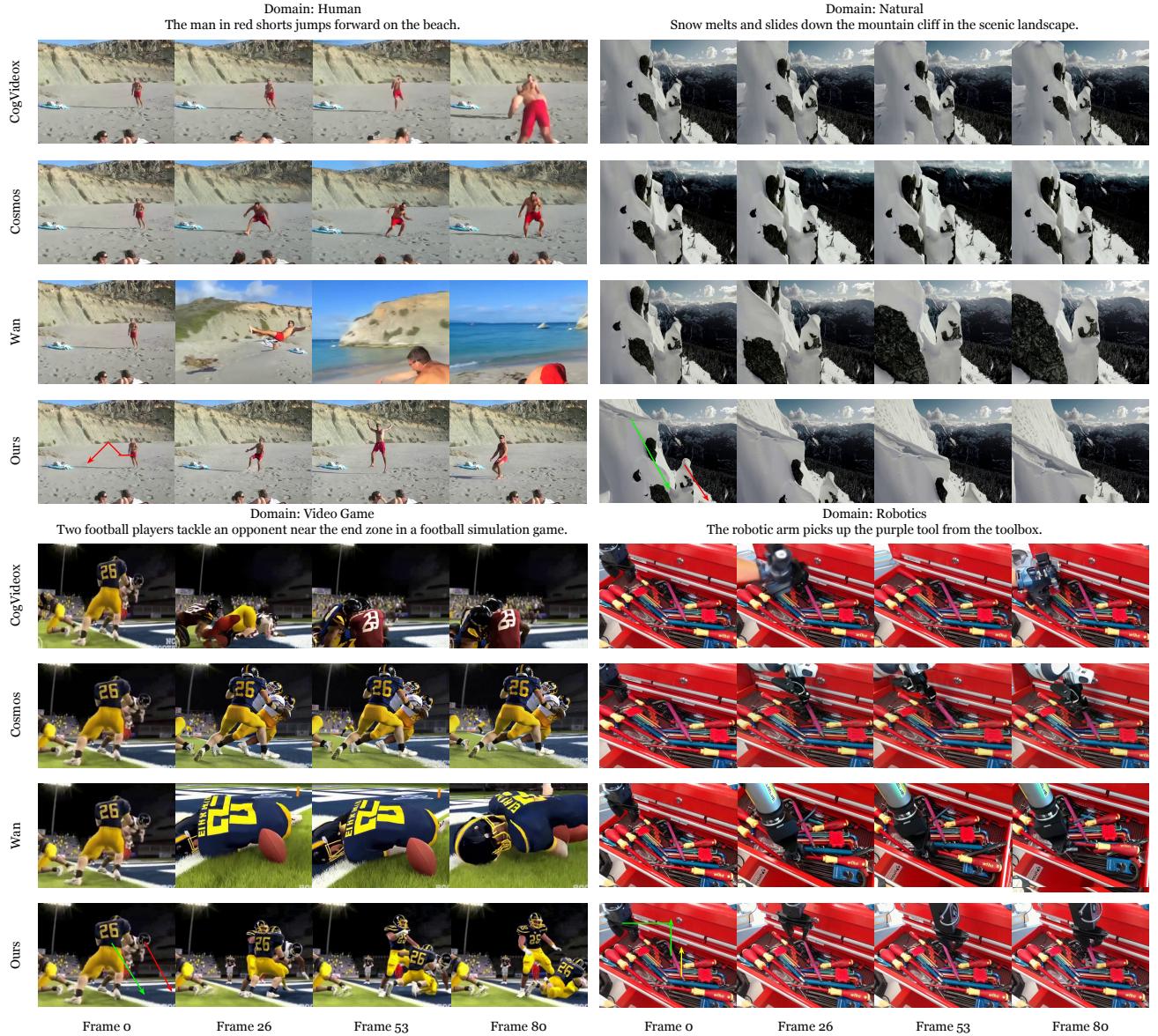


Figure 4. Qualitative comparison on four representative domains from WORLDMODELBENCH: Human, Natural, Video Game, and Robotics. Each group shows sampled frames from competing models given the same text prompt. Frames are uniformly sampled from each generated 81-frame video.

ulation, whereas our planner enables stable grasping and lifting trajectories. Together, these examples demonstrate that verifier-guided planning effectively reduces physical implausibilities and enhances temporal coherence across diverse environments. We show more qualitative results in Appendix A.3.

5.2. Ablation Study

We conduct ablation studies on WorldModelBench to examine how verifier type and test-time sampling affect motion planning quality.

Verifier Modality. We evaluate the impact of verifier modality in Table 3. Relying solely on language-based planning, such as VideoMSG [23], leads to suboptimal motion control, as the generated trajectories often lack spatial coherence and violate basic physical constraints. Adding a language-only verifier slightly improves overall scores by providing textual feedback, yet it remains limited in capturing fine-grained motion cues, as language models struggles to directly perceive object geometry, depth, or trajectory smoothness without visual context. In contrast, our multimodal verifier can directly visually assess motion consis-

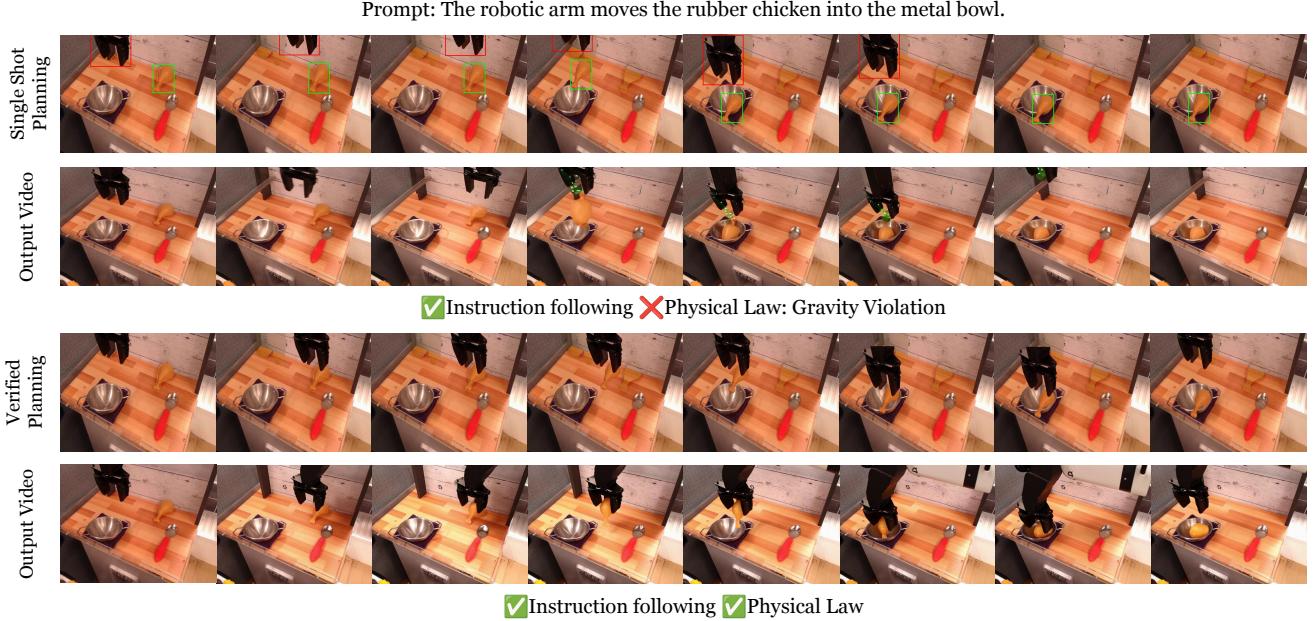


Figure 5. Ablation study on verifier modality. Introducing visual input to the verifier significantly improves both instruction following and physical plausibility, highlighting the importance of multimodal grounding for reliable trajectory evaluation.

Table 3. Ablation study on WORLDMODELBENCH showing the effect of verifier guidance and test-time sampling. “Lang-Only” uses a language-only verifier over trajectory descriptions, while “Ours” adds sampling-based trajectory selection.

Variant	Instr. Follow↑	Phys. Score↑
Single-Shot (no verifier)	1.46	4.55
Lang-Only Verifier	1.49	4.76
Ours (MLLM verifier)	2.08	4.81

tency and interactions, thereby offering stronger and more physically grounded guidance during test-time planning. As illustrated in Fig. 5, the multimodal verifier effectively identifies and rejects physically implausible trajectories (e.g., gravity violations), leading to more realistic and physically consistent motion generation.

Effect of Different Verifier Choices. We compare different MLLMs as verifiers in the plan ranking loop in Table 4. Using a smaller model such as Qwen-VL-3B provides limited improvements due to weaker spatial reasoning. In contrast, a stronger model like Gemini-2.5 yields more accurate motion selection and higher overall quality, highlighting the clear benefit of scaling the verifier’s reasoning capacity.

Effect of Different Planner Choices. We compare different MLLMs as verifiers in the plan-ranking loop (Table 5). Smaller models such as Qwen-VL-3B yield limited gains due to weaker spatial reasoning and poorer grounding of object dynamics. In contrast, stronger verifiers like

Table 4. Ablation of different verifier scale. Using stronger and larger multimodal verifiers improves semantic and physical reasoning during motion plan selection. All experiments share the same underlying inference pipeline.

Verifier	Instr. Follow↑	Phys. Score↑
Qwen2.5-VL-3B	1.62	4.68
Qwen2.5-VL-32B	1.83	4.72
Gemini (Default)	2.08	4.81

Table 5. Ablation of different planner choices. Using stronger MLLM planner improves semantic and physical reasoning during motion plan selection. All methods use the same pipeline.

Planner	Instr. Follow↑	Phys. Score↑
Qwen2.5-VL-3B	1.23	4.50
Qwen2.5-VL-72B	1.59	4.57
GPT-4.1(Default)	2.08	4.81

GPT-4.1 provide more reliable assessments of motion quality, enabling more accurate trajectory selection. Moreover, weaker open-source VLMs exhibit insufficient instruction-following ability, which further limits their effectiveness on multi-step, numerically precise planning tasks.

Effect of Sampling Budget K . Fig. 6 examines how the number of sampled candidate trajectories K influences the final motion quality. When there is no iterative genera-

Table 6. Ablation study comparing different verification strategies on WORLDMODELBENCH. The reported *plan time* measures only the duration of motion planning and verification before the final diffusion-based video generation.

Verification Strategy	Instruction↑	Physics↑	Plan time (min) ↓
Unverified	1.52	4.56	0
Generation-based	1.92	4.62	38.99
Sketch-based	1.90	4.66	4.08

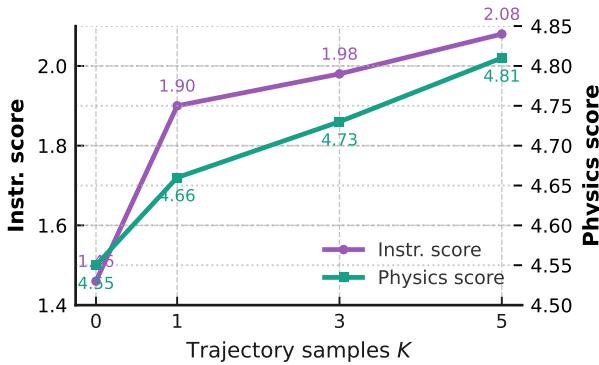


Figure 6. Ablation on the number of sampled trajectories K during planning. Larger K values enable stronger verifier-guided selection. $K = 0$ denotes the setting without a verifier, which is identical to the VideoMSG baseline.

tion (equivalent to VideoMSG [23]), the planner commits to a single trajectory without verification, resulting in limited instruction adherence (1.46) and lower physical consistency (4.55). Introducing even a small sampling budget ($K = 1$) with refinement from yields noticeable gains, as the verifier can reject implausible motions and select improved alternatives. Performance continues to increase with larger K , reflecting the benefit of exploring a broader trajectory space. Our full configuration ($K = 5$) achieves the strongest results across both instruction following and physics coherence, demonstrating that moderate test-time sampling is sufficient for robust trajectory selection without compromising efficiency.

Verification Strategy. To assess the efficiency of verifying on lightweight sketches versus fully generated videos, we compare our sketch-based verification with two alternatives: (1) Generation-based verification, which evaluates semantic and physical quality after full diffusion-based video synthesis, and (2) Verifier-regeneration, which regenerates videos based on verifier feedback from sketches. All methods use the same verifier model and multimodal planner. In both alternatives, the planner follows a generate–render–verify–refine loop to update the trajectory. Because full video generation is extremely expensive (rendering all ~ 100 candidates require over 30 GPU-hours), we adopt a practical two-round verify–regenerate scheme,

where the planner generates one trajectory per round, receives verifier feedback, then refines it into one updated trajectory in the next round. As shown in Table 6, verification improves motion quality across all settings, and our sketch-based approach matches or surpasses full video verification while achieving a nearly $10\times$ speedup. Unlike full videos which are expensive to render and often contain diffusion artifacts that mislead the verifier, sketches isolate motion from appearance, enabling cleaner and more reliable spatial–temporal evaluation.

6. Conclusion

We introduced a verifier-guided test-time planning framework for physically grounded video generation that decouples motion planning from synthesis. By integrating a multimodal verifier into the trajectory sampling loop and employing sketch-based proxy rendering, our approach enables efficient and reliable evaluation of candidate motions prior to video synthesis. This design allows the model to generate semantically coherent, physically plausible, and temporally smooth videos while reducing planning cost by nearly an order of magnitude compared to iterative refinement methods. Comprehensive experiments on WORLDMODELBENCH and PHYWORLDBENCH demonstrate substantial improvements in instruction following, physical law adherence, and motion realism over SoTA baselines.

7. Acknowledgements

We thank Justin Chih-Yao Chen, Jaemin Cho, Elias Stengel-Eskin, Zaid Khan, and Shoubin Yu for their helpful feedback. This work was supported by NSF-AI Engage Institute DRL-2112635, ARO Award W911NF2110220, ONR Grant N00014-23-1-2356, DARPA ECOLE Program No. HR00112390060, and a Capital One Research Award. The views contained in this article are those of the authors and not of the funding agency.

References

- [1] Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos world foundation model platform for physical ai. *arXiv preprint arXiv:2501.03575*, 2025. 6, 12
- [2] Pierfrancesco Ardino, Marco De Nadai, Bruno Lepri, Elisa Ricci, and Stéphane Lathuilière. Click to move: Controlling video generation with sparse motion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, page 14749–14758, 2021. 1
- [3] Hritik Bansal, Clark Peng, Yonatan Bitton, Roman Goldenberg, Aditya Grover, and Kai-Wei Chang. Videophy-2: A challenging action-centric physical commonsense evaluation in video generation, 2025. 3

- [4] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023. 1
- [5] Google AI / Google DeepMind. Gemini 2.5 (stable release). <https://ai.google.dev/models/gemini>, 2025. Multimodal large language model, model code: gemini-2.5-pro / gemini-2.5-flash. 5
- [6] Nate Gillman, Charles Herrmann, Michael Freeman, Daksh Aggarwal, Evan Luo, Deqing Sun, and Chen Sun. Force prompting: Video generation models can learn and generalize physics-based control signals. *arXiv preprint arXiv:2505.19386*, 2025. 3
- [7] Jiaxi Gu, Shicong Wang, Haoyu Zhao, Tianyi Lu, Xing Zhang, Zuxuan Wu, Songcen Xu, Wei Zhang, Yu-Gang Jiang, and Hang Xu. Reuse and diffuse: Iterative denoising for text-to-video generation. *arXiv preprint arXiv:2309.03549*, 2023. 1
- [8] Jing Gu, Xian Liu, Yu Zeng, Ashwin Nagarajan, Fangrui Zhu, Daniel Hong, Yue Fan, Qianqi Yan, Kaiwen Zhou, Ming-Yu Liu, et al. "phyworldbench": A comprehensive evaluation of physical realism in text-to-video models. *arXiv preprint arXiv:2507.13428*, 2025. 1, 5, 13
- [9] Lin Han, Abhay Zala, Jaemin Cho, and Mohit Bansal. Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning. *arXiv preprint arXiv:2309.15091*, 2023. 1
- [10] Yutong Hao, Chen Chen, Ajmal Saeed Mian, Chang Xu, and Daochang Liu. Enhancing physical plausibility in video generation by reasoning the implausibility. *arXiv preprint arXiv:2509.24702*, 2025. 3
- [11] Yingqing He, Menghan Xia, Haoxin Chen, Xiaodong Cun, Yuan Gong, Jinbo Xing, Yong Zhang, Xintao Wang, Chao Weng, Ying Shan, et al. Animate-a-story: Storytelling with retrieval-augmented video generation. *arXiv preprint arXiv:2307.06940*, 2023. 2
- [12] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 1
- [13] hpcatech. Open-sora: Democratizing efficient video production for all, 2024. 6, 12
- [14] Haoyang Huang, Guoqing Ma, Nan Duan, Xing Chen, Changyi Wan, Ranchen Ming, Tianyu Wang, Bo Wang, Zhiying Lu, Aojie Li, et al. Step-video-ti2v technical report: A state-of-the-art text-driven image-to-video generation model. *arXiv preprint arXiv:2503.11251*, 2025. 6, 12
- [15] Ziqi Huang, Ning Yu, Gordon Chen, Haonan Qiu, Paul Debevec, and Ziwei Liu. Vchain: Chain-of-visual-thought for reasoning in video generation. *arXiv preprint arXiv:2510.05094*, 2025. 1, 2
- [16] Sangwon Jang, Taekyung Ki, Jaehyeong Jo, Jaehong Yoon, Soo Ye Kim, Zhe Lin, and Sung Ju Hwang. Frame guidance: Training-free guidance for frame-level control in video diffusion models. *arXiv preprint arXiv:2506.07177*, 2025. 1
- [17] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 3, 12
- [18] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanyvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 6, 12
- [19] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 3, 5, 12, 13
- [20] Daeun Lee, Jaehong Yoon, Jaemin Cho, and Mohit Bansal. Videorepair: Improving text-to-video generation via misalignment evaluation and localized refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 1, 3
- [21] Chenyu Li, Oscar Michel, Xichen Pan, Sainan Liu, Mike Roberts, and Saining Xie. Pisa experiments: Exploring physics post-training for video diffusion models by watching stuff drop. *arXiv preprint arXiv:2503.09595*, 2025. 3
- [22] Dacheng Li, Yunhao Fang, Yukang Chen, Shuo Yang, Shiyi Cao, Justin Wong, Michael Luo, Xiaolong Wang, Hongxu Yin, Joseph E Gonzalez, et al. Worldmodelbench: Judging video generation models as world models. *arXiv preprint arXiv:2502.20694*, 2025. 1, 5, 6, 12
- [23] Jialu Li, Shoubin Yu, Han Lin, Jaemin Cho, Jaehong Yoon, and Mohit Bansal. Training-free guidance in text-to-video generation via multimodal planning and structured noise initialization. *arXiv preprint arXiv:2504.08641*, 2025. 1, 2, 6, 7, 9
- [24] Long Lian, Baifeng Shi, Adam Yala, Trevor Darrell, and Boyi Li. Llm-grounded video diffusion models. *arXiv preprint arXiv:2309.17444*, 2023. 2
- [25] Han Lin, Abhay Zala, Jaemin Cho, and Mohit Bansal. Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning. In *First Conference on Language Modeling*, 2024. 2
- [26] Wang Lin, Liyu Jia, Wentao Hu, Kaihang Pan, Zhongqi Yue, Wei Zhao, Jingyuan Chen, Fei Wu, and Hanwang Zhang. Reasoning physical video generation with diffusion timestep tokens via reinforcement learning. *arXiv preprint arXiv:2504.15932*, 2025. 3
- [27] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 3, 12
- [28] Shaowei Liu, Zhongzheng Ren, Saurabh Gupta, and Shenlong Wang. Physgen: Rigid-body physics-grounded image-to-video generation. In *European Conference on Computer Vision*, pages 360–378. Springer, 2024. 3
- [29] Do Xuan Long, Xingchen Wan, Hootan Nakhost, Chen-Yu Lee, Tomas Pfister, and Sercan Ö Arik. Vista: A test-time self-improving video generation agent. *arXiv preprint arXiv:2510.15831*, 2025. 3

- [30] Antonio Montanaro, Luca Savant Aira, Emanuele Aiello, Diego Valsesia, and Enrico Magli. Motioncraft: Physics-based zero-shot video generation. *Advances in Neural Information Processing Systems*, 37:123155–123181, 2024. 3
- [31] Saman Motamed, Laura Culp, Kevin Swersky, Priyank Jaini, and Robert Geirhos. Do generative video models understand physical principles?, 2025. 3
- [32] OpenAI. Gpt-4.1. <https://openai.com/research/gpt-4>, 2024. Large language model. 5
- [33] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kun-chang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks, 2024. 3, 5, 12
- [34] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data. In *International Conference on Learning Representations (ICLR)*, 2023. 1
- [35] Xiyang Tan, Ying Jiang, Xuan Li, Zeshun Zong, Tianyi Xie, Yin Yang, and Chenfanfu Jiang. Physmotion: Physics-grounded dynamics from a single image. *arXiv preprint arXiv:2411.17189*, 2024. 3
- [36] Wan Team. Wan: Open and advanced large-scale video generative models, 2025. 6, 12
- [37] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 1
- [38] Angtian Wang, Haibin Huang, Jacob Zhiyuan Fang, Yiding Yang, and Chongyang Ma. Ati: Any trajectory instruction for controllable video generation. *arXiv preprint arXiv:2505.22944*, 2025. 5, 12
- [39] Chen Wang, Chuhao Chen, Yiming Huang, Zhiyang Dou, Yuan Liu, Jiatao Gu, and Lingjie Liu. Physctrl: Generative physics for controllable and physics-grounded video generation. *arXiv preprint arXiv:2509.20358*, 2025. 3
- [40] Hanlin Wang, Hao Ouyang, Qiuyu Wang, Wen Wang, Ka Leong Cheng, Qifeng Chen, Yujun Shen, and Limin Wang. Levitor: 3d trajectory oriented image-to-video synthesis. *arXiv preprint arXiv:2412.15214*, 2024. 1
- [41] Peiyao Wang, Weining Wang, and Qi Li. Physcorr: Dual-reward dpo for physics-constrained text-to-video generation with automated preference selection. *arXiv preprint arXiv:2511.03997*, 2025. 3
- [42] Zun Wang, Jaemin Cho, Jialu Li, Han Lin, Jaehong Yoon, Yue Zhang, and Mohit Bansal. Epic: Efficient video camera control learning with precise anchor-video guidance. *arXiv preprint arXiv:2505.21876*, 2025. 3
- [43] Zun Wang, Jialu Li, Han Lin, Jaehong Yoon, and Mohit Bansal. Dreamrunner: Fine-grained storytelling video generation with retrieval-augmented motion adaptation. In *The AAAI Conference on Artificial Intelligence*, 2026. 2
- [44] Runpu Wei, Zijin Yin, Shuo Zhang, Lanxiang Zhou, Xueyi Wang, Chao Ban, Tianwei Cao, Hao Sun, Zhongjiang He, Kongming Liang, et al. Omnieraser: Remove objects and their effects in images with paired video-frame data. *arXiv preprint arXiv:2501.07397*, 2025. 3, 12
- [45] Qiyao Xue, Xiangyu Yin, Boyuan Yang, and Wei Gao. Phyt2v: Llm-guided iterative self-refinement for physics-grounded text-to-video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 1, 3, 6
- [46] Ling Yang, Zhaochen Yu, Chenlin Meng, Minkai Xu, Stefano Ermon, and Bin Cui. Mastering text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms. In *Forty-first International Conference on Machine Learning*, 2024. 3
- [47] Xindi Yang, Baolu Li, Yiming Zhang, Zhenfei Yin, Lei Bai, Liqian Ma, Zhiyong Wang, Jianfei Cai, Tien-Tsin Wong, Huchuan Lu, et al. Vlipp: Towards physically plausible video generation with vision and language informed physical prior. *arXiv preprint arXiv:2503.23368*, 2025. 3
- [48] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazhen Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 6, 12
- [49] Jaehong Yoon, Shoubin Yu, and Mohit Bansal. Raccoon: A versatile instructional video editing framework with auto-generated narratives. In *Conference on Empirical Methods in Natural Language Processing*, 2025. 2
- [50] Shoubin Yu, Jacob Zhiyuan Fang, Jian Zheng, Gunnar Sigurdsson, Vicente Ordonez, Robinson Piramuthu, and Mohit Bansal. Zero-shot controllable image-to-video animation via motion decomposition. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 3332–3341, 2024.
- [51] Shoubin Yu, Difan Liu, Ziqiao Ma, Yicong Hong, Yang Zhou, Hao Tan, Joyce Chai, and Mohit Bansal. Veggie: Instructional editing and reasoning of video concepts with grounded generation. *arXiv preprint arXiv:2503.14350*, 2025.
- [52] Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi Feng, and Qibin Hou. Storydiffusion: Consistent self-attention for long-range image and video generation. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [53] Shaobin Zhuang, Kunchang Li, Xinyuan Chen, Yaohui Wang, Ziwei Liu, Yu Qiao, and Yali Wang. Vlogger: Make your dream a vlog. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8806–8817, 2024. 2

A. Appendix

A.1. Detailed Implementations

A.1.1. SketchVerify Pipeline Specification

In this section, we detail the full SketchVerify implementation pipeline, including high-level planning, object detection and masking, background extraction, test-time trajectory search, sketch rendering, multimodal verification, and final video generation.

High-Level Planning. Given the input text prompt, we first perform a high-level decomposition of the described action into a sequence of structured sub-instructions. This step is carried out by GPT-4.1 (multimodal version) using a constrained prompt that requires the model to output: (a) a list of action segments, (b) their temporal ordering, and (c) the moving objects involved in each segment (detailed prompt in Fig. 9). The planner is asked to produce M sub-instructions ($M \in [1, 4]$), depending on the complexity of the prompt. Each sub-instruction corresponds to an independent phase of motion planning and receives its own temporal budget T_i , with all phase lengths summing to the fixed $\sum_{i=1}^M T_i = 41$ total frames used throughout the pipeline. All high-level plans are parsed through a strict JSON schema that enforces the required fields (action, duration, object ids). Malformed or under-specified outputs are automatically rejected, and the planner is re-sampled. This guarantees that downstream modules always operate on well-structured, machine-readable action plans.

Object Detection and Masks. We use Grounding DINO [27, 33] for text-conditioned object detection (based on the detected objects in high-level planning) with a confidence threshold of 0.3, followed by SAM-HQ [17] for segmentation. For each detected moving object, we retain the highest-scoring instance mask and compute the corresponding bounding box. Boxes are normalized to $[0, 1]^2$ following the image coordinate system.

Background Removal. To obtain a clean static background, we use FLUX.1-dev [19] with the Omnieraser [44] LoRA. All moving-object masks are combined into a single inpainting mask. The background is generated at the same resolution as the input image. We generate the background image with 28 diffusion steps and $cfg = 3.5$.

Test-Time Search. We adopt GPT-4.1 as the default multimodal planner. Each planning call produces $K = 5$ candidate trajectories, each of length T_i , where each trajectory contains per-frame bounding box coordinates for all moving objects. All planner outputs are validated through a structure-enforcing JSON parser, and malformed samples are automatically re-sampled. We set the temperature to 1.0. Diversity filtering is enforced by requiring an ℓ_2 distance of at least 0.05 (in normalized coordinates) between trajectories.

Sketch Rendering. For each candidate trajectory, we gen-

erate a lightweight video sketch of T_i frames using object sprites cropped from the first frame I_0 and composited onto the static background. All sketches are rendered at the input image resolution and saved as MP4 at 4 fps for verification.

Multimodal Verification. We use Gemini 2.5-Flash as the default verifier. Two scores are produced per candidate:

- a semantic alignment score from the first/last-frame comparison;
- a physics plausibility score from the full sketch video.

The planner uses a weighted combination with default weights $(\lambda_{sem}, \lambda_{phys}) = (0.5, 0.5)$. For all four physical laws, we also set $\lambda_l = 0.25$ for all $l \in \mathcal{L} = \{\text{Newton, Penetration, Gravity, Deformation}\}$.

For scoring, we use the following criteria:

- 1.0: Perfect plan alignment and physical-law coherence.
- 0.7–0.9: Good alignment with minor deviations.
- 0.4–0.6: Partial alignment with some correct aspects.
- 0.0–0.3: Poor alignment; does not achieve the goal or directly breaks physical laws.

Video Generation Model. We use the ATI-14B model [38] to generate 81-frame 480p videos. We generate with 40 steps and $cfg = 5.0$. The conditions include the input image, the input text prompt, and the trajectory plan obtained from the planner.

A.1.2. Baseline and Hardware Specifications

We use Wan2.1-14B-480p-I2V [36], CogVideoX-5B [48], Cosmos-Predict2-2B [1], HunyuanVideo-I2V [18], OpenSora-I2V [13], and Step-Video-I2V [14] as baseline open-source TI2V models. For these models, we directly use the official implementations and sample videos at 480p with 81 frames using 50 diffusion steps. For PhyT2V and VideoMSG, we replace the backbone video generation model with Wan2.1 for fair comparison. All experiments are conducted on NVIDIA A100 80G and NVIDIA RTX A6000 GPUs.

A.1.3. Per-step Runtime

On average, high-level planning takes 14.16 s. Object detection, segmentation, and background inpainting require 108 s. For each sub-instruction, test-time planning takes 72.5 s on average, consisting of 20.3 s for trajectory sampling and 52.2 s for multimodal verification. All timings are measured on a single NVIDIA A100 GPU using the standard-speed APIs of GPT-4.1 and Gemini-2.5.

A.2. Benchmark Details and Metric Definitions

A.2.1. Benchmark Details

WorldModelBench [22] evaluates video generation models as world models, focusing on instruction following, physical plausibility, and commonsense temporal behavior. It contains 7 domains and 56 subdomains across 350 image/text-conditioned tasks. It supports both I2V and T2V settings (we use I2V).

PhyWorldBench [8] focuses on fine-grained physical realism, testing whether videos obey Newtonian laws, gravitational motion, and object–interaction constraints. This benchmark includes 350 text prompts describing physically grounded events. It is a T2V benchmark; therefore, we generate a first frame using FLUX [19] and then perform I2V generation.

A.2.2. Metric Details

WorldModelBench provides a vision–language model (MLLM) scorer trained on 67K human annotations. Each generated video outputs scalar scores in three categories:

- Instruction Following: Measures how well the generated motion follows the input instruction. Scores range from 1 to 3: 3 = correct motion, 2 = partially correct, 1 = incorrect. The score is produced by the MLLM via a trained textual comparison head.
- Physics Coherence: Evaluates adherence to natural physical priors across six dimensions, each in $[0, 1]$: Newtonian motion, deformation consistency, fluid dynamics, object penetration, gravity coherence, and frame-level physics consistency.
- Commonsense Consistency: Includes per-frame visual realism and motion smoothness/continuity, both in $[0, 1]$.

All metrics are computed by the MLLM via prompt-driven scoring heads.

PhyWorldBench uses its own MLLM-based evaluator and reports three pass-rate metrics. For each video, eight frames are sampled uniformly and passed to a proprietary SoTA MLLM. Questions about the following criteria are asked:

- Obj+Evt: whether the described objects appear and the event occurs.
- Phys. Std: whether motion follows expected physical laws (gravity, collision response, continuous motion, no penetration).
- All: counted as correct only if both Obj+Evt and Phys. Std pass.

Each is computed as a binary decision per prompt and averaged into a percentage. We use GPT-5 as the proprietary MLLM evaluator.

A.3. Extra Qualitative Results

We provide additional qualitative examples in Fig. 7 and Fig. 8. These examples highlight the consistency of SketchVerify across diverse scenes and motion types. On WorldModelBench tasks (Fig. 7), our trajectories produce smoother and more semantically aligned motions than baseline models. On PhyWorldBench (Fig. 8), our method more reliably maintains physical plausibility, avoiding common failure modes such as objects floating against gravity or moving on their own in violation of Newton’s first law.

A.4. Prompt Template for SketchVerify Pipeline

We provide the full set of prompt templates used in our system, covering high-level planning in Fig. 9, object proposal in Fig. 10, trajectory generation in Fig. 11, semantic alignment verification in Fig. 12, and physical plausibility checking in Fig. 13. These prompts define the behavior of each module and ensure consistent outputs across tasks and benchmarks.

A.5. Limitations

While SketchVerify substantially improves motion planning quality, several limitations remain. Our verification module primarily evaluates coarse object motion and high-level physical plausibility; however, capturing fine-grained physics, such as frictional forces, collision responses, or other continuous dynamics that typically require differentiable simulation, would require additional modeling beyond our current verifier-based design. Moreover, because both the planner and verifier are external MLLMs, they may occasionally produce incorrect judgments, which can lead to suboptimal candidate selection. Finally, since motion is represented through 2D bounding boxes, the framework can struggle with fine-grained 3D interactions such as detailed affordances or fluid-like behavior, and the realism of the final video remains bounded by the capability of the underlying video generation model. We expect these limitations to diminish as stronger video generators and verification models continue to improve.

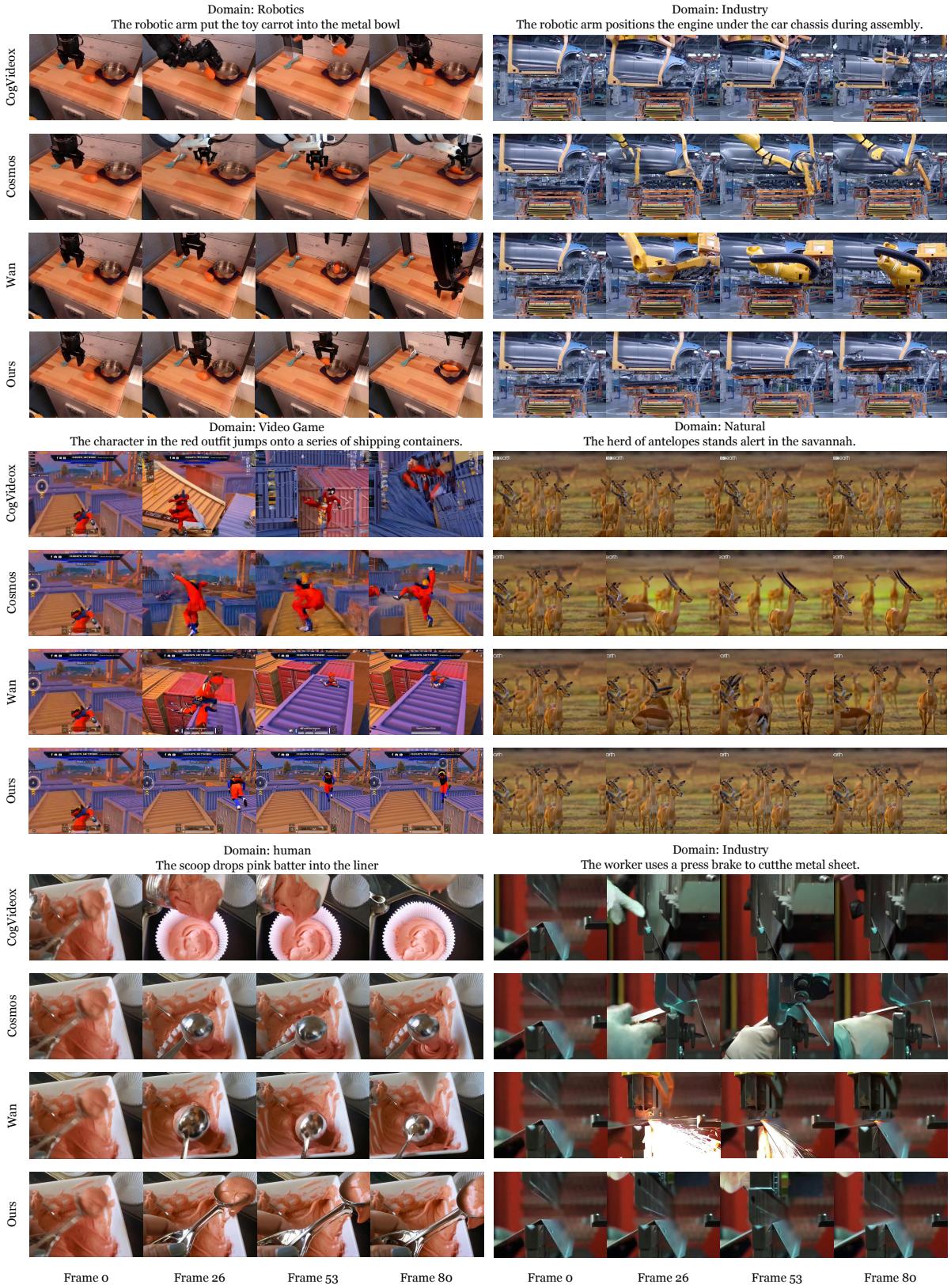


Figure 7. Qualitative comparison on WorldModelBench.

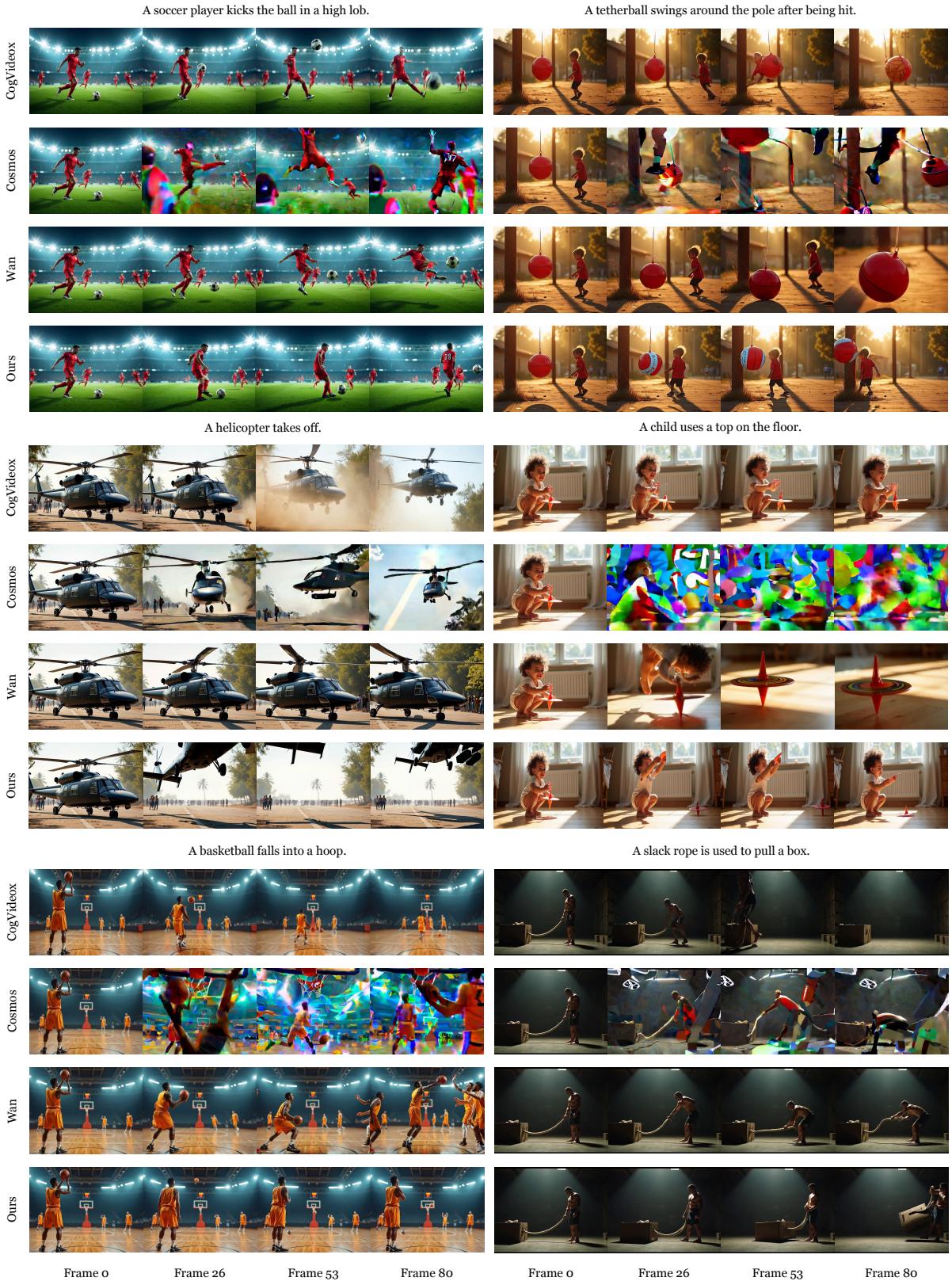


Figure 8. Qualitative comparison on PhyWoldBench.

Global Movement Planning Prompt (GPT-4.1 with Tool Calling)

System message:

You are a video motion planning expert. You have EXACTLY {total_frames} frames to complete the ENTIRE task.

User message:

Text prompt: "{text_prompt}"
Available frames: EXACTLY {total_frames}.
Current objects in frame 1:
- {label_1}: currently at [x_min, y_min, x_max, y_max]
- {label_2}: currently at [...] (etc.)

Return ONLY a call to submit_movement_plan with a complete plan that finishes by frame {total_frames}.

Function schema (submit_movement_plan):

```
{  
    "task_breakdown": {  
        "complete_objective",  
        "phase_1", "phase_2", "phase_3", ...  
        "success_criteria"  
    },  
  
    "frame_allocation": {  
        "phase_1",  
        "phase_2",  
        "phase_3", ...  
    },  
  
    "moving_objects": [...],  
    "static_objects": [...],  
  
    "detailed_timeline": {  
        frame_1, frame_3, frame_6, frame_8,  
        frame_10, frame_12, frame_15, frame_18, ...  
        frame_{total_frames}  
    },  
  
    "movement_plans": {  
        "<object_name>": {  
            "movement_type",  
            "total_distance",  
            "movement_phases": {  
                "phase_1_frames",  
                "phase_2_frames",  
                "phase_3_frames", ...  
            },  
        }  
    },  
  
    "completion_verification"  
}
```

Figure 9. Prompt used for high-level planning

Object Proposal Prompt (System + User)

System Message:

You are an expert in video generation and object-centric scene analysis. Given the first frame of a video and a text description, determine which objects should be added, moved, or animated to fulfill the described action.

Your responsibilities:

1. Analyze the given frame and identify all existing objects.
2. Based on the text prompt, determine which additional objects (if any) must be introduced to achieve the described event.
3. Identify which objects|either existing or newly added|must move or animate to satisfy the prompt.
4. Ensure that proposed object placement and motion are physically plausible and consistent with real-world interactions.
5. Focus on major objects that materially affect the scene. If multiple parts form a single rigid object, treat them as one entity.
6. For each object name, use minimal wording (1{3 words), concise and unambiguous.

Return the result strictly in the following JSON format:

```
{  
  "scene_analysis": "Brief description of the current frame",  
  "existing_objects": [...],  
  "objects_to_add": [  
    {  
      "name": "object_name",  
      "reasoning": "why this object is required",  
      "movement_type": "static / linear / curved / complex",  
      "priority": "high / medium / low"  
    }  
  ],  
  "moving_objects": [...],  
  "static_objects": [...]  
}
```

User Message:

Text prompt: "<TEXT_PROMPT>"

Using the first-frame image provided, analyze how the scene should be modified or animated to satisfy the text description.

Please determine:

- Which required objects are currently missing, based on the prompt.
- Which objects must move or animate to create the described action.
- Which objects remain static as part of the background.
- Realistic object placement and timing relative to the prompt's intent.

Produce the full JSON output exactly as specified in the system instructions.

Figure 10. Prompt used for object proposal

Sub-Instruction Trajectory Planning Prompt (GPT-4.1)

System Message:

You are a video motion planning expert generating trajectories for frames {CHUNK_START} to {CHUNK_END}.

Current phase: {PHASE_NAME}
Phase description: {PHASE_DESCRIPTION}
Total frames in video: {TOTAL_FRAMES_NUM}

COORDINATE SYSTEM:
{COORDS_GUIDE}

DIRECTIONAL MAPPINGS:
- RIGHT: x1 += delta; x2 += delta
- LEFT: x1 -= delta; x2 -= delta
- UP: y1 -= delta; y2 -= delta
- DOWN: y1 += delta; y2 += delta

Focus ONLY on moving objects: {MOVING_OBJECTS}
Ignore static objects in outputs: {STATIC_OBJECTS}

User Message:

Text prompt: "{TEXT_PROMPT}"

{HISTORY_TEXT}

Generate a smooth trajectory from frame {CHUNK_START} to frame {CHUNK_END} for this {PHASE_NAME}.

IMPORTANT: Since multiple trajectories will be generated, explore DIFFERENT valid motion paths. Consider variations in:

- Path shape (straight, curved, arc)
- Speed profile (constant, accelerating, decelerating)
- Intermediate waypoints (different approaches to the goal)

For each object, you should maintain its size (box dimensions) and only change its position unless you are specifically instructed to resize.

For EACH frame from {CHUNK_START} to {CHUNK_END}, output:
Frame_N: [{"object_name": [x1, y1, x2, y2]}, ...], caption: <description>

Requirements:

- Smooth motion (delta 0.03{0.08 per frame})
- Consistent with phase objectives
- Maintain object sizes
- Only include moving objects: {MOVING_OBJECTS}

Figure 11. Prompt used for sub-instruction trajectory planning

Plan-Alignment Verifier Prompt (GPT-4.1)

System Message:

You are an expert at evaluating video motion trajectories.
Your task is to verify if the motion from the first frame to the last frame aligns with the expected phase goal.

Rate the alignment on a scale of 0.0 to 1.0 where:

- 1.0 = Perfect alignment, the last frame clearly achieves the phase goal
- 0.7-0.9 = Good alignment with minor deviations
- 0.4-0.6 = Partial alignment, some aspects correct
- 0.0-0.3 = Poor alignment, does not achieve the goal

Return ONLY a JSON object with:

```
{  
  "score": <float between 0 and 1>,  
  "explanation": "<brief explanation of why this score was given>"  
}
```

User Message (paired with first/last-frame images):

Phase: {PHASE_NAME}
Phase Description: {PHASE_DESCRIPTION}
Expected End Goal: {END_GOAL}

Please compare the FIRST frame (starting state) with the LAST frame (ending state). Does the last frame show that the phase goal has been achieved?

Consider:

- Object positions relative to the goal
- Whether objects moved in the expected direction
- Whether the motion is consistent with the phase description

Figure 12. Prompt used for plan-alignment verification

Physical Plausibility Verifier Prompt (Gemini 2.5-Flash)

Text Prompt (sent together with the sketch video file):

Analyze this video sequence and evaluate whether the motion obeys physical laws.

IMPORTANT NOTE: This video is generated by copy & paste composition – each frame is created by pasting objects onto a background. Therefore, please focus on evaluating the movement trajectories and positions of individual objects across frames, not visual quality, shadows, or composition artifacts.

Consider for each moving object (one of the laws):

Newtonian Consistency: acceleration / deceleration should be physically plausible;

Penetration Violation: objects must not pass through static elements;

Gravitational Coherence: objects should not be floating in the air without anything holding it

Deformation Consistency: object size should remain stable unless specified.

Rate the physical plausibility on a scale of 0.0 to 1.0 where:

- 1.0 = Perfectly realistic, obeys this physical laws
- 0.7-0.9 = Mostly realistic with minor issues
- 0.4-0.6 = Some unrealistic aspects but acceptable
- 0.0-0.3 = Highly unrealistic, violates physics (teleportation, impossible speeds, etc.)

Return ONLY a JSON object:

```
{  
  "score": <float between 0 and 1>,  
  "explanation": "<brief explanation focusing on object movement quality,  
                 highlight any physics violations>"  
}
```

Example:

A example for the specific physical Law

Figure 13. Prompt used for physical plausibility verification