# Wow, wo, val! A Comprehensive Embodied World Model Evaluation Turing Test

Chun-Kai Fan[1,2,*]   Xiaowei Chi[2,3,*]   Xiaozhu Ju[2,†]   Hao Li[2]   Yong Bao[2]

Yu-Kai Wang[1,2]   Lizhang Chen[1]   Zhiyuan Jiang[2]   Kuangzhi Ge[1,2]   Ying Li[1]

Weishi Mi[2]   Qingpo Wuwu[1]   Peidong Jia[1,2]   Yulin Luo[1]   Kevin Zhang[1,2]

Zhiyuan Qin[2]   Yong Dai[2]   Sirui Han[3]   Yike Guo[3]   Shanghang Zhang[1,‡]   Jian Tang[2,‡]

[1]State Key Laboratory of Multimedia Information Processing, School of Computer Science, Peking University

[2]Beijing Innovation Center of Humanoid Robotics   [3]The Hong Kong University of Science and Technology

## Abstract

*As world models gain momentum in Embodied AI, an increasing number of works explore using video foundation models as predictive world models for downstream embodied tasks like 3D prediction or interactive generation. However, before exploring these downstream tasks, video foundation models still have two critical questions unanswered: (1) whether their generative generalization is sufficient to maintain perceptual fidelity in the eyes of human observers, and (2) whether they are robust enough to serve as a universal prior for real-world embodied agents. To provide a standardized framework for answering these questions, we introduce the Embodied Turing Test benchmark: **WoW-World-Eval (Wow,wo,val)**. Building upon 609 robot manipulation data, Wow-wo-val examines five core abilities, including perception, planning, prediction, generalization, and execution. We propose a comprehensive evaluation protocol with 22 metrics to assess the models' generation ability, which achieves a high Pearson Correlation between the overall score and human preference (>0.93) and establishes a reliable foundation for the Human Turing Test. On Wow-wo-val, models achieve only 17.27 on long-horizon planning and at best 68.02 on physical consistency, indicating limited spatiotemporal consistency and physical reasoning. For the Inverse Dynamic Model Turing Test, we first use an IDM to evaluate the video foundation models' execution accuracy in the real world. However, most models collapse to $\approx$ 0% success, while WoW maintains a 40.74% success rate. These findings point to a noticeable gap between the generated videos and the real world, highlighting the urgency and necessity of benchmarking World Model in Embodied AI.*

## 1. Introduction

World models – which capture an agent's understanding of how the world changes with actions [35] – have emerged as a pivotal concept in robotics and Embodied AI. In embodied settings, a world model allows a robot to understand and predict its environment [21, 32, 49, 56], and can function as the robot's "internal brain", enabling it to simulate future scenarios for planning and decision-making [17, 43, 72, 85], or operate as an environment simulator [47, 52, 88].

Compared with cutting-edge spatial-prediction world models [3, 79], embodied world models operate in context-rich environments that demand a deeper understanding of physical common sense. In robotics, the complexity and lack of standardization across setups further lead to a broad and diverse landscape of embodied-world-model designs. These models vary widely in their control conditions—ranging from approaches conditioned solely on images and language instructions [15–17, 72], to those incorporating keypoints [84], trajectories [9, 31, 34, 41, 51, 69], depth, semantics, and other modalities [50, 92, 99]. They also differ in camera configurations, requiring either single-view inputs [9, 15, 16, 31, 51, 69, 72, 84, 100] or multi-view setups [11, 17, 34, 52, 66, 73, 92]. Moreover, several recent efforts have begun to pursue cross-embodiment generalization [17, 37, 98]. Therefore, despite the broader research of world models in general-purpose robotics, two core questions still remain:

1. Can these models generalize well enough to maintain perceptual fidelity from a human perspective?
2. Are they robust and expressive enough to serve as universal priors for real-world embodied agents?

Existing video-generation benchmarks [29, 40, 46, 53, 54, 95, 97] largely target general-purpose settings or isolated dimensions and overlook the unique requirements of robotic world models. Most evaluations emphasize visual fidelity or coarse task success, but rarely assess deeper embodied abilities such as physical plausibility, planning ra-

---

[*]Equal contribution.

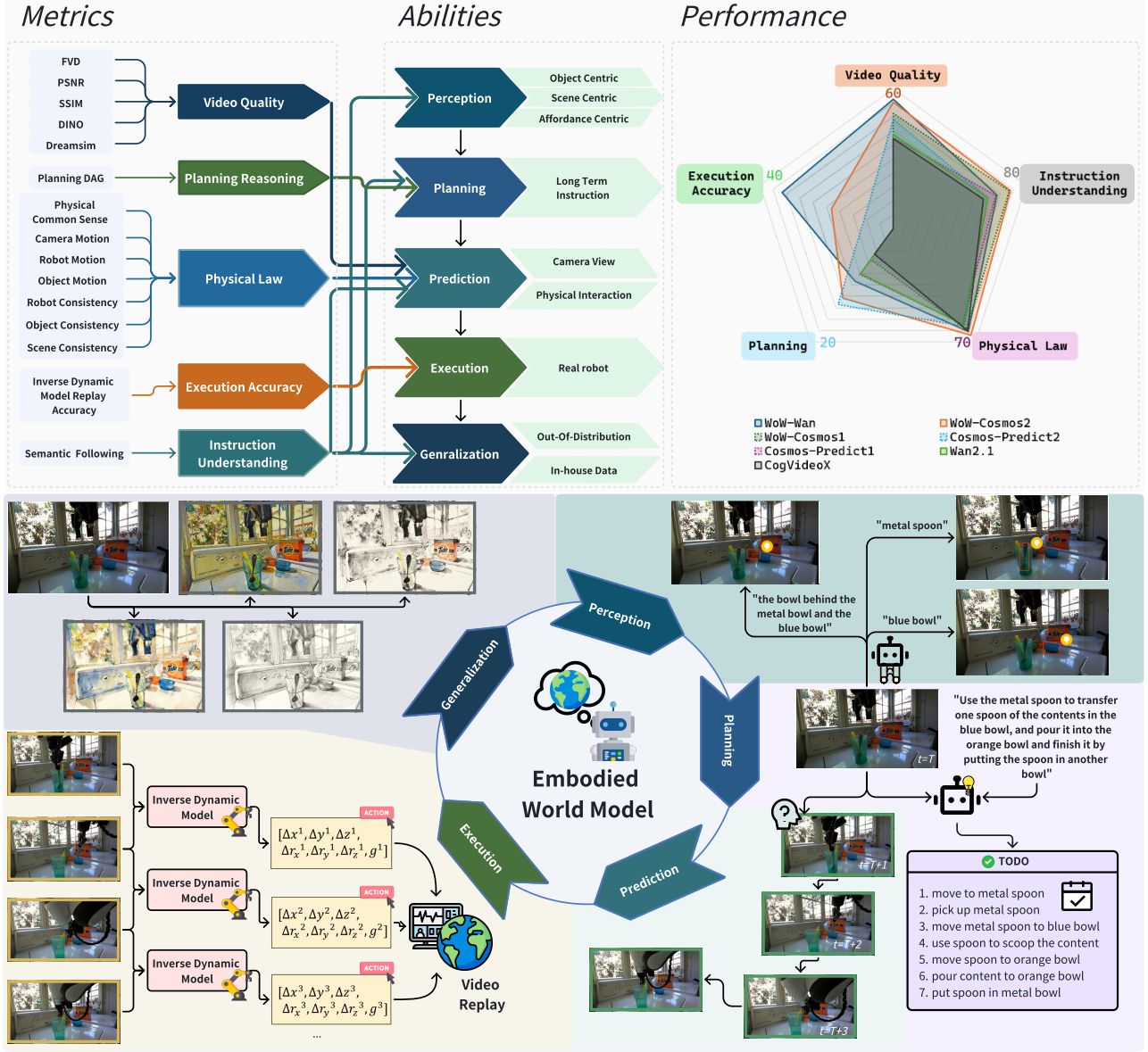[†]Project leader.

[‡]Corresponding author.

Figure 1. **The Overview of WoW-World-Eval. (Top-left)** A multi-faceted *Metrics* suite evaluates generated videos across five dimensions: Video Quality, Instruction Understanding, Planning Reasoning, Physical Law, and Execution Accuracy. **(Top-center)** These metrics align with five core *Abilities* of embodied world models: Perception, Planning, Prediction, Execution, and Generalization. **(Top-right)** Performance gaps across state-of-the-art models. **(Bottom)** The benchmark follows the embodied world model pipeline from Perception to Generalization.

tionality, and actionability. This gap makes progress difficult to measure: a model may score well on conventional video metrics [27, 65, 81, 86] yet produce physically impossible or contextually incorrect predictions in robotic scenarios. Our results further confirm this misalignment—standard video-quality scores correlate poorly with human judgments in embodied settings—highlighting the need for more reliable evaluation standards.

Consequently, in this paper, we address these challenges by proposing a new comprehensive benchmark for embodied world models, and use it to systematically evaluate foundational models under the simplest *image-to-video* setting. Our benchmark, **WoW-World-Eval**, as illustrated in Figure 1, is designed around the core capabilities that an embodied world model should possess: perceiving the environment, understanding and planning based on task instruc-

tions, predicting and simulating future world states, executing real-world interactions, and generalizing across diverse scenarios and embodiments. It contains about 609 robot manipulation samples with meticulous cleaning and annotation by human annotators. Also we incorporate 22 evaluation metrics aligned with our core dimensions across Video Quality, Instruction Understanding, Planning Reasoning, Physical Law, and Execution Accuracy.

Moreover, **WoW-World-Eval** follows the two-alternative forced-choice (2AFC) [28] methodology from psychophysics to establish the evaluation as a standardized Turing Test for generative video models. By collecting fine-grained human answers distinguishing real and generated videos, we compute the proportion of generated videos from each model that successfully fool human evaluators. Notably, an overall human preference alignment score exceeding Pearson Correlation = 0.93 demonstrates the effectiveness of our benchmark in evaluating high-quality generations and serves as a reliable proxy for the *Human Turing Test*.

In addition to the human-centered evaluation, we also introduce a *Machine Turing Test*—specifically, an *Inverse Dynamics Model (IDM) Turing Test*. In this setting, we assess whether videos generated by a model can "fool" an IDM that has only been trained on real-world execution sequences. If the generated videos lead the IDM to output plausible actions that are executable in the real world, it indicates that the model's outputs are indistinguishable from real data in terms of physical and action plausibility.

By evaluating existing models under this new benchmark, we reveal which models already exhibit credible world understanding and where they fall short. We believe this benchmark and the accompanying Turing Test criterion will provide a much-needed standard for the field, driving research towards embodied world models that can truly imagine the world with the accuracy and fidelity that robotics applications demand. Our contributions are threefold as follows:

- A comprehensive World Model Benchmark, **WoW-World-Eval**, focuses on the Embodied AI domain, introducing a new perspective with a novel framework for the five core abilities in the embodied world model.
- Based on **WoW-World-Eval**, we propose two novel Turing tests aligned with our benchmark, the Human Turing test and IDM Turing Test, that can distinguish the models' true ability as an embodied world model to simulate and interact with the real world.
- We curate 609 high-quality robot manipulation samples with careful human annotations in **WoW-World-Eval**. Using this benchmark, we conduct a comprehensive evaluation of various world models, whose performance provides new insights into the strengths, limitations, and generalization gaps of current embodied world models.

## 2. Related work

Table 1. **Comparison of benchmark features.** A checkmark (✓) indicates that the benchmark explicitly includes the corresponding evaluation dimension or metric, while a cross (×) indicates that it does not. In Core Dimensions, Perception(Percept); Prediction(Pred); Planning(Plan); Execution(Exec); Generalization(General); in Metrics Design, Video Quality(VQ); Instruction Understanding(IU); Physical Law(PL); Planning Reasoning(PR); Execution Accuracy(EA).

| Benchmark | Core Dimensions | | | | | Metrics Design | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Percept | Pred | Plan | Exec | General | VQ | IU | PL | PR | EA |
| Physics-IQ [61] | ✓ | ✓ | × | × | × | × | × | ✓ | × | × |
| PhyGenBench [59] | ✓ | ✓ | × | × | × | × | × | ✓ | × | × |
| T2V-CompBench [78] | ✓ | ✓ | × | × | × | × | ✓ | ✓ | × | × |
| VBench-2.0 [97] | ✓ | ✓ | × | × | ✓ | ✓ | ✓ | ✓ | × | × |
| WorldModelBench [46] | ✓ | ✓ | × | × | ✓ | × | ✓ | ✓ | × | × |
| EWMBench [95] | ✓ | ✓ | × | × | × | × | ✓ | ✓ | × | × |
| Ours | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

The evaluation of generative world models and video generation models is rapidly evolving. Existing benchmarks can be grouped into: (1) those targeting general video quality, (2) those assessing physical reasoning, and (3) benchmarks for holistic world model evaluation.

**Video Generation Benchmarks.** Early benchmarks mainly assessed visual fidelity and temporal coherence. TC-Bench [29] studies temporal compositionality with transition-aware metrics. [53] measures perception-aligned motion quality. T2V-CompBench [78] evaluates compositionality via MLLM- and tracking-based methods, challenging metrics like FVD. [15, 40, 54] introduce tasks related to controllability and anticipation. VBench-2.0 [97] focuses on "intrinsic faithfulness," adding *Physics* and *Commonsense* dimensions using VLM/LLM evaluators.

**Physical Reasoning Benchmarks.** Physics-IQ [61] introduces physics-based diagnostics. PhyGenBench [59] targets semantic and physical plausibility. PhysBench [18] evaluates VLM physical reasoning via QA tasks rather than pixel metrics. Videophy [4] benchmarks physical commonsense using rule-based and learned evaluators.

**World Model Benchmarks.** WorldModelBench [46] and EWMBench [95] evaluate core dimensions like Perception, Prediction, and Generalization. Critically, as shown in Table 1, these benchmarks do not assess the crucial core dimension of Planning and Execution. Furthermore, no prior benchmark provides a metrics design that covers Planning Reasoning or Execution Accuracy in the robotics domain. Our proposed **WoW-World-Eval** extends these efforts with a robotics-centric design, making it the most comprehensive embodied world model benchmark to date.
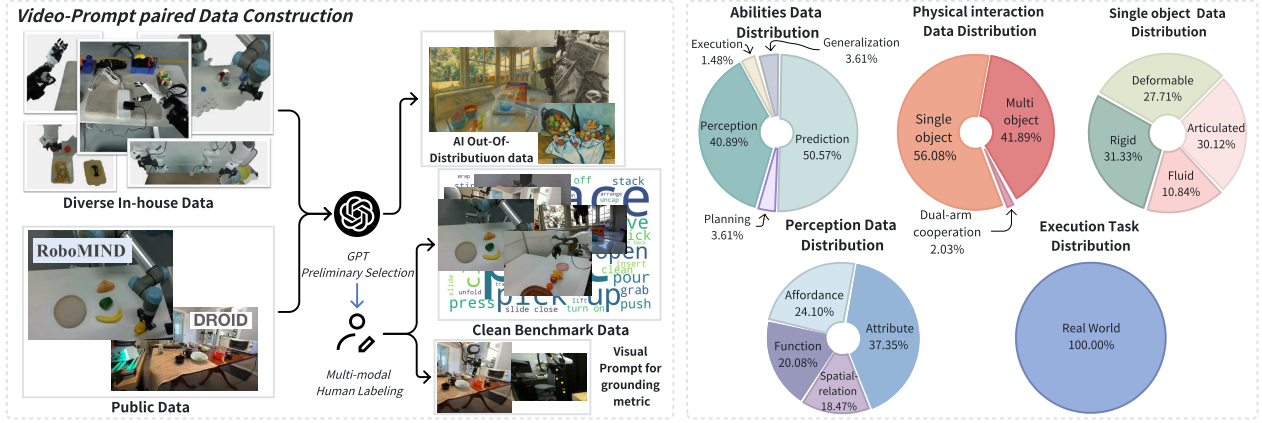
Figure 2. **Overview of WoW-World-Eval Data Construction Pipeline and Data Statistics. (Left)** Public and in-house data are cleaned sequentially by GPT and human annotators to produce high-quality samples consisting of an initial image, prompt, ground-truth video, and annotated keypoints. **(Right)** Data distribution across five different dimensions from overall to fine-grained.

## 3. WoW-World-Eval: A Multi-faceted Benchmark for Embodied World Models

In this section, we introduce **WoW-World-Eval**, a novel benchmark designed to evaluate embodied world models rigorously. We formulate the core evaluation task as **conditional video generation from an initial image and a text instruction (Image-Text-to-Video)**, a setting that directly probes a model's ability to understand a given state and execute a specified action. Moving beyond metrics of visual appeal, **WoW-World-Eval** specifically assessed a model's instruction understanding, planning, observation perception abilities, and understanding of physically grounded dynamics in embodied settings. For introducing our benchmark, we structure our presentation by first outlining the five core capabilities that a competent embodied world model must possess (Section 3.1), then detailing our data curation pipeline (Section 3.2), followed by our comprehensive evaluation metrics (Section 3.3).

### 3.1. Core Evaluation Dimensions

We posit that a truly effective embodied world model must demonstrate mastery across five fundamental and orthogonal dimensions, as outlined in our benchmark design (Figure 1, Top-center).

**Perception Understanding.** A world model must first accurately perceive and represent the environment to enable more reliable subsequent prediction and planning. We assess this through tasks requiring fine-grained **object recognition** [12, 13, 18, 58](attributes like color, shape, number, and size), **spatial understanding** [14, 23, 77, 90] (relative positions and arrangements), and **affordance recognition** [14, 33, 63] (identifying interactive parts of objects).

**Decision-making and Planning.** Embodied agents must execute long-horizon tasks [10, 12, 48, 74, 91]. Therefore, we assess a model's planning ability by challenging it to generate coherent video sequences for complex instructions. This requires implicitly understanding **task decomposition** [15, 80] into key sub-goals and respecting their **causal dependencies**. Hence, we collect 25 samples to fill in this long-term planning task, and transform the text instruction into a suitable description for the world model to plan. In order to evaluate the planning ability in the world model, we refer to the metric from RoboBench [57]. The detailed metric will be elaborated in section 3.3.

**Predictive Reasoning.** This dimension evaluates the model's internal physics engine. Given an initial state and an action, the model must generate a future that respects core physical principles such as **object permanence** [4, 18], **collision dynamics** [4, 18, 59, 61], and **trajectory plausibility** [24, 46, 68, 95, 97]. This directly probes the model's capacity to function as a world simulator [7]. Therefore, we design several sub-dimensions that focus on these principles, as illustrated by the pie chart in the right of Figure 2.

**Interactive Execution.** Interacting with the real world and executing on a real robot is the ultimate goal of an embodied world model. To examine this, we collect 9 different tasks from our in-house data that are compatible with the real robot experiments, for the model to generate. Furthermore, we will have the model generate videos based on this in-house data and use the Gripper-Centric Inverse Dynamics Model(GC-IDM) [17], to interpret generated videos into actions that can be executed for a real robot, finding out the **execution ability** in the world model. The detailed task names and the implementation will be described in the Appendix 10.3.

4

**Generative Generalization.** A universal world model should not only perform well on the In-Distribution data, but it should also generalize beyond the data it has seen before to demonstrate its generalization ability. For this reason, we test generalization on the in-house robot data that we collected, by using GPT-5 [39] to perform style transfer or image editing on it, and generating images that the world model had never seen before. We also collected some world-famous masterpiece paintings, such as *"Girl with a Pearl Earring"*, and asked the world model to execute the task instructions that humans created. These two types of images constitute our **In-house Data** and **Out-of-Distribution (OOD)** dimension.

### 3.2. Data Curation Pipeline

To systematically evaluate these capabilities, we build a principled, semi-automated data curation pipeline (Figure 2 Left). Our dataset combines open-source robotics data (e.g., RoboMIND [87], DROID [42]), in-house trajectories, and AI-generated OOD samples to ensure coverage and diversity. GPT-4o [39] is initially used as an intelligent annotator, scoring the matching level of video–instruction pairs based on our five capability dimensions and their subdivisions, specifically examining which instructions match which dimensions, thereby achieving large-scale filtering and coarse categorization. Human experts then verify all samples to guarantee category accuracy and resolve edge cases. Five additional annotators selected the best initial frames for generation (both the robotic arm and the manipulated object were in the same frame) and key point annotations (the robotic arm gripper, joints, and the manipulated object) on the initial frame for evaluation metrics.

Each benchmark entry contains: (1) a natural-language instruction, (2) an initial image, (3) a ground-truth video, and (4) annotated keypoints (in Prediction). In total, the dataset comprises 609 samples across all dimensions (Figure 2 Right). Prediction (50.57%) and Perception (40.89%) dominate the ability distribution. The perception subdivisions includes object attribute, object affordance, object function, and spatial-relation tasks (249 samples). Physical interaction covers single-object manipulation (56.08%), multi-object interaction (41.89%), and dual-arm cooperation (2.03%). We further include 107 non-occluded views and 54 semi-occluded views in the initial frames to test robustness. With execution data, we choose 9 real-world tasks from easy to hard for the settings.

### 3.3. Multi-faceted Evaluation Metrics

Our evaluation protocol is a suite of metrics designed to be as comprehensive as the capabilities we measure. We introduce several novel metrics alongside standard ones, grouped by the property they assess. For detailed information of all the metrics, see Appendix 9.

**Visual Fidelity.** To fully evaluate visual fidelity, we report standard video quality metrics spanning pixel-, perceptual-, and distribution-level quality perspectively: **PSNR**[27] measures pixel-level fidelity via the log ratio between the maximum signal value and the mean-squared error. To capture higher-level content consistency beyond pixel statistics, we compute **SSIM**[86] for a structural measure and **DINO** [65] as a semantic/instance-alignment signal, combine with **Dreamsim** [30], a human-aligned perceptual similarity metric trained from human triplet judgments. **FVD**[81] compares the distributions of real and generated videos by computing a Fréchet distance, which reflecting overall realism and temporal dynamics at the dataset level.

**Instruction Semantic Alignment.** We use GPT-4o [39] as a scalable evaluator to assess semantic alignment between the given instruction and the generated videos. Depending on whether ground-truth (GT) video is available, we adopt two methods:

- **With Ground-Truth:** We first prompt GPT-4o to extract structured descriptions (Initial-, Processing-, Final-state) from both the generated and GT videos. A vision–language model then scores their **Caption Score**, which scales from 1-to-5. In addition, GPT-4o also evaluates the generated video action-object pairs against the instruction to produce a **Sequence Match Score** (0–1 scale), which stands for the correctness of the order of the action-object pairs, and an **Execution Quality Score** (1–5 scale) for the correctness of the action-object pairs. *We report all three metrics in this setting.*
- **Without Ground-Truth (in Generalization):** When GT video is unavailable, we only assess instruction adherence: GPT-4o directly analyzes the generated video and the instruction to output the **Sequence Match Score** and the **Execution Quality Score**. *For this setting, only these two metrics are reported.*

**Physical Consistency and Causal Reasoning.** To quantify the physical plausibility of generated videos, we compute three kinds of metrics as follows:
- **Mask-guided Regional Consistency.** To better disentangle inconsistencies in the background, robot arm, and manipulated object, we propose **Mask-guided Regional Consistency**. We first use the GroundedSAM2 [71] with human annotation to obtain masks for the robot arm, the manipulated object(s), and the background in each frame. We then compute region-specific embeddings using a vision foundation model (e.g., DINOv3 [76]) and measure cosine similarity across time for each region separately. This allows us to pinpoint the source of temporal flaws—for instance, identifying a "jittery" robot arm even when the object and background are stable.
- **Trajectory Consistency:** For comparing the trajectory between generated videos and ground-truth counterparts,

5

we track both the end-effector and object trajectories. In our benchmark, we leverage SAM2 [70], given a few representative key points in the initial frame that humans annotated, to follow the motion of objects in both videos. Trajectory similarity is then evaluated using a complementary set of metrics: **Mean Euclidean Distance (MED)** [22] to capture average deviation, **Dynamic Time Warping (DTW)** [62] to assess temporal alignment, and **Fréchet Distance** [25] to measure worst-case path similarity.

- **Physical common sense:** Physical common sense covers dimensions ranging from object interaction and properties to temporal consistency, lighting, fluid dynamics, and local anomalies. To automatically score these six distinct dimensions, we fine-tuned Qwen-2.5-VL [2] using a two-stage GRPO [75] process. First, we used several physical and temporal datasets [5, 18, 19, 44, 55, 96] for an initial GRPO fine-tuning to enhance the model's understanding of video content, temporal causality, and physical laws. Second, we used a human-annotated dataset of 1,297 ratings (a mix of synthetic and real videos) for further training, aligning the model to human preferences. The fine-tuned model was finally used to score these dimensions on a 1-to-5 scale.

**Planning and Task Decomposition.** To evaluate long-horizon planning, we refer to the metric of RoboBench [57] based on Directed Acyclic Graphs (DAGs). We first parse the natural language instruction and ground-truth video into a ground-truth plan DAG, where nodes are atomic actions parameterized by $\langle \text{skill}, \text{object}, \text{args} \rangle$ and edges represent dependencies. This representation flexibly handles non-unique but valid action orderings. We then compare the model-generated plan (which also uses the same approach from the video) to the ground-truth DAG using two scores:

1. **Node Correctness:** The fraction of correctly predicted nodes aligns with the ground-truth nodes.
2. **Task Completion:** Using the MLLM to conduct a lightweight world-simulation rollout using: (i) first-frame image, (ii) reference action list, and (iii) reference DAG, and record the stage changes that model-predicted DAGs have correctly executed.

The final planning score **LongHorizon** integrates these aspects to reward both correctness and completeness:

$$S_{\text{LongHorizon}} = (S_{\text{NodeCorrectness}} + S_{\text{TaskCompletion}}) * 50.$$

### 3.4. Overall Benchmark Score

**Setup.** For each model $i$ and metric $m$, we map the raw measurement $x_{i,m}$ to a common desirability score $s_{i,m} \in (0, 100)$ via a monotone parametric mapping applied after an absolute pre-scaling to $[0, 1]$. We then aggregate desirability scores by weighted arithmetic means at both metric-group and overall levels.

**Pre-scale to** $[0, 1]$ **with absolute anchors.** Let $L_m < U_m$ be fixed anchors for metric $m$ (documented per-metric). Define the clipping operator $\text{clip}(u; a, b) = \min\{\max\{u, a\}, b\}$. We first clamp raw values to $[L_m, U_m]$, then linearly map to $[0, 1]$; for "higher-is-better" (HIB) metrics:

$$\hat{x}_{i,m}^{\text{HIB}} = \frac{\text{clip}(x_{i,m}; L_m, U_m) - L_m}{U_m - L_m} \in [0, 1],$$

And for "lower-is-better" (LIB) metrics:

$$\hat{x}_{i,m}^{\text{LIB}} = 1 - \hat{x}_{i,m}^{\text{HIB}} \in [0, 1].$$

We use absolute anchors for two common metrics:

$$\textbf{PSNR (HIB):} \quad L_{\text{PSNR}} = 0, \ U_{\text{PSNR}} = 50$$

$$\textbf{FVD (LIB):} \quad L_{\text{FVD}} = 0, \ U_{\text{FVD}} = 2000 \ .$$

For other metrics, $(L_m, U_m)$ are fixed per protocol (e.g., theoretical bounds for bounded scales, task-specific absolute targets for unbounded ones).

**Monotone parametric mappings.** After pre-scaling, we apply a single-parameter monotone transform $f_m(\cdot; \theta_m)$ and then rescale to $(0, 100)$:

$$s_{i,m} = 100 \, f_m(\hat{x}_{i,m}; \theta_m), \qquad s_{i,m} \in (0, 100).$$

We consider the following families (all are strictly increasing on $[0, 1]$):

$$\textbf{Simple:} \quad f_(x) = x,$$

$$\textbf{Power (Gamma):} \quad f_\gamma(x) = x^\gamma, \ \gamma > 0,$$

$$\textbf{Logit temperature:} \quad f_T(x) = \sigma(\text{logit}(x)/T), \ T > 0,$$
$$\sigma(t) = \tfrac{1}{1+e^{-t}},$$

$$\textbf{Tanh slope:} \quad f_\kappa(x) = \tfrac{1}{2}(\tanh(\kappa(2x - 1)) + 1),$$
$$\kappa > 0.$$

Particularly detailed information can be found in the Appendix 9.6 .

**Aggregation (weighted arithmetic mean).** Metrics are grouped into categories $g$ (e.g., *quality*, *instruction*, *physical*, *planning*). For each model $i$, let $\mathcal{M}_g$ denote the set of metrics available in group $g$, and $|\mathcal{M}_g| = N_{i,g} > 0$ be its cardinality. We aggregate first within each group by averaging and then across groups with normalized weights.

Let nonnegative group weights $\{W_g\}$ sum to one over the groups available for model $i$. The overall score is

$$O_i = \sum_{g \in \mathcal{G}_i} \frac{W_g}{\sum_{h \in \mathcal{G}_i} W_h} \left( \frac{1}{|\mathcal{M}_g|} \sum_{m \in \mathcal{M}_g} s_{i,m} \right).$$

where $\mathcal{G}_i = \{g : N_{i,g} > 0\}$ is the set of groups for which model $i$ has at least one valid metric. Setting $W_g \equiv 1$ yields an unweighted overall mean across all available groups.

6

Table 2. **Comparative analysis of different video generation models in Video Quality, Instruction Understanding and Planning Reasoning.** All metrics: higher is better. Best results are **bold** with green highlight. Seq. stands for Sequence; Exec. stands for Execution.

| Model | Video Quality | | | | | | Instruction Understanding | | | | Planning Reasoning |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | FVD | PSNR | SSIM | DINO | DreamSim | Overall | Caption Score | Seq. Match Score | Exec. Quality Score | Overall | Planning DAG |
| **Closed-Source Model** | | | | | | | | | | | |
| **Kling** | 80.41 | 5.41 | 85.00 | 36.95 | 58.43 | 53.24 | 81.91 | 0.25 | 1.75 | 27.97 | 2.50 |
| **Hailuo** | 85.01 | **6.06** | **85.97** | **41.45** | 61.95 | **56.09** | **88.78** | 51.07 | **70.48** | 70.11 | **17.27** |
| **Open-Source Model** | | | | | | | | | | | |
| **CogVideoX** | 50.49 | 2.48 | 80.48 | 17.06 | 42.09 | 38.52 | 79.66 | 38.44 | 44.18 | 54.09 | 4.55 |
| **Cosmos-Predict1** | 62.99 | 2.91 | 79.86 | 15.62 | 33.92 | 39.06 | 83.60 | 50.89 | 49.89 | 61.46 | 8.10 |
| **Wan2.1** | 62.95 | 1.85 | 74.09 | 20.50 | 41.76 | 40.23 | 82.28 | 41.43 | 46.84 | 56.85 | 7.95 |
| **Cosmos-Predict2** | 78.77 | 2.27 | 75.69 | 24.20 | 53.12 | 46.81 | 84.16 | 55.51 | 30.72 | 56.80 | 13.41 |
| **WoW-cosmos1** | 83.03 | 4.88 | 83.22 | 21.89 | 53.71 | 49.35 | 85.34 | **63.33** | 60.38 | 69.68 | 5.45 |
| **WoW-wan** | 84.89 | 5.95 | 85.76 | 37.74 | 62.57 | 55.38 | 84.85 | 52.45 | 49.18 | 62.16 | 9.32 |
| **WoW-cosmos2** | **85.71** | 3.52 | 78.33 | 38.63 | **64.42** | 54.12 | 85.85 | 59.04 | 66.18 | **70.36** | 12.27 |

Table 3. **Comparative analysis of different video generation models in Physical Law.** Con. stands for Consistency; Traj. stands for Trajectory; Obj. stands for Object; Cam. stands for Camera.

| Model | Physical Law | | | | | | | | | | | | | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Robot Con. | Obj. Con. | Scene Con. | Robot Traj. L2norm | Robot Traj. DTW | Robot Traj. FD | Obj. Traj. L2norm | Obj. Traj. DTW | Obj. Traj. FD | Physical Score | Cam. ATE | Cam. RPE | Overall | |
| **Closed-Source Model** | | | | | | | | | | | | | | |
| **Kling** | **57.62** | 16.09 | 89.26 | 56.08 | 86.05 | 25.27 | 87.14 | 87.31 | 41.77 | 71.72 | 98.12 | 99.86 | **68.02** | 37.93 |
| **Hailuo** | 56.99 | 10.08 | **89.95** | 56.79 | 86.53 | 22.79 | **87.87** | **88.96** | **44.90** | 57.64 | 98.26 | 99.86 | 66.72 | **52.55** |
| **Open-Source Model** | | | | | | | | | | | | | | |
| **CogVideoX** | 59.17 | 5.82 | 82.12 | 47.67 | 81.31 | 16.62 | 83.69 | 83.14 | 32.16 | 70.51 | 97.55 | 99.85 | 63.30 | 40.12 |
| **Cosmos-Predict1** | 39.02 | 4.40 | 71.69 | 54.87 | 84.60 | 23.55 | 83.93 | 83.48 | 36.99 | 30.21 | 96.23 | 99.66 | 59.05 | 41.93 |
| **Wan2.1** | 38.57 | 4.72 | 78.38 | 46.61 | 80.51 | 14.57 | 81.91 | 81.98 | 30.71 | **74.98** | 83.69 | 99.36 | 59.66 | 41.18 |
| **Cosmos-Predict2** | 34.61 | 5.83 | 68.34 | 53.15 | 84.99 | 22.24 | 84.82 | 84.92 | 40.74 | 48.78 | 98.41 | **99.87** | 60.56 | 44.40 |
| **WoW-cosmos1** | 44.72 | 13.46 | 77.78 | 56.00 | 86.77 | 25.10 | 86.00 | 86.49 | 37.86 | 37.10 | 96.42 | 99.66 | 62.28 | 46.70 |
| **WoW-wan** | 50.27 | 9.83 | 87.10 | 53.61 | 85.15 | 21.19 | 85.84 | 86.46 | 40.59 | 47.96 | 97.21 | 99.73 | 63.75 | 47.66 |
| **WoW-cosmos2** | 49.01 | **16.76** | 85.20 | **59.97** | **87.56** | **30.53** | 86.93 | 87.06 | 44.46 | 48.17 | **98.64** | **99.87** | 66.18 | 50.74 |

## 4. Experiments

**Models.** We evaluate our benchmark across a diverse set of commercial closed-source, open-source video generation models, and the World Foundation Model using direct Image-Text-to-Video generation. The closed-source models include Kling-2.1 [45] and Hailuo I2V-02 [60], two state-of-the-art proprietary systems optimized for zero-shot text-to-video and image-to-video generation. For open-source baselines, we include CogVideoX1.5-I2V-5B [94], one of the earliest large-scale video diffusion models, and Wan2.1-I2V-14B [83], an advanced text-to-video generator emphasizing temporal coherence. Furthermore, we evaluate the World Foundation Model, Cosmos-Predict1-7B [1] and Cosmos-Predict2-2B [64], which focus on physics, scene, and world modeling under multimodal prompts. Finally, we evaluate the Embodied World Foundation Model called WoW [17], which was trained for physically consistent and instruction-aware video generation in the embodied domain on different backbones. Together, these models cover a broad spectrum of architectures and training paradigms, enabling a comprehensive comparison across diverse dimensions in our benchmark.

### 4.1. Quantitative Evaluation Results

As shown in Figure 1 (Top-left and center), each ability dimension is evaluated by a targeted subset of metrics. We benchmark core dimensions with four categories of metrics—Video Quality, Instruction Understanding, Physical Law, and Planning Reasoning—with all scores normalized to a 0–100 scale. Results are summarized in Tables 2 and 3.

**Overall Performance.** Overall score balanced the four capabilities. Under this score, the best closed-source model, Hailuo, achieves the highest score (52.55). Among open-source models, WoW-cosmos2 consistently ranks highest, reaching 50.74, with balanced improvements in Video Quality, Instruction Understanding, and Physics Law, whereas earlier baselines, such as CogVideoX, lag in stability and physical consistency. The importance of overall balance is further highlighted by Kling, which is strong in Video Quality and achieves the best Physical Law score, yet attains a low score (37.93) due to severe deficiencies in Instruction Understanding and Planning. In conclusion, these four dimensions are indispensable for embodied world model benchmarks.

**Video Quality.** In Table 2, on the Video Quality axis, Hailuo obtains the highest closed-source overall score (56.09), while WoW-wan leads the open-source group, achieving a nearly matched performance (55.38). WoW-cosmos2 attains the best DreamSim (64.42) and the strongest FVD (85.71), even slightly exceeding Hailuo. These results suggest that WoW mainly focuses on high-level perceptual and distributional realism. For commercial models, they gain an excellent balance between low-level fidelity and high-level perceptual and distributional realism.

**Instruction Understanding.** For Instruction Understanding, WoW-cosmos2 achieves the best overall score among open-source models (70.36), surpassing the strongest closed-source model, Hailuo (70.11). When decomposing the metrics, we observe complementary strengths across sub-dimensions. Hailuo excels at Caption Score (88.78) and Exec. Quality (70.48), while WoW-cosmos1 shows the best Seq. Match (63.33), highlighting that both of them benefit from large-scale text–video alignment for instruction comprehension. Especially, WoW-cosmos2 better enforces procedural consistency—i.e., correctly ordering and realizing multi-step instructions over time—while remaining competitive in execution quality.

**Physical Law.** Table 3 shows that closed-source models establish the upper bound on overall score, with Kling and Hailuo scoring 68.02 and 66.72, respectively. Among open-source models, WoW-cosmos2 leads open-source models (66.18), with strong robot trajectory and camera-motion consistency. Firstly, from consistency perspective, the differentiation in Scene Consistency is not as intense as other two, which means all models performed well across this consistency. Secondly, WoW-cosmos2 attains the strongest Robot Trajectory and Hailuo attains the strongest Object Trajectory, indicates that commercial models focus more on object manipulation, while WoW-cosmos2 is more accurate in handling robot arm movements. Thirdly, camera motion is basically saturated; that is, the camera is very stable across all models, and sudden, erratic movements are rare. The last, these results all show that open-source models closely track commercial models in physical realism.

**Planning Reasoning.** Planning Reasoning remains a primary bottleneck for current world models. As shown in Table 2, Planning DAG scores occupy a markedly lower and more compressed range than others. The best performance is achieved by Hailuo (17.27), whereas Cosmos-Predict2 unexpectedly leads the open-source group (13.41), followed by WoW-cosmos2 (12.27) and WoW-wan (9.32). Models such as CogVideoX and Kling exhibit weak long-term planning, indicating that long-horizon planning and temporal reasoning still represent underexplored and require more explicit planning representations or control.

## 4.2. Dense prompts quantitative results

Prior work [26] has shown that the textual prompt has a substantial impact on video generation quality, particularly for models that rely on language-conditioned. In our setting, Image-Text-to-Video generation—the initial frame is fixed and cannot be altered, making the text prompt the primary controllable input. Our main results, Table 2 and 3, therefore, use concise prompts that specify only the task goal, without detailed descriptions of the scene, or step-by-step execution. Such under-specified prompts may limit the model's ability to ground the intended task semantics.

Table 4. **Comparative analysis of different video generation models with dense prompts.** All metrics: higher is better. Best results are **bold** with green highlight.

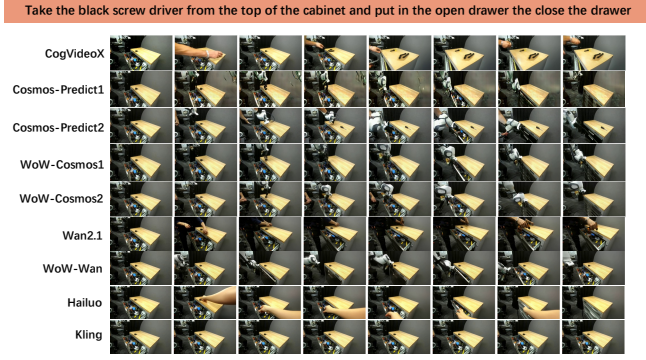| Model | VQ | IU | PL | PR | Overall |
|---|---|---|---|---|---|
| **Cosmos-Predict1** | 35.48 | 61.07 | 53.78 | 7.5 | 39.46 |
| **Cosmos-Predict2** | 49.75 | 75.96 | 64.66 | **13.86** | 51.06 |
| **WoW-cosmos1** | 59.41 | 72.54 | **69.71** | 11.14 | **53.20** |
| **WoW-wan** | **60.55** | 50.83 | 67.48 | 8.64 | 46.88 |
| **WoW-cosmos2** | 56.86 | **76.16** | 67.15 | 9.09 | 52.32 |

To study this effect, we use InternVL3-78B [101] to expand each short prompt into a dense prompt containing explicit environment and object references, sub-goals, and physical constraints. We then re-evaluate models on **WoW-World-Eval** using dense prompts in the above Table 4.

Across all evaluation dimensions, dense prompts lead to consistent and sometimes substantial improvements. In **Video Quality**, every model benefits from richer textual conditioning, with the best embodied world model surpassing the closed-source commercial baseline (60.55 vs. 56.09). **Instruction Understanding** exhibits even larger gains: nearly all models improve, and Cosmos-Predict2 increases by almost 20, although WoW-cosmos2 remains the most stable performer across both prompt types. Improvements in **Physical Law** are similarly widespread, with WoW-cosmos1 showing the largest jump (+7), surpassing the previous short-prompt SOTA. In contrast, in **Planning Reasoning**, DAG metric shows minimal change, indicating that detailed prompting alone cannot compensate for the current limitations of world models in long-horizon reasoning and structured task decomposition. Overall, dense prompts provide richer semantic and physical cues that meaningfully enhance generation quality, while exposing fundamental gaps in planning capability that require deeper architectural advances beyond prompt refinement.
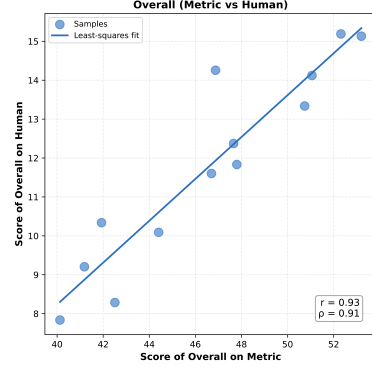
## 4.3. Qualitative Evaluation and Human Evaluation

**Qualitative Evaluation.** Figure 3a presents a qualitative comparison across models on a representative manipulation instruction, revealing failure patterns consistent with our quantitative results. General video generators such as CogVideoX, Cosmos-Predict1, and Wan2.1 exhibit visual artifacts, unintended camera motion, and even hallucinated human arms, reflecting dataset biases toward human demonstrations.

In terms of instruction understanding and planning, Cosmos-Predict1, WoW-Cosmos1, and WoW-Cosmos2 correctly execute the multi-step task, whereas CogVideoX and Wan2.1 fail despite producing motion. Cosmos-Predict2 and WoW-Wan show physics violations such as object duplication or motion without force, while Kling

(a) Qualitative Evaluation Example.



(b) Correlation between WoW-World-Eval Overall Score and Human Eval.

Figure 3. **Side-by-side comparison:** (a) qualitative results by different models; (b) display a metric–preference correlation.

remains static and ignores the task, demonstrating limited planning capability.

Overall, the videos generated by WoW- Cosmos2 is more likely to be indistinguishable from the real world in all aspects. These results highlight the need to evaluate models across all four dimensions—Video Quality, Instruction Understanding, Physical Law, and Planning Reasoning—as each captures a distinct failure mode in embodied video generation. Additional qualitative examples are included in the Appendix 11.

**Human Evaluation.** To assess the validity of our evaluation metric, we conduct a human study with 15 domain experts, who rated over 1,200 videos, both real and generated, along four core dimensions that mirror our benchmark and examine the alignment between **WoW-World-Eval** and human judgment. As shown in Figure 3b, both Pearson(r) and Spearman ($\rho$) correlations are reported for the Overall Score and Human Eval Score.

Results show a strong consistency between metric-based and human evaluations. The Overall Score demonstrates very strong Pearson and Spearman correlation with human judgment (r = 0.93, $\rho$ = 0.91), confirming that our benchmark provides a reliable and effective human-aligned evaluation of video generation performance.

### 4.4. Turing Test

To demonstrate that our benchmark can serve as a principled foundation for Turing-test evaluations—and that performance on the benchmark is strongly correlated with outcomes in such Turing tests—while also addressing two remaining questions about embodied world models, we design and implement two complementary Turing Tests: *Human Turing Test* and *Inverse Dynamic Model Turing Test*.

**Human Turing Test.** Figure 4 examines how **WoW-World-Eval** scores relate to the Deceive Human Ratio, measured by asking 13 participants to distinguish real videos from generated ones, which shows that models with

higher overall scores consistently achieve higher deceive-human ratios, indicating stronger perceptual realism. For example, Hailuo and the WoW-cosmos2 variants are the most effective at fooling humans, consistent with their strong performance on our benchmark. This trend is statistically validated: the deceive-human ratio is strongly correlated with our Overall score (r = 0.679), meaning models that perform well on the benchmark are also more likely to fool humans. Breaking down by dimensions, Video Quality (r = 0.874) and Physical Law (r = 0.753) exhibit the highest correlations, showing that realism in both appearance and physical dynamics is what most convinces humans.

These results reveal a core insight: both visual quality and physical plausibility are essential for generated videos to appear real to humans. Videos lacking either fidelity or correct physics fail to deceive observers. This strong alignment with human perception underscores that quality and physics are indispensable components of an embodied world model benchmark.
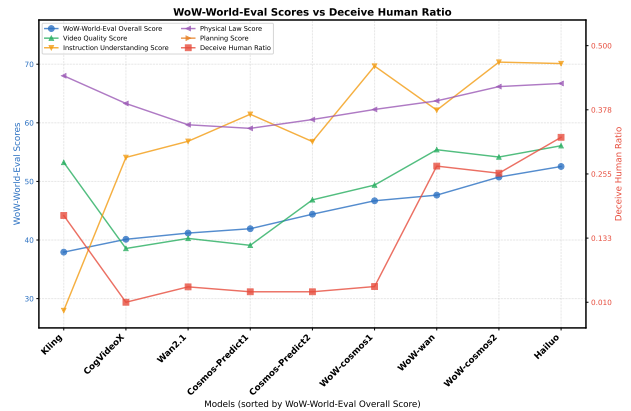


Figure 4. **Trending between WoW-World-Eval Overall Score and Deceive Human Ratio.**

**Inverse Dynamic Model Turing Test.** We further introduce a Machine Turing Test using a real-world–trained Inverse Dynamics Model (IDM), GC-IDM [17], to evaluate whether generated videos exhibit physically executable dynamics. We collect 9 different manipulation tasks, from easy to hard, for evaluation. Details are in Appendix 10.3

As shown in Table 5, most high-scoring models on **WoW-World-Eval** still fail this test: Kling (9.88%), Hailuo (2.47%), and early open-source models all perform near 0%, revealing that visual realism alone is insufficient for embodied execution. Models with real-robot data transfer significantly better—WoW-wan (40.74%) and WoW-cosmos2 (18.52%)—consistent with their strong Physical Law and Instruction Understanding scores.

These results highlight that: (i) The success rate of the generated videos on real robots is also correlated with our benchmark scores. (ii) Passing the IDM Turing Test requires both physics-grounded modeling and real-world exposure, emphasizing that true embodied competence goes beyond appearance quality. (iii) The videos generated by current embodied world models still have a significant gap from the real physical world.

Table 5. **Execution Accuracy on real world by Inverse Dynamic Model.** All metrics: higher is better. Best results are **bold** with green highlight. Succ. Rate stands for Real-world success rate

| Model | Real-world Succ. Rate (%) |
|---|---|
| **Kling** | 9.88 |
| **Hailuo** | 2.47 |
| **CogVideoX** | 0.00 |
| **Cosmos-Predict1** | 0.00 |
| **Wan2.1** | 0.00 |
| **Cosmos-Predict2** | 8.64 |
| **WoW-wan** | **40.74** |
| **WoW-cosmos2** | 18.52 |

## 5. Conclusion

In this work, we present **WoW-World-Eval**, a comprehensive benchmark designed as a Turing Test for evaluating video foundation models in embodied AI. By combining fine-grained human preference assessments and the Inverse Dynamics Model (IDM), we provide a dual-perspective evaluation framework that examines both perceptual realism and physical executability. Through extensive experiments on 609 real-world robot manipulation data, we reveal that while some models achieve high perceptual fidelity, they struggle with long-horizon reasoning and real-world execution, often collapsing under the IDM Turing Test. These findings highlight a substantial gap between visual generation and embodied action, pointing to the urgent need

for more physically grounded, generalizable world models for robotics. We hope our benchmark will serve as a foundation for future progress in building robust, general-purpose embodied intelligence.

## 6. Acknowledge

## 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. 7, 1

[2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. 6

[3] Philip J. Ball, Jakob Bauer, Frank Belletti, Bethanie Brownfield, Ariel Ephrat, Shlomi Fruchter, Agrim Gupta, Kristian Holsheimer, Aleksander Holynski, Jiri Hron, Christos Kaplanis, Marjorie Limont, Matt McGill, Yanko Oliveira, Jack Parker-Holder, Frank Perbet, Guy Scully, Jeremy Shar, Stephen Spencer, Omer Tov, Ruben Villegas, Emma Wang, Jessica Yung, Cip Baetu, Jordi Berbel, David Bridson, Jake Bruce, Gavin Buttimore, Sarah Chakera, Bilva Chandra, Paul Collins, Alex Cullum, Bogdan Damoc, Vibha Dasagi, Maxime Gazeau, Charles Gbadamosi, Woohyun Han, Ed Hirst, Ashyana Kachra, Lucie Kerley, Kristian Kjems, Eva Knoepfel, Vika Koriakin, Jessica Lo, Cong Lu, Zeb Mehring, Alex Moufarek, Henna Nandwani, Valeria Oliveira, Fabio Pardo, Jane Park, Andrew Pierson, Ben Poole, Helen Ran, Tim Salimans, Manuel Sanchez, Igor Saprykin, Amy Shen, Sailesh Sidhwani, Duncan Smith, Joe Stanton, Hamish Tomlinson, Dimple Vijaykumar, Luyu Wang, Piers Wingfield, Nat Wong, Keyang Xu, Christopher Yew, Nick Young, Vadim Zubov, Douglas Eck, Dumitru Erhan, Koray Kavukcuoglu, Demis Hassabis, Zoubin Gharamani, Raia Hadsell, Aäron van den Oord, Inbar Mosseri, Adrian Bolton, Satinder Singh, and Tim Rocktäschel. Genie 3: A new frontier for world models. 2025. 1

[4] Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chenfanfu Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. Videophy: Evaluating physical commonsense for video generation. *arXiv preprint arXiv:2406.03520*, 2024. 3, 4

[5] Hritik Bansal, Clark Peng, Yonatan Bitton, Roman Goldenberg, Aditya Grover, and Kai-Wei Chang. Videophy-2: A challenging action-centric physical commonsense evalua-

tion in video generation. *arXiv preprint arXiv:2503.06800*, 2025. 6

[6] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 1

[7] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, and OpenAI. Video generation models as world simulators. Technical report, OpenAI, 2024. Sora. 4, 1

[8] Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024. 1

[9] Jun Cen, Chaohui Yu, Hangjie Yuan, Yuming Jiang, Siteng Huang, Jiayan Guo, Xin Li, Yibing Song, Hao Luo, Fan Wang, et al. Worldvla: Towards autoregressive action world model. *arXiv preprint arXiv:2506.21539*, 2025. 1

[10] Yi Chen, Yuying Ge, Yixiao Ge, Mingyu Ding, Bohao Li, Rui Wang, Ruifeng Xu, Ying Shan, and Xihui Liu. Egoplan-bench: Benchmarking multimodal large language models for human-level planning, 2024. 4

[11] Zixuan Chen, Jing Huo, Yangtao Chen, and Yang Gao. Robohorizon: An llm-assisted multi-view world model for long-horizon robotic manipulation, 2025. 1

[12] Sijie Cheng, Kechen Fang, Yangyang Yu, Sicheng Zhou, Bohao Li, Ye Tian, Tingguang Li, Lei Han, and Yang Liu. Videgothink: Assessing egocentric video understanding capabilities for embodied ai, 2024. 4

[13] Sijie Cheng, Zhicheng Guo, Jingwen Wu, Kechen Fang, Peng Li, Huaping Liu, and Yang Liu. Egothink: Evaluating first-person perspective thinking capability of vision-language models, 2024. 4

[14] Zhili Cheng, Yuge Tu, Ran Li, Shiqi Dai, Jinyi Hu, Shengding Hu, Jiahao Li, Yang Shi, Tianyu Yu, Weize Chen, Lei Shi, and Maosong Sun. Embodiedeval: Evaluate multimodal llms as embodied agents, 2025. 4

[15] Xiaowei Chi, Chun-Kai Fan, Hengyuan Zhang, Xingqun Qi, Rongyu Zhang, Anthony Chen, Chi-min Chan, Wei Xue, Qifeng Liu, Shanghang Zhang, et al. Eva: An embodied world model for future video anticipation. *arXiv preprint arXiv:2410.15461*, 2024. 1, 3, 4

[16] Xiaowei Chi, Kuangzhi Ge, Jiaming Liu, Siyuan Zhou, Peidong Jia, Zichen He, Yuzhen Liu, Tingguang Li, Lei Han, Sirui Han, Shanghang Zhang, and Yike Guo. Mind: Learning a dual-system world model for real-time planning and implicit risk analysis, 2025. 1

[17] Xiaowei Chi, Peidong Jia, Chun-Kai Fan, Xiaozhu Ju, Weishi Mi, Kevin Zhang, Zhiyuan Qin, Wanxin Tian, Kuangzhi Ge, Hao Li, et al. Wow: Towards a world omniscient world model through embodied interaction. *arXiv preprint arXiv:2509.22642*, 2025. 1, 4, 7, 10, 2

[18] Wei Chow, Jiageng Mao, Boyi Li, Daniel Seita, Vitor Guizilini, and Yue Wang. Physbench: Benchmarking and enhancing vision-language models for physical world understanding. *arXiv preprint arXiv:2501.16411*, 2025. 3, 4, 6

[19] Daniel Cores, Michael Dorkenwald, Manuel Mucientes, Cees G. M. Snoek, and Yuki M. Asano. Lost in time: A new temporal benchmark for videollms, 2025. 6

[20] Sudeep Dasari, Frederik Ebert, Stephen Tian, Suraj Nair, Bernadette Bucher, Karl Schmeckpeper, Siddharth Singh, Sergey Levine, and Chelsea Finn. Robonet: Large-scale multi-robot learning. *arXiv preprint arXiv:1910.11215*, 2019. 1

[21] Jingtao Ding, Yunke Zhang, Yu Shang, Yuheng Zhang, Zefang Zong, Jie Feng, Yuan Yuan, Hongyuan Su, Nian Li, Nicholas Sukiennik, et al. Understanding world or predicting future? a comprehensive survey of world models. *ACM Computing Surveys*, 58(3):1–38, 2025. 1

[22] Ivan Dokmanic, Reza Parhizkar, Juri Ranieri, and Martin Vetterli. Euclidean distance matrices: Essential theory, algorithms, and applications. *IEEE Signal Processing Magazine*, 32(6):12–30, 2015. 6

[23] Mengfei Du, Binhao Wu, Zejun Li, Xuanjing Huang, and Zhongyu Wei. Embspatial-bench: Benchmarking spatial understanding for embodied tasks with large vision-language models, 2024. 4

[24] Haoyi Duan, Hong-Xing Yu, Sirui Chen, Li Fei-Fei, and Jiajun Wu. Worldscore: A unified evaluation benchmark for world generation, 2025. 4

[25] Thomas Eiter and Heikki Mannila. Computing discrete fréchet distance. In *Technical Report CD-TR 94/64, Christian Doppler Laboratory for Expert . . .* , 1994. 6

[26] Tiehan Fan, Kepan Nan, Rui Xie, Penghao Zhou, Zhenheng Yang, Chaoyou Fu, Xiang Li, Jian Yang, and Ying Tai. Instancecap: Improving text-to-video generation via instance-aware structured caption. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 28974–28983, 2025. 8

[27] Fernando A. Fardo, Victor H. Conforto, Francisco C. de Oliveira, and Paulo S. Rodrigues. A formal evaluation of psnr as quality measurement parameter for image segmentation algorithms, 2016. 2, 5

[28] Gustav Theodor Fechner. *Elemente der psychophysik*. Breitkopf u. Härtel, 1860. 3

[29] Weixi Feng, Jiachen Li, Michael Saxon, Tsu-jui Fu, Wenhu Chen, and William Yang Wang. Tc-bench: Benchmarking temporal compositionality in text-to-video and image-to-video generation. *arXiv preprint arXiv:2406.08656*, 2024. 1, 3

[30] Stephanie Fu, Netanel Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. Dreamsim: Learning new dimensions of human visual similarity using synthetic data, 2023. 5

[31] Xiao Fu, Xintao Wang, Xian Liu, Jianhong Bai, Runsen Xu, Pengfei Wan, Di Zhang, and Dahua Lin. Learning video generation for robotic manipulation with collaborative trajectory control. *arXiv preprint arXiv:2506.01943*, 2025. 1

[32] Pascale Fung, Yoram Bachrach, Asli Celikyilmaz, Kamalika Chaudhuri, Delong Chen, Willy Chung, Emmanuel Dupoux, Hongyu Gong, Hervé Jégou, Alessandro Lazaric, et al. Embodied ai agents: Modeling the world. *arXiv preprint arXiv:2506.22355*, 2025. 1

[33] James J Gibson. The theory of affordances:(1979). In *The people, place, and space reader*, pages 56–60. Routledge, 2014. 4

[34] Yanjiang Guo, Lucy Xiaoyang Shi, Jianyu Chen, and Chelsea Finn. Ctrl-world: A controllable generative world model for robot manipulation. *arXiv preprint arXiv:2510.10125*, 2025. 1

[35] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2(3), 2018. 1

[36] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023. 1

[37] Zihao He, Bo Ai, Tongzhou Mu, Yulin Liu, Weikang Wan, Jiawei Fu, Yilun Du, Henrik I Christensen, and Hao Su. Scaling cross-embodiment world models for dexterous manipulation. *arXiv preprint arXiv:2511.01177*, 2025. 1

[38] Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shotton, and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving. *arXiv preprint arXiv:2309.17080*, 2023. 1

[39] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 5

[40] Pengliang Ji, Chuyang Xiao, Huilin Tai, and Mingxiao Huo. T2vbench: Benchmarking temporal dynamics for text-to-video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5325–5335, 2024. 1, 3

[41] Yuxin Jiang, Shengcong Chen, Siyuan Huang, Liliang Chen, Pengfei Zhou, Yue Liao, Xindong He, Chiming Liu, Hongsheng Li, Maoqing Yao, et al. Enerverse-ac: Envisioning embodied environments with action condition. *arXiv preprint arXiv:2505.09723*, 2025. 1

[42] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, Peter David Fagan, Joey Hejna, Masha Itkina, Marion Lepert, Yecheng Jason Ma, Patrick Tree Miller, Jimmy Wu, Suneel Belkhale, Shivin Dass, Huy Ha, Arhan Jain, Abraham Lee, Youngwoon Lee, Marius Memmel, Sungjae Park, Ilija Radosavovic, Kaiyuan Wang, Albert Zhan, Kevin Black, Cheng Chi, Kyle Hatch, Shan Lin, Jingpei Lu, Jean Mercat, Abdul Rehman, Pannag R. Sanketi, Archit Sharma, Cody Simpson, Quan Vuong, Homer Rich Walke, Blake Wulfe, Ted Xiao, Jonathan Heewon Yang, Arefeh Yavary, Tony Z. Zhao, Christopher Agia, Rohan Baijal, Mateo Guaman Castro, Daphne Chen, Qiuyu Chen, Trinity Chung, Jaimyn Drake, Ethan Paul Foster, Jensen Gao, David Antonio Herrera, Minho Heo, Kyle Hsu, Jiaheng Hu, Donovon Jackson, Charlotte Le, Yunshuang Li, Kevin Lin, Roy Lin, Zehan Ma, Abhiram Maddukuri, Suvir Mirchandani, Daniel Morton, Tony Nguyen, Abigail O'Neill, and Rosario Scalise. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*, 2024. 5

[43] Po-Chen Ko, Jiayuan Mao, Yilun Du, Shao-Hua Sun, and Joshua B Tenenbaum. Learning to Act from Actionless Videos through Dense Correspondences. *arXiv:2310.08576*, 2023. 1, 10

[44] Benno Krojer, Mojtaba Komeili, Candace Ross, Quentin Garrido, Koustuv Sinha, Nicolas Ballas, and Mahmoud Assran. A shortcut-aware video-qa benchmark for physical understanding via minimal video pairs, 2025. 6

[45] Kuaishou. Kling, 2025. https://app.klingai.com/cn/. 7, 1

[46] 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, 3, 4

[47] Hengtao Li, Pengxiang Ding, Runze Suo, Yihao Wang, Zirui Ge, Dongyuan Zang, Kexian Yu, Mingyang Sun, Hongyin Zhang, Donglin Wang, et al. Vla-rft: Vision-language-action reinforcement fine-tuning with verified rewards in world simulators. *arXiv preprint arXiv:2510.00406*, 2025. 1

[48] Jinming Li, Yichen Zhu, Zhiyuan Xu, Jindong Gu, Minjie Zhu, Xin Liu, Ning Liu, Yaxin Peng, Feifei Feng, and Jian Tang. Mmro: Are multimodal llms eligible as the brain for in-home robotics?, 2024. 4

[49] Xinqing Li, Xin He, Le Zhang, and Yun Liu. A comprehensive survey on world models for embodied ai. *arXiv preprint arXiv:2510.16732*, 2025. 1

[50] Ying Li, Xiaobao Wei, Xiaowei Chi, Yuming Li, Zhongyu Zhao, Hao Wang, Ningning Ma, Ming Lu, and Shanghang Zhang. Manipdreamer: Boosting robotic manipulation world model with action tree and visual guidance. *arXiv preprint arXiv:2504.16464*, 2025. 1

[51] Ying Li, Xiaobao Wei, Xiaowei Chi, Yuming Li, Zhongyu Zhao, Hao Wang, Ningning Ma, Ming Lu, and Shanghang Zhang. Manipdreamer3d: Synthesizing plausible robotic manipulation video with occupancy-aware 3d trajectory. *arXiv preprint arXiv:2509.05314*, 2025. 1

[52] Yue Liao, Pengfei Zhou, Siyuan Huang, Donglin Yang, Shengcong Chen, Yuxin Jiang, Yue Hu, Jingbin Cai, Si Liu, Jianlan Luo, Liliang Chen, Shuicheng Yan, Maoqing Yao, and Guanghui Ren. Genie envisioner: A unified world foundation platform for robotic manipulation, 2025. 1

[53] Xinran Ling, Chen Zhu, Meiqi Wu, Hangyu Li, Xiaokun Feng, Cundian Yang, Aiming Hao, Jiashu Zhu, Jiahong Wu, and Xiangxiang Chu. Vmbench: A benchmark for perception-aligned video motion generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13087–13098, 2025. 1, 3

[54] Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu, Tieyong Zeng, Raymond Chan, and Ying Shan. Evalcrafter: Benchmarking and evaluating large video generation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22139–22149, 2024. 1, 3

[55] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Tempcompass: Do video llms really understand videos?, 2024. 6

12

[56] Xiaoxiao Long, Qingrui Zhao, Kaiwen Zhang, Zihao Zhang, Dingrui Wang, Yumeng Liu, Zhengjie Shu, Yi Lu, Shouzheng Wang, Xinzhe Wei, et al. A survey: Learning embodied intelligence from physical simulators and world models. *arXiv preprint arXiv:2507.00917*, 2025. 1

[57] Yulin Luo, Chun-Kai Fan, Menghang Dong, Jiayu Shi, Mengdi Zhao, Bo-Wen Zhang, Cheng Chi, Jiaming Liu, Gaole Dai, Rongyu Zhang, et al. Robobench: A comprehensive evaluation benchmark for multimodal large language models as embodied brain. *arXiv preprint arXiv:2510.17801*, 2025. 4, 6

[58] Arjun Majumdar, Anurag Ajay, Xiaohan Zhang, Pranav Putta, Sriram Yenamandra, Mikael Henaff, Sneha Silwal, Paul Mcvay, Oleksandr Maksymets, Sergio Arnaud, Karmesh Yadav, Qiyang Li, Ben Newman, Mohit Sharma, Vincent Berges, Shiqi Zhang, Pulkit Agrawal, Yonatan Bisk, Dhruv Batra, Mrinal Kalakrishnan, Franziska Meier, Chris Paxton, Alexander Sax, and Aravind Rajeswaran. Openeqa: Embodied question answering in the era of foundation models. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16488–16498, 2024. 4

[59] Fanqing Meng, Jiaqi Liao, Xinyu Tan, Wenqi Shao, Quanfeng Lu, Kaipeng Zhang, Yu Cheng, Dianqi Li, Yu Qiao, and Ping Luo. Towards world simulator: Crafting physical commonsense-based benchmark for video generation. *arXiv preprint arXiv:2410.05363*, 2024. 3, 4

[60] MiniMax. Hailuo, 2025. https://hailuoai.video/. 7, 1

[61] Saman Motamed, Laura Culp, Kevin Swersky, Priyank Jaini, and Robert Geirhos. Do generative video models understand physical principles? *arXiv preprint arXiv:2501.09038*, 2025. 3, 4

[62] Meinard Müller. *Dynamic Time Warping*, pages 69–84. Springer Berlin Heidelberg, Berlin, Heidelberg, 2007. 6

[63] Soroush Nasiriany, Sean Kirmani, Tianli Ding, Laura Smith, Yuke Zhu, Danny Driess, Dorsa Sadigh, and Ted Xiao. Rt-affordance: Affordances are versatile intermediate representations for robot manipulation. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8249–8257. IEEE, 2025. 4

[64] Nvidia. Cosmos-predict2, 2025. https://research.nvidia.com/labs/dir/cosmos-predict2/. 7, 1

[65] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, and Julien Mairal. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2, 5

[66] Jing-Cheng Pang, Nan Tang, Kaiyuan Li, Yuting Tang, Xin-Qiang Cai, Zhen-Yu Zhang, Gang Niu, Masashi Sugiyama, and Yang Yu. Learning view-invariant world models for visual robotic manipulation. In *The Thirteenth International Conference on Learning Representations*, 2025. 1

[67] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 1

[68] Yiran Qin, Zhelun Shi, Jiwen Yu, Xijun Wang, Enshen Zhou, Lijun Li, Zhenfei Yin, Xihui Liu, Lu Sheng, Jing Shao, Lei Bai, Wanli Ouyang, and Ruimao Zhang. Worldsimbench: Towards video generation models as world simulators, 2024. 4

[69] Julian Quevedo, Percy Liang, and Sherry Yang. Evaluating robot policies in a world model. *arXiv preprint arXiv:2506.00613*, 2025. 1

[70] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 6

[71] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang 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. 5

[72] Unitree Robotics. Unifolm-wma-0: A world-model–action (wma) framework under the unifolm family. https://unigen-x.github.io/unifolm-world-model-action.github.io/, 2025. Accessed: YYYY-MM-DD. 1

[73] Younggyo Seo, Junsu Kim, Stephen James, Kimin Lee, Jinwoo Shin, and Pieter Abbeel. Multi-view masked world models for visual robotic manipulation. In *Proceedings of the 40th International Conference on Machine Learning*, pages 30613–30632. PMLR, 2023. 1

[74] Pierre Sermanet, Tianli Ding, Jeffrey Zhao, Fei Xia, Debidatta Dwibedi, Keerthana Gopalakrishnan, Christine Chan, Gabriel Dulac-Arnold, Sharath Maddineni, Nikhil J Joshi, Pete Florence, Wei Han, Robert Baruch, Yao Lu, Suvir Mirchandani, Peng Xu, Pannag Sanketi, Karol Hausman, Izhak Shafran, Brian Ichter, and Yuan Cao. Robovqa: Multimodal long-horizon reasoning for robotics, 2023. 4

[75] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. 6

[76] Oriane Siméoni, Huy V. Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, Francisco Massa, Daniel Haziza, Luca Wehrstedt, Jianyuan Wang, Timothée Darcet, Théo Moutakanni, Leonel Sentana, Claire Roberts, Andrea Vedaldi, Jamie Tolan, John Brandt, Camille Couprie, Julien Mairal, Hervé Jégou, Patrick Labatut, and Piotr Bojanowski. Dinov3, 2025. 5

[77] Chan Hee Song, Valts Blukis, Jonathan Tremblay, Stephen Tyree, Yu Su, and Stan Birchfield. Robospatial: Teaching spatial understanding to 2d and 3d vision-language models for robotics, 2025. 4

[78] Kaiyue Sun, Kaiyi Huang, Xian Liu, Yue Wu, Zihan Xu, Zhenguo Li, and Xihui Liu. T2v-compbench: A comprehensive benchmark for compositional text-to-video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 8406–8416, 2025. 3

[79] HunyuanWorld Team, Zhenwei Wang, Yuhao Liu, Junta Wu, Zixiao Gu, Haoyuan Wang, Xuhui Zuo, Tianyu Huang, Wenhuan Li, Sheng Zhang, et al. Hunyuanworld 1.0: Generating immersive, explorable, and interactive 3d worlds from words or pixels. *arXiv preprint arXiv:2507.21809*, 2025. 1

[80] Wanxin Tian, Shijie Zhang, Kevin Zhang, Xiaowei Chi, Yulin Luo, Junyu Lu, Chunkai Fan, Qiang Zhou, Yiming Zhao, Ning Liu Siyu Lin, et al. Seea-r1: Tree-structured reinforcement fine-tuning for self-evolving embodied agents. *arXiv preprint arXiv:2506.21669*, 2025. 4

[81] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018. 2, 5

[82] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017. 1

[83] 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. 7, 1

[84] Boyang Wang, Nikhil Sridhar, Chao Feng, Mark Van der Merwe, Adam Fishman, Nima Fazeli, and Jeong Joon Park. This&that: Language-gesture controlled video generation for robot planning. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 12842–12849. IEEE, 2025. 1

[85] Zixing Wang and Ahmed H Qureshi. Anypose: Anytime 3d human pose forecasting via neural ordinary differential equations. *arXiv preprint arXiv:2309.04840*, 2023. 1

[86] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4): 600–612, 2004. 2, 5

[87] Kun Wu, Chengkai Hou, Jiaming Liu, Zhengping Che, Xiaozhu Ju, Zhuqin Yang, Meng Li, Yinuo Zhao, Zhiyuan Xu, Guang Yang, Shichao Fan, Xinhua Wang, Fei Liao, Zhen Zhao, Guangyu Li, Zhao Jin, Lecheng Wang, Jilei Mao, Ning Liu, Pei Ren, Qiang Zhang, Yaoxu Lyu, Mengzhen Liu, Jingyang He, Yulin Luo, Zeyu Gao, Chenxuan Li, Chenyang Gu, Yankai Fu, Di Wu, Xingyu Wang, Sixiang Chen, Zhenyu Wang, Pengju An, Siyuan Qian, Shanghang Zhang, and Jian Tang. Robomind: Benchmark on multiembodiment intelligence normative data for robot manipulation. In *Proceedings of Robotics: Science and Systems (RSS) 2025*, 2025. 5

[88] Junjin Xiao, Yandan Yang, Xinyuan Chang, Ronghan Chen, Feng Xiong, Mu Xu, Wei-Shi Zheng, and Qing Zhang. World-env: Leveraging world model as a virtual environment for vla post-training. *arXiv preprint arXiv:2509.24948*, 2025. 1

[89] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021. 1

[90] Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in space: How multimodal large language models see, remember, and recall spaces, 2025. 4

[91] Rui Yang, Hanyang Chen, Junyu Zhang, Mark Zhao, Cheng Qian, Kangrui Wang, Qineng Wang, Teja Venkat Koripella, Marziyeh Movahedi, Manling Li, Heng Ji, Huan Zhang, and Tong Zhang. Embodiedbench: Comprehensive benchmarking multi-modal large language models for vision-driven embodied agents, 2025. 4

[92] Xiuyu Yang, Bohan Li, Shaocong Xu, Nan Wang, Chongjie Ye, Zhaoxi Chen, Minghan Qin, Yikang Ding, Xin Jin, Hang Zhao, et al. Orv: 4d occupancy-centric robot video generation. *arXiv preprint arXiv:2506.03079*, 2025. 1

[93] Ze Yang, Yun Chen, Jingkang Wang, Sivabalan Manivasagam, Wei-Chiu Ma, Anqi Joyce Yang, and Raquel Urtasun. Unisim: A neural closed-loop sensor simulator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1389–1399, 2023. 1

[94] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 7, 1

[95] Hu Yue, Siyuan Huang, Yue Liao, Shengcong Chen, Pengfei Zhou, Liliang Chen, Maoqing Yao, and Guanghui Ren. Ewmbench: Evaluating scene, motion, and semantic quality in embodied world models. *arXiv preprint arXiv:2505.09694*, 2025. 1, 3, 4

[96] Yuan Zhang, Fei Xiao, Tao Huang, Chun-Kai Fan, Hongyuan Dong, Jiawen Li, Jiacong Wang, Kuan Cheng, Shanghang Zhang, and Haoyuan Guo. Unveiling the tapestry of consistency in large vision-language models. *Advances in Neural Information Processing Systems*, 37: 118632–118653, 2024. 6

[97] Dian Zheng, Ziqi Huang, Hongbo Liu, Kai Zou, Yinan He, Fan Zhang, Lulu Gu, Yuanhan Zhang, Jingwen He, Wei-Shi Zheng, et al. Vbench-2.0: Advancing video generation benchmark suite for intrinsic faithfulness. *arXiv preprint arXiv:2503.21755*, 2025. 1, 3, 4

[98] Hongyan Zhi, Peihao Chen, Siyuan Zhou, Yubo Dong, Quanxi Wu, Lei Han, and Mingkui Tan. 3dflowaction: Learning cross-embodiment manipulation from 3d flow world model. *arXiv preprint arXiv:2506.06199*, 2025. 1

[99] Siyuan Zhou, Yilun Du, Jiaben Chen, Yandong Li, Dit-Yan Yeung, and Chuang Gan. Robodreamer: Learning compositional world models for robot imagination. *arXiv preprint arXiv:2404.12377*, 2024. 1

[100] Fangqi Zhu, Hongtao Wu, Song Guo, Yuxiao Liu, Chilam Cheang, and Tao Kong. Irasim: Learning interactive real-robot action simulators. *arXiv preprint arXiv:2406.14540*, 2024. 1

[101] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 8

# Wow, wo, val! A Comprehensive Embodied World Model Evaluation Turing Test

## Supplementary Material

## 7. More Related Works

**Video Generation Models.** The quest for scalable video synthesis has converged on two dominant architectures, both serving as backbones for predictive world modeling.

**Autoregressive Transformers.** Inspired by Large Language Models, architectures like VideoGPT [89] treat video generation as a discrete sequence modeling task. By tokenizing frames via VQ-VAE [82], these models predict visual token sequence-by-sequence. This paradigm, akin to next-word prediction, naturally aligns with the state-transition logic of world models. Notable examples like Genie [8] utilize this spatiotemporal autoregressive objective to build interactive environments and predictive models from diverse video data. However, the discrete quantization often imposes an upper bound on visual fidelity, resulting in the loss of high-frequency details.

**The Diffusion Paradigm.** Diffusion Models (DMs) have established a new state-of-the-art by formulating generation as iterative denoising. Early attempts like SVD [6] adapted 2D U-Nets with temporal layers, excelling at texture but often treating time as a secondary dimension. The recent shift to Diffusion Transformers (DiT) [67], exemplified by Wan [83], Cosmos [64], Sora [7], operates on unified spacetime patches. These models demonstrate emergent capabilities akin to a physics engine—simulating fluid dynamics and occlusion—making them ideal candidates for high-fidelity world simulation.

**Embodied World Models.** Unlike passive video generation, embodied world models function as predictive engines $(s_{t+1}|s_t, a_t)$ to facilitate agent planning.

**From Latent States to Pixel Space.** Early research balanced between efficiency and detail. The Dreamer series [36] prioritized computational tractability by learning dynamics in a compact latent space, enabling efficient RL planning but producing schematic visual reconstructions. Conversely, pixel-space approaches like RoboNet [20] predicted future frames directly from interaction logs. While capturing texture, these models often struggled with the stochastic nature of real-world physics, leading to blurry predictions over long horizons.

**Generalist Simulators.** Recent efforts converge generative AI with robotics to create "Generalist Simulators." Foundation models like UniSim [93] and GAIA-1 [38] leverage internet-scale data to learn high-level physics for zero-shot generalization.

## 8. Models Detail

Since the models have different versions and parameter sizes, and also generate videos of different durations and resolutions, we introduce the models we use, their specific versions and parameter counts, as well as the generated durations and resolutions, to facilitate comparison and reproduction.

**Kling. [45]** A large-scale diffusion-transformer video generator known for photorealistic long-form outputs, strong motion control, and realistic physics. It employs a latent-space DiT architecture with 3D VAE compression and is optimized for high-consistency content such as sports or human motion. We used the latest Klingv2.1 to generate a 5-second 720p video.

**Hailuo. [60]** Hailuo is a commercial text-to-video and image-to-video model designed for high-quality short-form generation (6–10 seconds, up to 1080p). It emphasizes photorealism, smooth motion, and subject consistency, making it a strong closed-source baseline for perceptual quality. We use the publicly accessible Hailuo-02 for 5-second 768p video generation and evaluation.

**CogVideoX. [94]** An open diffusion-transformer video generator available in 2B/5B variants. It can produce 10-second clips with improved motion expressiveness using a 3D causal VAE and MoE-style transformer backbone. The model supports both text-to-video and image-to-video generation and serves as a strong open-source baseline for long-duration video synthesis. For evaluation, we use CogVideoX-I2V-5B to generate 5-second 1360x768 videos.

**Cosmos-Predict. [1, 64]** NVIDIA's open-source world-model family for video-based state prediction. It supports Text-/Image-/Video-to-World forecasting and can generate long-horizon, physics-aware video rollouts. The model is trained on large-scale real-world videos and emphasizes dynamics, object permanence, and 3D geometry consistency. We evaluate both Cosmos-Predict1 and the stronger Cosmos-Predict2 variants. Both Predict1&2 are 5-second 720p, but the parameters of the model are 7B and 2B, respectively.

**Wan. [83]** Wan is a series of open-source Diffusion Transformer video models (e.g., Wan 2.1/2.2) trained on

large-scale image–video corpora. They support text-to-video and image-to-video generation, producing 5–10 second clips with strong spatial–temporal coherence. Wan models use 3D VAE compression and DiT backbones to improve motion consistency and are widely used as high-quality open video generation baselines. We fix the resolution for a 5-second 832x480 generation with Wan2.1-I2V-14B.

**WoW. [17]** WoW is a 14B-parameter open-source embodied world model designed for physically grounded video prediction in robotic manipulation. It generates future video rollouts conditioned on visual observations and task instructions, with inductive biases (e.g., DINOv2 token distillation) that encourage spatial and physical consistency. Unlike generic video models, WoW explicitly targets planning, contact dynamics, and multi-step task execution. We use the different backbones of WoW, including Wan, Cosmos-Predict1, and Cosmos-Predict2, for evaluation; all generated video settings are the same as the backbone models.

## 9. Detailed Metrics

### 9.1. Visual Fidelity.

To evaluate the visual fidelity of generated videos, we adopt a suite of complementary metrics that capture low-level reconstruction accuracy, perceptual similarity, and distributional realism. These include FVD, SSIM, PSNR, DINO, and DreamSim.

**Fréchet Video Distance (FVD).** FVD is a distribution-level metric that measures how closely generated video features match those of real videos. Given feature means and covariances $(\mu_r, \Sigma_r)$ for real videos and $(\mu_g, \Sigma_g)$ for generated videos extracted using an I3D network, FVD is defined as:

$$\text{FVD} = \|\mu_r - \mu_g\|_2^2 + \text{Tr}\left(\Sigma_r + \Sigma_g - 2\left(\Sigma_r\Sigma_g\right)^{1/2}\right).$$

Lower FVD indicates higher spatial-temporal coherence and distributional realism.

**Structural Similarity Index measures (SSIM).** The Structural Similarity Index (SSIM) measures perceptual degradation in luminance, contrast, and structural information. For two frames $x$ and $y$:

$$\text{SSIM}(x,y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)},$$

where $\mu$ and $\sigma$ denote means and variances of local patches. Higher SSIM corresponds to better perceptual consistency.

**Peak signal-to-noise ratio (PSNR).** PSNR evaluates pixel-level fidelity via the mean squared error (MSE) between reference and generated frames:

$$\text{PSNR}(x,y) = 10\log_{10}\left(\frac{MAX^2}{\text{MSE}(x,y)}\right).$$

Higher PSNR generally indicates better reconstruction, although it is less aligned with human perception than SSIM or embedding-based metrics.

**DINO Score.** DINOv2 is a self-supervised visual foundation model trained with large-scale teacher–student distillation. We compute the cosine similarity between DINOv2 embeddings of generated and reference frames:

$$\text{DINO}(g_t, r_t) = \frac{\langle f(g_t), f(r_t)\rangle}{\|f(g_t)\|_2\,\|f(r_t)\|_2},$$

where $f(\cdot)$ denotes the DINOv2 encoder. Higher similarity indicates closer semantic and structural alignment.

**DreamSim.** DreamSim combines CLIP, OpenCLIP, and DINO embeddings and is fine-tuned on human perceptual judgments to produce similarity scores aligned with human perception. For frames $x$ and $y$:

$$\text{DreamSim}(x,y) = 1 - \|E(x) - E(y)\|_2,$$

where $E(\cdot)$ is the fused embedding.

### 9.2. Instruction Semantic Correctness.

These three scores—Caption Score, Sequence Match Score, and Execution Quality Score—allow us to evaluate instruction following in both supervised (with ground truth) and fully OOD settings. Empirically, all three metrics show strong agreement with human preference judgments, validating the reliability of our evaluation protocol.

**Caption Score.** For samples paired with ground-truth videos, we obtain a structured semantic representation of each video. Using a templated prompt, GPT extracts a standardized caption containing: initial state, intermediate process, final state, object descriptions, and action descriptions. The same procedure is applied to the model-generated video. A vision–language model (VLM) compares the two captions along these five predefined components, assigning a score of 1 (match), 0.5 (partial match), or 0 (mismatch). The Caption Score is computed as the mean over the five components.

**Sequence Match Score.** In this score, we directly analyze the alignment between the instruction and the actions depicted in the generated video. GPT identifies the sequence
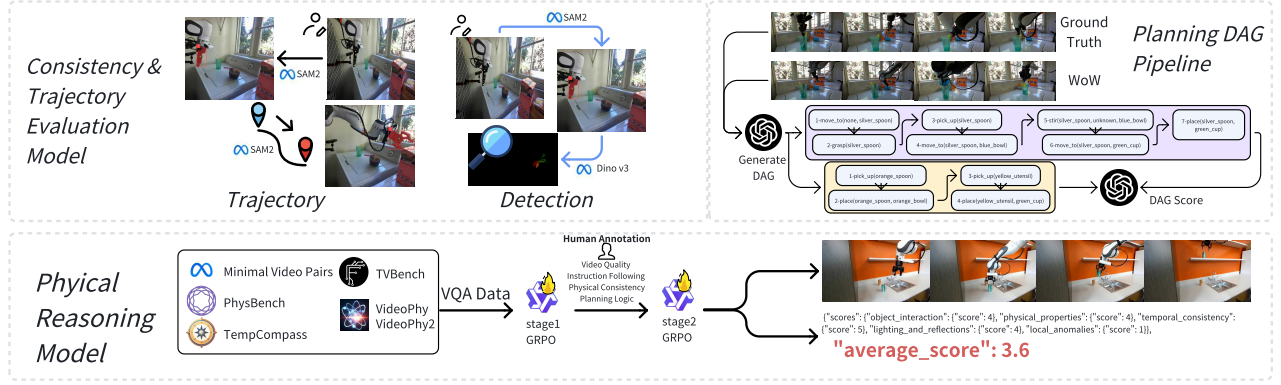
Figure 5. **Overview of WoW-World-Eval Metrics.**

of action–object pairs present in the video and aligns them with the intended execution order specified by the instruction. We measure the degree to which the extracted sequence preserves the correctness of both action selection and temporal ordering, again using a discrete scale of 1 (match), 0.5 (partial match), or 0 (mismatch). This yields the Sequence Match Score.

**Execution Quality Score.** To further evaluate execution fidelity, we assess how completely the depicted actions are carried out. We employ a five-level rubric measuring the extent of object motion and action completion—from no meaningful change to full task completion.
1. object does not move, and the action is not executed
2. object slightly moves, or the action is partially executed
3. object slightly moves, and the action is partially executed
4. object reaches goal or action is fully executed
5. object reaches goal and action is fully executed

The score ranges from 1 to 5, producing the Execution Quality Score. This score captures the degree to which the video evidences a coherent and physically plausible execution of the instructed task.

### 9.3. Mask-guided Regional Consistency.

**Setup.** We evaluate the mask-guided regional consistency of each model based on three core components, as illustrated in Figure 5: (1) *GroundingSAM-2* for point-guided segmentation and multi-frame tracking; (2) *DINOv3-Large (standard)* as a frozen visual encoder to extract patch-level features; (3) *Run-length Encoding (RLE)* for storing region masks efficiently across all frames.

**Implementation.** We begin by manually identifying the first frame in which both the manipulated object and the robot gripper are clearly visible. A human annotator

then places *3–5 boundary points* along the visible contour of each region (object and gripper). These point sets provide fine-grained localization cues and are passed to GroundingSAM-2, which returns segmentation masks for the first frame and tracks them throughout the video. For each frame $t$, GroundingSAM-2 outputs an RLE mask for the object and an RLE mask for the gripper. If a region cannot be tracked in frame $t$ (due to occlusion or motion blur), we record its mask as an *all-zero RLE mask*. The background region is defined as the complement of the union of the two tracked regions:

$$M_t^{\mathrm{bg}} = \mathbf{1} - \left( M_t^{\mathrm{obj}} \vee M_t^{\mathrm{arm}} \right).$$

Each RGB frame $I_t$ is encoded by DINOv3-Large to obtain a grid of patch features

$$F_t \in \mathbb{R}^{H_p \times W_p \times d}.$$

All masks are decoded from RLE and downsampled to $(H_p, W_p)$ so they can serve as per-patch weights. For each region $r \in \{\mathrm{obj}, \mathrm{arm}, \mathrm{bg}\}$ we compute a normalized mask:

$$w_t^{\mathrm{r}}(i,j) = \frac{M_t^{\mathrm{r}}(i,j)}{\sum_{i',j'} M_t^{\mathrm{r}}(i',j') + \epsilon},$$

which becomes all zero when the region is missing. A region feature is obtained via mask-weighted averaging:

$$\mathbf{f}_t^{\mathrm{r}} = \sum_{i,j} w_t^{\mathrm{r}}(i,j)\, F_t(i,j,:),$$

followed by $\ell_2$ normalization if $\mathbf{f}_t^{\mathrm{r}} \neq 0$.

**Consistency Score.** We compute temporal consistency for each region separately using cosine similarity. Given

3

normalized region features $\tilde{\mathbf{f}}_t^r$, the consistency between two frames $a$ and $b$ for region $r$ is

$$\text{Consist}^r(a, b) = \begin{cases} \tilde{\mathbf{f}}_a^r \cdot \tilde{\mathbf{f}}_b^r, & \|\tilde{\mathbf{f}}_a^r\|_2 > 0, \|\tilde{\mathbf{f}}_b^r\|_2 > 0, \\ 0, & \text{otherwise.} \end{cases}$$

For each frame $t \geq 2$ and region $r \in \{\text{obj}, \text{arm}, \text{bg}\}$, we combine long-range and short-range temporal coherence:

$$s_t^r = \frac{1}{2} \text{Consist}^r(1, t) + \frac{1}{2} \text{Consist}^r(t-1, t).$$

The video-level Mask-guided Regional Consistency for region $r$ is then

$$\text{MRC}^r = \frac{1}{T-1} \sum_{t=2}^{T} s_t^r,$$

where $T$ is the number of frames. In practice, we report three separate scores

$$\left( \text{MRC}^{\text{obj}}, \text{MRC}^{\text{arm}}, \text{MRC}^{\text{bg}} \right),$$

corresponding to the manipulated object, the gripper, and the background, respectively.

### 9.4. Trajectory-level Consistency.

We evaluate the trajectory-level consistency of each model by comparing the 2D motion trajectories of both the robot end-effector and manipulated objects in the generated videos against those in ground-truth (GT) videos. Human annotators first label representative keypoints on the robot arm and objects in the first frame; these serve as positive point prompts for SAM2 tracking, yielding dense and stable point tracks throughout the sequence, as illustrated in Figure 5. The per-frame trajectory is defined as the mean of the tracked keypoints, normalized to the range $[0, 1]^2$. We align the number of frames of the generated video and the GT video to the smaller one by uniformly sampling. We then compute three complementary metrics—L2Norm Error, Dynamic Time Warping (DTW), and Fréchet Distance—for both the robot end-effector and the manipulated object.

**Keypoint labeling and SAM2 tracking.** Given a video with frames $\{I_t\}_{t=1}^T$, human annotators place $N$ representative points on the robot end-effector or an object in the first frame:

$$\mathcal{K}_1 = \left\{ \mathbf{u}_k^{(1)} \in \mathbb{R}^2 \right\}_{k=1}^N, \qquad \mathbf{u}_k^{(1)} = (x_k^{(1)}, y_k^{(1)}).$$

**Keypoint prompting and SAM2 mask tracking.** Given a video with frames $\{I_t\}_{t=1}^T$, human annotators place $N$ representative points on the robot end-effector or an object

in the first frame. These points serve as positive prompts for SAM2, which returns a binary segmentation mask for the object at each frame:

$$\mathbf{M}^{(t)} \in \{0, 1\}^{H \times W}, \qquad t = 1, \ldots, T.$$

Each mask $\mathbf{M}^{(t)}$ defines the set of foreground pixels

$$\Omega^{(t)} = \left\{ (x, y) \mid \mathbf{M}^{(t)}(x, y) = 1 \right\}.$$

**Mask-to-point trajectory conversion.** To obtain a 2D trajectory from the mask sequence, we compute the centroid of the foreground region:

$$\mathbf{p}_t = (x_t, y_t) = \frac{1}{|\Omega^{(t)}|} \sum_{(x,y) \in \Omega^{(t)}} (x, y).$$

Let the video resolution be $(W, H)$. We then normalize all coordinates to $[0, 1]^2$:

$$\hat{\mathbf{p}}_t = \left( \frac{x_t}{W}, \frac{y_t}{H} \right), \qquad t = 1, \ldots, T.$$

This centroid simplification provides a stable 2D trajectory independent of resolution for each tracked component.

**Camera-motion–aware trajectory correction.** In generated videos with unstable viewpoints, the observed end-effector trajectory entangles true robot motion with camera drift or jitter, leading to biased evaluations. To recover the actual end-effector motion, we subtract the estimated camera trajectory $\mathbf{c}_t$ from the observed end-effector trajectory $\hat{\mathbf{p}}_t$:

$$\mathbf{p}_t^{\text{true}} = \hat{\mathbf{p}}_t - \mathbf{c}_t.$$

This correction yields a cleaner and more comparable motion signal across videos. Empirically, smaller camera metrics ATE/RPE correspond to more reliable corrected trajectories, whereas larger camera errors introduce evaluation noise and indicate the need for improved video generation or stabilization.

**Camera trajectory estimation.** Given a video, we uniformly sample frames $\{I_t\}_{t=1}^T$ and detect sparse Shi–Tomasi keypoints on the image boundary of the first frame

$$\mathcal{P}_1 = \left\{ \mathbf{u}_k^{(1)} \in \mathbb{R}^2 \right\}_{k=1}^N, \qquad \mathbf{u}_k^{(1)} = (x_k^{(1)}, y_k^{(1)}),$$

where boundary regions are assumed to be dominated by a static background. We then track these keypoints across subsequent frames with pyramidal Lucas–Kanade optical flow, obtaining correspondences between a reference frame and frame $t$:

$$\mathcal{P}_{\text{ref}} = \{\mathbf{u}_k^{\text{ref}}\}, \qquad \mathcal{P}_t = \{\mathbf{u}_k^{(t)}\}.$$

4

For each frame $t$, we robustly estimate a 2D affine transform (restricted to rotation/scale/translation) via RANSAC:

$$\mathbf{A}_t = \begin{bmatrix} a_{11} & a_{12} & t_x \\ a_{21} & a_{22} & t_y \end{bmatrix}, \quad \mathbf{u}_k^{(t)} \approx \mathbf{A}_t \begin{bmatrix} \mathbf{u}_k^{\text{ref}} \\ 1 \end{bmatrix},$$

where $\mathbf{t}_t = (t_x, t_y)$ captures the apparent translation of the static background in the image plane. Because camera motion is opposite to background motion, we define the per-frame camera displacement as

$$\Delta \mathbf{c}_t = - \mathbf{t}_t = (-t_x, -t_y).$$

Starting from $\mathbf{c}_1 = (0, 0)$, we accumulate these displacements over time, with an additional drift clipping to suppress implausibly large jumps between adjacent frames:

$$\mathbf{c}_t = \mathbf{c}_{t-1} + \Delta \mathbf{c}_t, \quad t = 2, \ldots, T,$$

where $\mathbf{c}_t = (x_t, y_t)$ denotes the cumulative camera offset (in pixels) at frame $t$ relative to the first frame. The original frame resolution $(W, H)$ is stored together with the trajectory.

**Temporal alignment and normalization across videos.** For a GT video and a generated video, we obtain two normalized trajectories:

$$\{\hat{\mathbf{p}}_t^{\text{gt}}\}_{t=1}^{T_{\text{gt}}}, \qquad \{\hat{\mathbf{p}}_t^{\text{gen}}\}_{t=1}^{T_{\text{gen}}}.$$

To compare two trajectories of different lengths, we first determine the target number of frames:

$$T = \min(T_{\text{gt}}, T_{\text{gen}}).$$

We uniformly sample the longer trajectory to obtain exactly $T$ frames. Let the aligned trajectories be

$$\mathbf{q}_t^{\text{gt}} = \hat{\mathbf{p}}_{s^1(t)}^{\text{gt}}, \qquad \mathbf{q}_t^{\text{gen}} = \hat{\mathbf{p}}_{s^2(t)}^{\text{gen}}, \qquad t = 1, \ldots, T.$$

where $s^1$ and $s^2$ are the sampled indexes.

**Absolute Trajectory Error (ATE).** ATE measures the absolute discrepancy between the generated and GT camera trajectories over the entire sequence. We define

$$\text{ATE} = \sqrt{\frac{1}{T} \sum_{t=1}^{T} \left\| \hat{\mathbf{c}}_t^{\text{gen}} - \hat{\mathbf{c}}_t^{\text{gt}} \right\|_2^2},$$

where $\| \cdot \|_2$ denotes the L2 norm. A lower ATE indicates that the overall camera path in the generated video more closely follows the ground-truth trajectory in the image plane.

**Relative Pose Error (RPE).** To evaluate the fidelity of *local* camera motion (e.g., instantaneous velocity and direction), we compute the Relative Pose Error. For a temporal offset $\Delta$ (we use $\Delta = 1$ frame in our experiments), we define normalized relative motion as

$$\mathbf{v}_t^{\text{gt}} = \hat{\mathbf{c}}_{t+\Delta}^{\text{gt}} - \hat{\mathbf{c}}_t^{\text{gt}}, \qquad \mathbf{v}_t^{\text{gen}} = \hat{\mathbf{c}}_{t+\Delta}^{\text{gen}} - \hat{\mathbf{c}}_t^{\text{gen}},$$

$$t = 1, \ldots, T - \Delta.$$

The RPE is then defined as

$$\text{RPE} = \sqrt{\frac{1}{T - \Delta} \sum_{t=1}^{T-\Delta} \left\| \mathbf{v}_t^{\text{gen}} - \mathbf{v}_t^{\text{gt}} \right\|_2^2}.$$

RPE focuses on frame-to-frame motion consistency: lower values indicate that the generated video better matches the ground-truth in terms of local camera dynamics (e.g., smoothness, acceleration patterns, and motion direction).

**L2Norm Error.** The per-frame L2Norm error measures instantaneous positional discrepancy:

$$\text{L2norm} = \sqrt{\frac{1}{T} \sum_{t=1}^{T} \left\| \mathbf{q}_t^{\text{gen}} - \mathbf{q}_t^{\text{gt}} \right\|_2^2}.$$

A lower ED indicates that the generated trajectory more closely follows the GT path frame by frame.

**Dynamic Time Warping (DTW).** DTW preserves the shape similarity of two trajectories while allowing non-linear temporal alignment. The DTW cost is defined as

$$\text{DTW}(\mathbf{q}^{\text{gen}}, \mathbf{q}^{\text{gt}}) = \min_{\pi} \sum_{(t,s) \in \pi} \left\| \mathbf{q}_t^{\text{gen}} - \mathbf{q}_s^{\text{gt}} \right\|_2,$$

where $\pi$ is a monotonic warping path. Lower DTW values indicate higher similarity in trajectory geometry, regardless of speed variations.

**Fréchet Distance.** The Fréchet Distance captures the minimum leash-length needed for two curves to be traversed in order, providing a strict measure of geometric similarity:

$$\text{FD}(\mathbf{q}^{\text{gen}}, \mathbf{q}^{\text{gt}}) = \inf_{\alpha, \beta} \max_{t \in [0,1]} \left\| \mathbf{q}_{\alpha(t)}^{\text{gen}} - \mathbf{q}_{\beta(t)}^{\text{gt}} \right\|_2,$$

where $\alpha$ and $\beta$ are continuous, non-decreasing reparameterizations. Unlike DTW, the Fréchet metric enforces simultaneous forward progression along the curves.

### 9.5. Physical and Causal Reasoning.

**Base Model and Optimization Algorithm**

To implement automated evaluation, we selected **Qwen-2.5-VL (7B)** [2] as our base vision-language model. We employed the **GRPO** algorithm, optimizing the model through a two-stage fine-tuning process to equip it with precise physical commonsense understanding and scoring capabilities, as illustrated in Figure 5.

**The GRPO Algorithm**

GRPO optimizes the policy model $\pi_\theta$ by maximizing an objective that encourages high-reward outputs while maintaining stability via a clipping mechanism and KL-divergence regularization.

For each input query $q$, GRPO samples a group of outputs $\{o_1, o_2, \cdots, o_G\}$ from the old policy $\pi_{\theta_{old}}$. The optimization objective is defined as:

$$
\begin{aligned}
\mathcal{J}_{GRPO}(\theta) = \mathbb{E}\Bigg[ \frac{1}{G}\sum_{i=1}^{G}\bigg( &\min\bigg[ \\
&\frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}\hat{A}_i, \\
&\text{clip}\left(\frac{\pi_\theta(o_i|q)}{\pi_{\theta_{old}}(o_i|q)}, 1-\varepsilon, 1+\varepsilon\right)\hat{A}_i\bigg] \\
&- \beta D_{KL}[\pi_\theta||\pi_{ref}]\bigg)\Bigg]
\end{aligned}
\tag{1}
$$

where $\varepsilon$ is the clipping hyperparameter, and $\beta$ controls the KL-divergence penalty with respect to the reference model $\pi_{ref}$.

The advantage $\hat{A}_i$ for the $i$-th output is calculated based on the relative rewards within the sampled group. Specifically, assuming output $o_i$ receives a reward score $r_i$, the advantage is computed by standardizing the rewards:

$$
\hat{A}_i = \frac{r_i - \text{mean}(\{r_1, \ldots, r_G\})}{\text{std}(\{r_1, \ldots, r_G\}) + \epsilon}
\tag{2}
$$

This group-relative formulation serves as a dynamic baseline, effectively reducing variance without requiring a separate value network.

**Stage 1: Foundational Video Understanding Fine-tuning**

**Objective:** To imbue the model with a foundational understanding of physical events and causal relationships. **Data and Method:** We utilized a dataset of approximately 50,000 samples from six video understanding benchmarks. All data was formatted as multiple-choice Video Question Answering (VQA) tasks.

- **Sampling:** For each video-question pair $q$, we sampled a group of outputs of size $G = 8$.

- **Reward Calculation:** We implemented a rule-based reward function. If a sampled output $o_i$ correctly matched the ground-truth option (e.g., "Option A"), it received a reward $r_i = 1.0$; otherwise, $r_i = 0.0$.
- **Optimization:** The model was updated using the calculated group advantages to increase the probability of generating correct answers.

**Results:** Stage 1 fine-tuning significantly improved the model's video comprehension abilities. On a comprehensive held-out test set from these benchmarks, the model's average accuracy increased from **60.83%** (the Qwen-2.5-VL 7B backbone) to **71.51%**.

**Stage 2: Scoring Alignment Fine-tuning**

**Objective:** To train the model to follow instructions and output quantitative 1-5 scores for four dimensions in a JSON format, aligning with human-annotated standards. **Method:** We continued to use the GRPO algorithm, fine-tuning on our internally annotated dataset of 1,297 data points. We used the following prompt template to guide the model in generating the scoring JSON. **Reward Function:** In Stage 2, the GRPO optimization objective was to maximize the expectation of a reward function $R(y_{out}, y_{gt})$, which quantifies the alignment between the model-generated JSON $y_{out}$ and the ground-truth human-annotated JSON $y_{gt}$.

The reward $R$ is calculated as follows:
1. Parse the JSON objects from $y_{out}$ and $y_{gt}$.
2. Let $\mathcal{K}$ be the set of four evaluation keys (i.e., `video_quality`, `instruction_following`, `physical_consistency`, `planning_logic`).
3. For each key $k \in \mathcal{K}$ present in $y_{gt}$, if $k$ is also present in $y_{out}$, calculate the normalized absolute error $e_k$.
4. Before calculating $e_k$, the scores $s_k^{gt}$ from $y_{gt}$ and $s_k^{out}$ from $y_{out}$ are clipped to the valid range $[1.0, 5.0]$, yielding $s_k'^{gt}$ and $s_k'^{out}$.
5. The normalized absolute error $e_k$ is defined as:

$$
e_k = \frac{|s_k'^{gt} - s_k'^{out}|}{4.0}
\tag{3}
$$

The denominator 4.0 (i.e., $5 - 1$) normalizes the error to the range $[0.0, 1.0]$.
6. The mean error $\bar{e}$ is calculated over all matched keys $\mathcal{K}_{match} \subseteq \mathcal{K}$:

$$
\bar{e} = \frac{1}{|\mathcal{K}_{match}|}\sum_{k \in \mathcal{K}_{match}} e_k
\tag{4}
$$

7. The final reward $R$ is defined as $1.0 - \bar{e}$ and clipped to $[0.0, 1.0]$:

$$
R(y_{out}, y_{gt}) = \max(0.0, \min(1.0, 1.0 - \bar{e}))
\tag{5}
$$

If $y_{out}$ or $y_{gt}$ are not valid JSON, or if $\mathcal{K}_{match}$ is empty, the reward $R = 0.0$.

**Final Inference for Evaluation**

After completing the two-stage fine-tuning, we obtained an evaluation model with a strong understanding of physical common sense (from S1) and the ability to perform structured scoring tasks (from S2). For the final automated evaluation in our paper, we employed a detailed, zero-shot inference prompt. This prompt instructs the model to focus *exclusively* on physical plausibility, refining its evaluation criteria into the six dimensions used in our main analysis.

## 9.6. Overall Benchmark Score.

We normalize all metrics to a unified three-digit scale of 0–100 to ensure consistent interpretation and cross-metric comparability. Metrics that naturally fall within this range are retained as-is. For the remaining metrics (including FVD and PSNR, we fix their boundaries up to 2000 and 50, respectively), we first map their values to [0,1], then apply a monotonic transformation, and finally scale the result to [0,100].

To determine the optimal mapping method and parameter for each metric, we use a subset of human evaluation scores as supervision. Specifically, we apply different values within a grid over [0,5], compute the Pearson correlation between the transformed metric scores and human ratings, and select the method and parameter that maximizes this correlation.

**Parametric mappings.** After pre-scaling, we apply a single-parameter monotone transform $f_m(\cdot; \theta_m)$ and then rescale to $(0, 100)$:

$$s_{i,m} = 100\, f_m(\hat{x}_{i,m}; \theta_m), \qquad s_{i,m} \in (0, 100).$$

We consider the following families (all are strictly increasing on $[0, 1]$):

$$\text{Simple:} \quad f_{(}x) = x,$$
$$\text{Power (Gamma):} \quad f_\gamma(x) = x^\gamma,\ \gamma > 0,$$
$$\text{Logit temperature:} \quad f_T(x) = \sigma(\text{logit}(x)/T),\ T > 0,$$
$$\sigma(t) = \tfrac{1}{1+e^{-t}},$$
$$\text{Tanh slope:} \quad f_\kappa(x) = \tfrac{1}{2}(\tanh(\kappa(2x - 1)) + 1),$$
$$\kappa > 0.$$

In practice, $\gamma > 1$ accentuates the high end, while $T < 1$ or $\kappa > 1$ expands the mid-range and compresses extremes. For numerical stability with logit we use a small $\varepsilon$ (e.g., $10^{-6}$) and replace $x$ by $\text{clip}(x; \varepsilon, 1 - \varepsilon)$ only inside $\text{logit}(\cdot)$.

**Parameter selection and freezing.** For each metric $m$, $\theta_m \in \{\gamma, T, \kappa\}$ is selected on a fixed development set by maximizing the Fisher-$z$ averaged Pearson correlation between $f_m(\hat{x}; \theta)$ and human ratings across $K$-fold CV; Spearman correlation is used as a tie-breaker. The chosen $\theta_m$ is then *frozen* and applied to all evaluations.

The resulting mapping methods and $\theta_m$ values for all metrics $m$ are summarized in Table 6.

Table 6. **The mapping parameters for all the metrics.**

| Metric | Mapping | Parameter |
|---|---|---|
| **FVD** | gamma | 1.52 |
| **PSNR** | tanh | 4.71 |
| **SSIM** | gamma | 0.61 |
| **DINO** | gamma | 3.06 |
| **DreamSim** | gamma | 2.94 |
| **Caption Score** | gamma | 0.12 |
| **Seq. Match Score** | gamma | 2.45 |
| **Exec. Qual Score** | gamma | 2.97 |
| **Planning DAG** | simple | - |
| **Robot Con.** | gamma | 2.93 |
| **Obj. Con.** | tanh | 4.93 |
| **Scene Con.** | gamma | 3.94 |
| **Robot Traj. L2norm** | gamma | 2.86 |
| **Robot Traj. DTW** | gamma | 1.87 |
| **Robot Traj. FD** | gamma | 4.00 |
| **Obj. Traj. L2norm** | gamma | 1.27 |
| **Obj. Traj. DTW** | gamma | 2.99 |
| **Obj. Traj. FD** | gamma | 3.52 |
| **Physical Score** | simple | - |
| **Cam. ATE** | simple | - |
| **Cam. RPE** | simple | - |

# 10. Extended Experiment Analysis

## 10.1. WoW-World-Eval Human Evaluation Rule

We evaluate generated videos along four independent dimensions: Video Quality, Instruction Understanding, Physical Law, and Planning. Each dimension is scored on a 1–5 scale, and the Overall Score is defined as the sum of the four scores (range: 4–20). All dimensions are assessed independently.

**Video Quality.** This dimension assesses the perceptual integrity of the video itself, independent of task correctness or physical feasibility. The focus is strictly on the clarity, stability, and visibility of the content presented. Scoring Criteria will be as follows:

**5 — Excellent**

- High spatial clarity with sharp and stable rendering
- Proper exposure and color balance; no notable artifacts or noise
- All key elements remain continuously visible
- No distracting visual irregularities

**4 — Good**

- Minor imperfections (e.g., slight blur, mild exposure fluctuation, low-level noise)
- Overall visibility remains unaffected

### 3 — Fair
- Multiple noticeable issues, including moderate blur, frequent autofocus/exposure changes, or significant noise
- Occasional hindrance to understanding the scene
- Overall content remains interpretable

### 2 — Poor
- Significant visual defects: strong blur, over/under-exposure, instability, occlusion of key regions
- Compression artifacts or frame drops are obvious and disruptive
- Core task content is frequently hard to interpret

### 1 — Unusable
- Severe degradation: extreme darkness/brightness, pervasive blur, heavy artifacts, or constant frame loss
- Critical task elements cannot be identified
- The video cannot support meaningful evaluation

**Instruction Understanding.** This dimension measures whether the final outcome satisfies the given instruction. It is strictly goal-oriented, disregarding process efficiency or motion plausibility. Scoring Criteria will be as follows:

### 5 — Fully Achieved
- The instruction is completed precisely and entirely
- All expected final conditions are satisfied
- No essential elements are missing

### 4 — Mostly Achieved
- The majority of the instructions is correctly executed
- Minor deviations exist, but do not compromise overall task success

### 3 — Partially Achieved
- Only some components of the instruction are fulfilled
- Significant omissions or inaccuracies remain
- The intended goal is only partially realized

### 2 — Minimally Achieved
- Most of the required outcome is not completed
- The behavior exhibits weak or inconsistent correspondence to the instruction

### 1 — Not Achieved
- The instruction is not fulfilled at all
- Actions are irrelevant or contradictory to the task
- A fully static video always receives a score of 1 in this dimension

**Physical Law.** This dimension evaluates whether the depicted motions and object interactions adhere to real-world physical principles, including dynamics, continuity, and plausible human-robot interaction. Scoring Criteria will be as follows:

### 5 — Fully Physical
- All motions are dynamically coherent and physically plausible
- Object interactions follow realistic physical responses
- No penetration artifacts or implausible transitions
- A completely static video (identical to the first frame) is also scored as 5, as it contains no physical violations

### 4 — Mostly Physical
- Minor deviations from realistic physics, but overall behavior remains plausible
- No significant physical inconsistencies

### 3 — Moderately Physical
- Noticeable but non-catastrophic physical inconsistencies
- Partial deviations from expected object dynamics
- Movements occasionally appear unnatural

### 2 — Poorly Physical
- Frequent or severe physics violations
- Object motion or robot kinematics contradict basic real-world principles
- Motion discontinuities or penetrations are common

### 1 — Physically Impossible
- Major violations of fundamental physics (e.g., teleportation, reverse gravity, object behavior without force causality)
- Substantial geometric inconsistencies (large penetrations)
- If a generated human performs part or all of the task, the Physical Law score is supposed to be 1, regardless of physical plausibility

**Planning Reasoning.** This dimension assesses the logical structure and coherence of the robot's action sequence, independent of final task success. It captures whether the robot demonstrates purposeful, ordered planning. Scoring Criteria will be as follows:

### 5 — Well-Structured Planning
- Action sequence is coherent, efficient, and well aligned with the task
- No extraneous or aimless motions
- Exhibits deliberate, task-driven progression

### 4 — Reasonable Planning
- Overall logical sequence with minor redundant actions
- Task progression remains clear and purposeful

### 3 — Weak Planning
- Action sequence contains irregularities or inefficiencies
- Task intent is visible, but execution lacks structure
- Multiple corrective attempts may be present

### 2 — Poor Planning
- Disorganized or inconsistent action sequence
- Repeated irrelevant behaviors or unnecessary back-and-forth motions
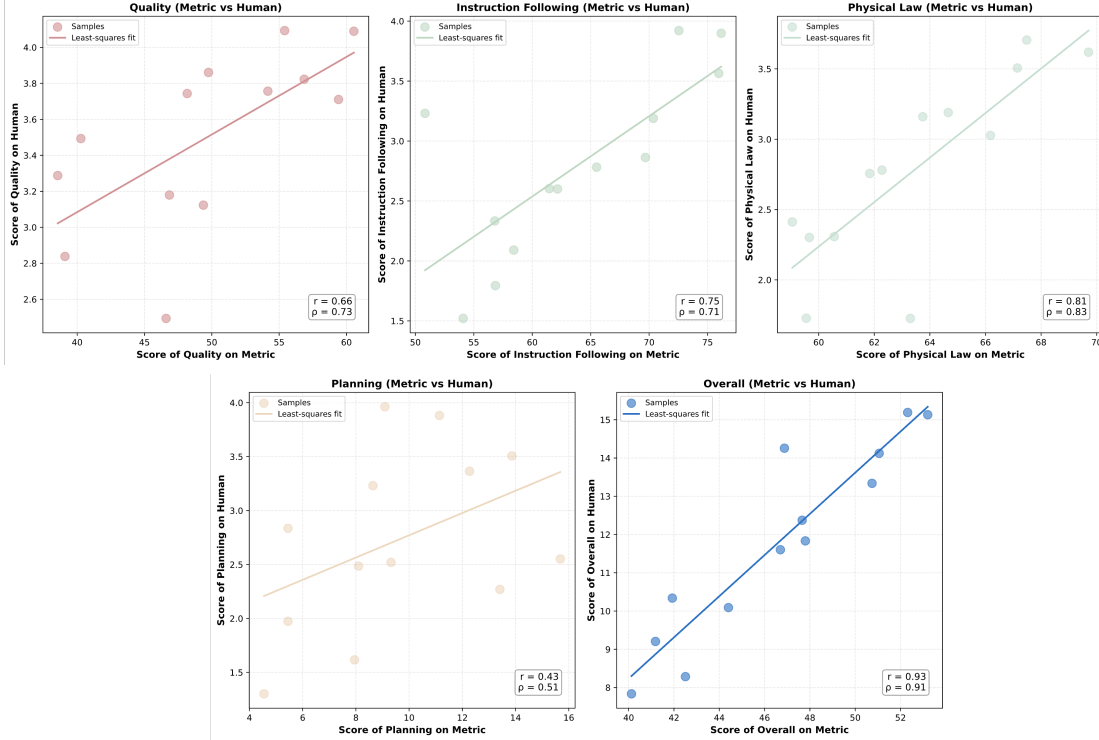- Weak alignment between steps and task objective

Figure 6. **All correlation for WoW-World-Eval and Human Preference.**

**1 — No Planning**
- Actions are random, chaotic, or entirely absent
- A fully static video always receives a score of 1 in this dimension

**Overall Score.** The Overall Score is defined as the sum of the four scores. All dimensions are weighted equally, supporting transparent and reproducible evaluation.

## 10.2. All correlation of WoW-World-Eval score and Human Preference

In section 4.1, we report the correlation between our overall scores and overall human preferences. However, to quantify how well WoW-World-Eval reflects human perception in all perspectives, we also provide the correlation between these four metrics—Video Quality, Instruction Understanding, Physical Law, and Planning. Figure 6 reports both strong Pearson (r) and Spearman ($\rho$) correlations between all groups of benchmark scores and human preferences, which, importantly, confirms that metric-based evaluation can reliably approximate human judgment, supporting the use of WoW-World-Eval as a scalable alternative to costly human evaluations.

**Video Quality.** Quality metrics exhibit a quite strong correlation with human assessments (r = 0.66, $\rho$ = 0.73). This

confirms that perceptual sharpness, temporal consistency, and high-level structure—captured jointly by FVD, PSNR, SSIM, DINO, and DreamSim—align well with how humans evaluate visual sensitivity.

**Instruction Understanding.** The Instruction Understanding score shows robust alignment with human preferences (r = 0.75, $\rho$ = 0.71). Models that accurately capture state transitions, object–action relationships, and multi-step task semantics tend to be judged by humans as better at "following instructions," validating our caption-, sequence- and execution-based IU metrics.

**Physical Law.** Correlation is highest for Physical Law (r = 0.81, $\rho$ = 0.83). This indicates that humans are highly sensitive to violations of physics—object irregular motion, penetration, trajectory discontinuities, implausible state changes—and that our physics metrics effectively capture such failures.

**Planning.** Planning demonstrates only moderate correlation with human ratings (r = 0.43, $\rho$ = 0.51). This reflects the inherent difficulty of automatically evaluating long-horizon reasoning and structured action sequences, as well as the limited planning ability of current video models.

9

## 10.3. Ground-truth video replay of IDM

To ensure that our Inverse Dynamics Model (IDM) provides reliable supervision in the IDM Turing Test, we first validate the accuracy of a reproduced Gripper-Centric IDM (GC-IDM) following WoW [17]. Before applying the IDM to model-generated videos, it is essential to verify that its action inference is sufficiently accurate on ground-truth real-world data; otherwise, low execution accuracy could stem from model weakness rather than the generated video being physically unrealistic.

We therefore evaluate GC-IDM across 9 manipulation tasks, each with 10 ground-truth execution videos. They are Easy (Pick bread and place on the plate; Close the drawer; Move the milk in front of the cup), Medium (Pick bread and place it in the upper drawer; Take the cup off the cup holder; Open the drawer) and Hard (Flip the button to the correct position; Hang the cup on the cup holder; Insert chopsticks into the bamboo tube). For each trial, we feed the real video sequence into GC-IDM, infer the corresponding gripper-centric action sequence, and replay it in the real world to measure Replay Execution Accuracy. As shown in Table 7, we showcased four of the tasks, GC-IDM substantially outperforms standard baselines such as ResNet-based inverse dynamics model and AVDC [43] across all tasks, achieving an overall replay accuracy of 90%, demonstrating strong fidelity in mapping real videos to executable actions. Representative frames from ground-truth video replays are provided in Figure 7.

This high replay accuracy confirms that GC-IDM is a reliable evaluator of physical plausibility. Consequently, when GC-IDM fails on model-generated videos, we can attribute the failure to deficiencies in video realism rather than inaccuracies in the IDM itself.

Table 7. **Comparison of Success Rate (%) on various tasks across different models on ground-truth video replay.**

| Model / Task | Bread to Plate | Close Drawer | Move Milk | Bread in Drawer |
|---|---|---|---|---|
| **ResNet-MLPs** | 6/10 | 7/10 | 5/10 | 4/10 |
| **AVDC** [43] | 3/10 | 4/10 | 4/10 | 2/10 |
| **GC-IDM** [17] | **10/10** | **9/10** | **9/10** | **8/10** |

## 11. Case Results Visualization

We also collect more generated results of different models for comparision, as visualized in Figure 8, 9, 10.

## 12. Prompts

For the prompts that we used in GPT for evaluation and data selection, or our own trained Physical MLLM.

---

### GPT Data Selection Prompt

```
You are a robotics data selection and
    instruction rewriting assistant.
Input: one task instruction (text) and
    one first-frame image.
Output: JSON only, describing which
    data dimensions this sample
    belongs to, with evidence,
    confidence scores, and a rewritten
    instruction only for Perception
    dimensions.

Absolute Rules

Use only what is visible in the first
    frame and what is stated in the
    instruction. Do not hallucinate or
    import outside knowledge.

A dimension is valid only if:

The image provides clear visual
    evidence, and

The original instruction can be
    unambiguously rewritten (if
    Perception) or directly preserved
    (if Prediction/Planning).

If evidence is insufficient or
    ambiguous, do not assign the
    dimension; instead, explain in
    issues.

Multiple dimensions can be assigned;
    each must include its own score
    and evidence.

Instruction rewriting requirement:

Perception including rewrite
    instruction into an equivalent
    form using attributes / functions
    / spatial relations / affordances.

Prediction and Planning will always
    keep the original instruction
    unchanged.

Dimension Checklist
A. Perception

object-centric / object-attribute (
    color, number, shape, size, type)
```
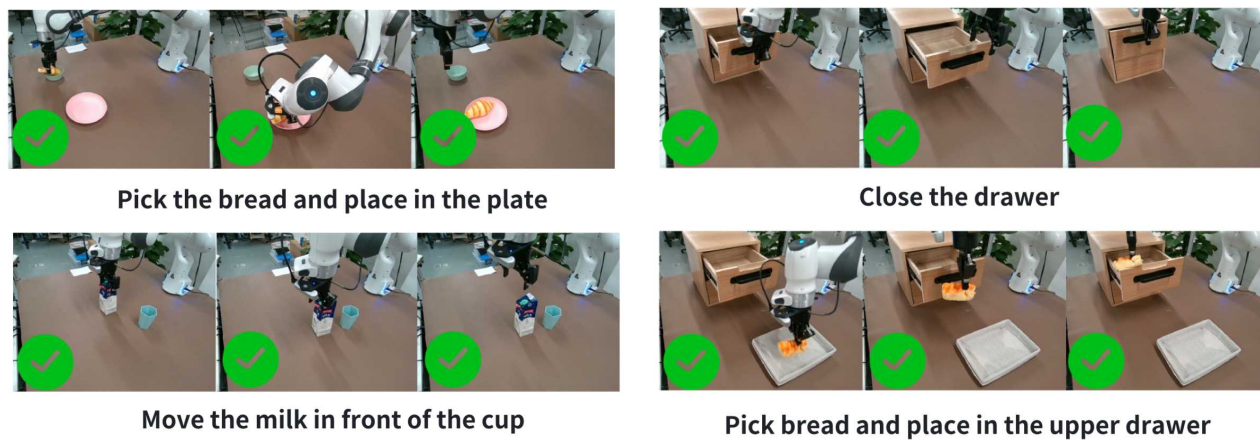
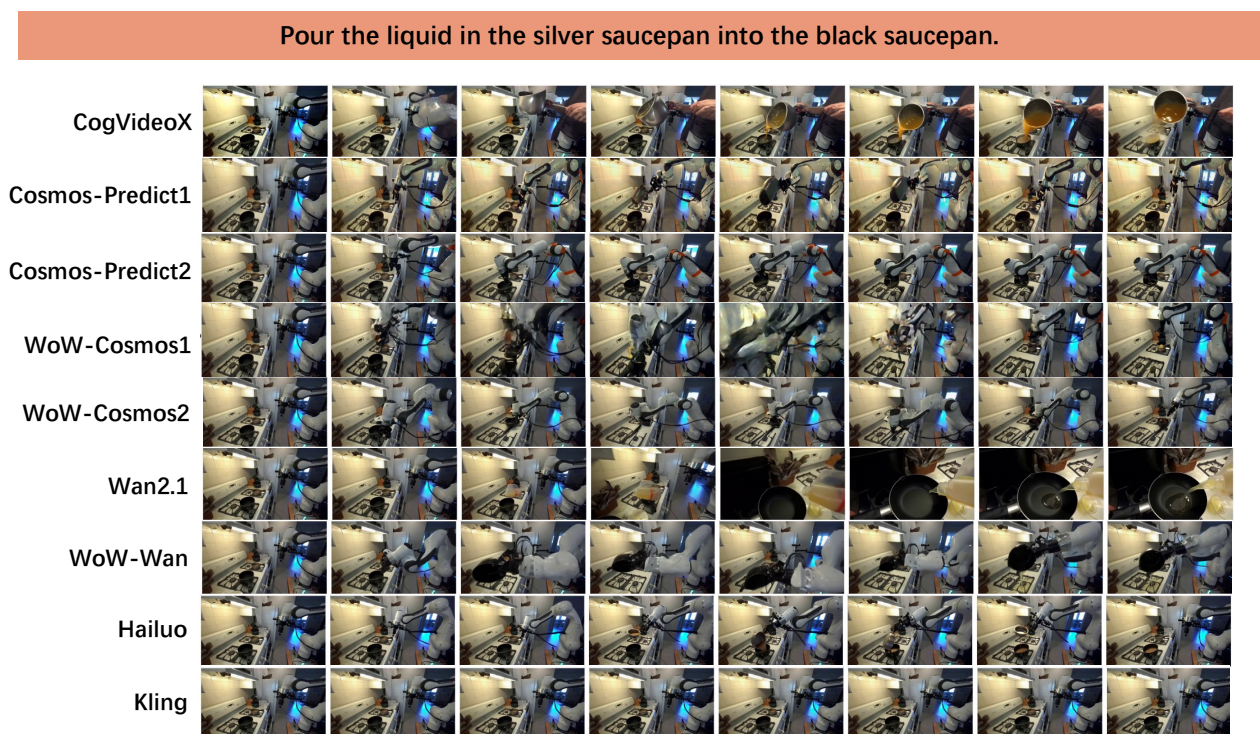Figure 7. **Real-World ground-truth video replay on real robot.**



Figure 8. **More Case Visualization in Perception.**

```
Paths must follow the format:
'Perception/object-centric/object-
    attribute/color'
'Perception/object-centric/object-
    attribute/type'
```

```
    etc.

Rewrite into attribute-based
    description.
```

**Pull up the waste collector liner.**

CogVideoX
Cosmos-Predict1
Cosmos-Predict2
WoW-Cosmos1
WoW-Cosmos2
Wan2.1
WoW-Wan
Hailuo
Kling

Figure 9. **More Case Visualization in Prediction.**

```
object-centric / object-function

Rewrite into function-based
    description.

scene-centric / spatial-relation

Rewrite into relation-based
    description.

affordance-centric / object-affordance

Rewrite into affordance-based
    description.

B. Prediction

camera-view including no-occlude |
    semi-occlude.

physical-interaction including single-
    object/rigid | single-object/fluid
    | single-object/deformable |
    single-object/articulated | multi-
    object/rigid-rigid | multi-object/
```

```
    rigid-fluid | multi-object/rigid-
    deformable | multi-object/fluid-
    fluid | dual-arm-cooperate

corner-case.

Keep original instruction (no
    rewriting).

C. Planning

long-term-instruction.

Keep original instruction (no
    rewriting).

JSON Output Format
{
  "input_instruction": "<original
      instruction>",
  "dimensions": [
    {
      "path": "Perception/object-
          centric/object-attribute/color
          ",
```

Figure 10. **More Case Visualization in Generalization.**

```
    "score": 0.9,
    "evidence": "blue plate, red
        tomato, yellow pepper visible
        ",
    "rewrite": "move the red and
        yellow objects to the blue
        object"
  },
  {
    "path": "Prediction/camera-view",
    "score": 0.8,
    "evidence": "all target objects
        fully visible",
    "rewrite": "keep original
        instruction"
  }
],
"prediction": {
 "camera-view": {
  "label": "no-occlude | semi-
      occlude | unknown",
  "evidence": "<short evidence>",
  "score": 0.0
  },
```

```
      "physical-interaction": [
        {
          "label": "single-object/rigid |
              ...",
          "evidence": "<short evidence>",
          "score": 0.0
        }
      ]
    },
    "planning": [
      {
        "label": "long-term-instruction",
        "evidence": "<short evidence>",
        "score": 0.0,
        "rewrite": "keep original
            instruction"
      }
    ],
    "issues": [
      "<if there is missing evidence or
          ambiguity, explain here>"
    ]
}
```

```
For Perception dimensions, rewrite
    contains the modified instruction.

For Prediction and Planning dimensions
    , rewrite must be exactly "keep
    original instruction".

score is confidence for 0 to 1.

evidence is short, factual, visual
    clues only.

**Example**

Input:

instruction: move tomato and pepper to
    the plate

image: blue plate, red tomato, yellow
    pepper; no occlusion.

Output:

{
  "input_instruction": "move tomato
     and pepper to the plate",
  "dimensions": [
    {
      "path": "Perception/object-
          centric/object-attribute/color
          ",
      "score": 0.95,
      "evidence": "blue plate; red
          tomato; yellow pepper visible
          ",
      "rewrite": "move the red and
          yellow objects to the blue
          object"
    },
    {
      "path": "Prediction/camera-view",
      "score": 0.9,
      "evidence": "all target objects
          fully visible",
      "rewrite": "keep original
          instruction"
    }
  ],
  "prediction": {
    "camera-view": {
    "label": "no-occlude",
    "evidence": "all objects clearly
        visible",
```

```
    "score": 0.9
    },
    "physical-interaction": [
      {
        "label": "single-object/rigid",
        "evidence": "rigid tabletop
            objects",
        "score": 0.8
      }
    ]
  },
  "planning": [],
  "issues": []
}
```

## Instruction Semantic Correctness: Extract Video Semantic Caption

```
Please analyze the instruction that I
    used to guide the model to
    generate video and the content of
    the following groundtruth video(
    there maybe no groundtruth video,
    if so please use only the
    instruction), extract semantic
    descriptions related to the robot'
    s behavior. Focus on the following
    five aspects and provide one
    clear, concise sentence for each:
1. Initial State: The state of the
    object(s) and environment before
    any action (e.g., position,
    posture, quantity)
2. Processing State: The interaction
    process while the robot arm is
    operating (e.g., how it moves or
    manipulates the object)
3. Final State: The outcome after the
    action (e.g., new position or
    state of the object, whether the
    task goal was achieved)
4. Action: The main type of operation
    performed (e.g., grasping, pushing
    , placing)
5. Object: The object being
    manipulated (e.g., its color,
    shape, category)

Please output your result in the
    following format:
```
Initial State: ...
```

14

```
Processing State: ...
Final State: ...
Action: ...
Object: ...
```

Guidelines:
- Keep each sentence objective,
  specific, and based only on what
  is visible in the video.
- Do not add assumptions or imagined
  information.
- If any aspect is unclear or not
  visible, write "Unknown" for that
  field.
- Avoid overly abstract descriptions
  such as "successfully completed
  the task."

## Instruction Semantic Correctness: Compare Video Semantic Caption

```
You are given structured semantic
    captions from two videos: the
    original ground truth and a
    generated version. Compare them
    across the following five
    dimensions and evaluate their
    semantic consistency.
"comparison_dimensions": {
    "initial_state": "The object/
        environment status before any
        action",
    "processing_state": "How the robot
        interacts with the object
        during operation",
    "final_state": "The outcome or
        resulting object state",
    "action": "The type of manipulation
        performed",
    "object": "The object being
        manipulated"
},
"scoring_guidelines": {
    "level_1": {
        "score": 1,
        "description": "Fully Consistent
            , Descriptions convey the
            same semantic content with
            only stylistic or wording
            differences"
    },
```

```
    "level_2": {
        "score": 0.5,
        "description": "Partially
            Consistent, Descriptions
            differ in minor semantic
            details or omit secondary
            information, but task intent
             remains aligned"
    },
    "level_3": {
        "score": 0,
        "description": "Inconsistent,
            Descriptions conflict
            semantically, refer to
            different actions, goals, or
            object properties"
    }
},
"output_format": {
    "initial_state": "Score = [0 / 0.5
        / 1] – Reason: ...",
    "processing_state": "Score = [0 /
        0.5 / 1] – Reason: ...",
    "final_state": "Score = [0 / 0.5 /
        1] – Reason: ...",
    "action": "Score = [0 / 0.5 / 1] –
        Reason: ...",
    "object": "Score = [0 / 0.5 / 1] –
        Reason: ...",
    "overall": "Score = [0 / 0.5 / 1] –
        Summary Reason: ..."
},
"evaluation_principles": [
"Focus strictly on the semantic
    content of the captions, do not
    infer or assume actions or states
    that are not explicitly stated.",
"The overall score may reflect the
    average of the five dimensions,
    but should also consider whether
    the task intent and outcome remain
     semantically aligned."
]
Ground truth video caption is: {ground
    _truth_caption}
Generated video caption is: {generated
    _caption}
```

## Instruction Semantic Correctness: Evaluate Video Action Execution

You are given a language instruction
    describing a series of robot
    manipulation actions, and a
    generated video that attempts to
    follow that instruction.
Your task has two parts:
PART 1: Action Sequence Consistency (
    Proportional Scoring)
Parse the instruction and extract the
    intended object-action pairs in
    order.
Analyze the generated video to
    determine the actual order of the
    executed object-action pairs.
    Compare the two sequences and
    compute a sequence match score:
    Sequence Match Scoring:
        Score = (number of matching
            object-action pairs in
            correct order) / (total
            number of instruction pairs)
        Matching means: the object-
            action pair appears in the
            video and in the same
            relative order as in the
            instruction. Partial matches
            or out-of-order actions
            reduce the score.
    Output format:
        Instruction Sequence: [(object1,
            action1), (object2, action
            2), ...]
        Video Sequence: [(object1,
            action1), (object2, action2)
            , ...]
        Sequence Match Score: X.XX (
            range: 0.00-1.00) - Reason:
            ...

PART 2: Action Execution Quality (
    Numerical Evaluation)
For each (object, action) pair in the
    instruction, evaluate execution
    quality using a 1-5 numerical
    scale:
    Execution Quality Scoring:
        1 - Object did not move AND
            action was not attempted
        2 - Object moved slightly OR
            action was partially
            attempted
        3 - Object moved slightly AND
            action was partially
            attempted

            4 - Object reached target state
                OR action completed
                successfully
            5 - Object reached target state
                AND action completed
                successfully

    Base this strictly on visible
        evidence in the video. Be
        conservative in scoring.

    Output format:
        Execution Quality:
        - (object1, action1): Score =
            [1-5] - Reason: ...
        - (object2, action2): Score =
            [1-5] - Reason: ...

Final Combined Output:
    Instruction Sequence: ...
    Video Sequence: ...
    Sequence Match Score: X.XX (range:
        0.00-1.00) - Reason: ...
    Execution Quality:
    - (object1, action1): Score = [1-5]
        - Reason: ...
    - (object2, action2): Score = [1-5]
        - Reason: ...
    Instruction is: {instruction}

## Physical Common Sense: Stage 2: Scoring Alignment Fine-tuning Prompt

You are an expert AI video evaluator.
    Your task is to analyze the
    provided video, which was
    generated from the accompanying
    text prompt. You must score the
    video on a scale of 1 to 5 across
    four dimensions based on the
    precise criteria listed below.

Evaluation Criteria:

1. Video Quality (1-5):

    5: Clear, stable, natural colors,
        no generation artifacts.

    4: Minor flaws (e.g., occasional
        flicker, slight blur) that don'
        t detract from the overall view

.

3: Noticeable flaws (e.g.,
   persistent noise, subject out
   of focus) but core content is
   understandable.

2: Severe defects (e.g., unstable
   subject form, heavy jitter)
   that significantly impair
   understanding.

1: Extremely poor quality, making
   the subject or content
   unrecognizable.

2. Instruction Following (1-5):

5: Perfectly matches the prompt;
   all core elements and details
   are accurately represented.

4: Core elements are correct, but a
   minor detail (e.g., color,
   location) is missed or
   incorrect.

3: The core element (subject or
   action) deviates significantly
   but is still partially related
   to the prompt.

2: The core element is
   fundamentally misinterpreted;
   the content is seriously
   inconsistent with the prompt.

1: Completely unrelated to the
   prompt; a random generation.

3. Physical Consistency (1-5):

5: Fully consistent with physics;
   interactions are believable, no
   object clipping or floating.

4: A momentary, almost unnoticeable
   physical anomaly (e.g., slight
   clipping) that doesn't affect
   the task.

3: A noticeable physical error (e.g
   ., arm passes through a table)
   that doesn't break the core
   task logic.

2: A severe physical error that
   causes the task to fail or its
   logic to break (e.g., a hand
   passes through the target
   object).

1: A complete lack of physical
   logic; chaotic and disorderly
   interactions.

4. Planning Logic (1-5):

5: The plan is logical, efficient,
   with no redundant actions, and
   successfully achieves the goal.

4: The plan is effective and
   successful but contains minor
   redundant or inefficient
   movements.

3: The goal is achieved, but the
   process is convoluted or clumsy
   ; the plan is illogical.

2: The plan has critical errors or
   missing steps, leading to task
   failure.

1: Actions are random and chaotic,
   with no discernible goal or
   plan.

Analyze the following inputs:

Original Prompt: {instruction}
Generated Video: <video>
Your final answer must be a single,
   valid JSON object and nothing else
   . Do not include any explanations,
    conversational text, or markdown
    formatting around the JSON block.
Your JSON Answer:
{
  "video_quality": <integer_score>,
  "instruction_following": <integer_
     score>,
  "physical_consistency": <integer_
     score>,
  "planning_logic": <integer_score>
}

## Physical Common Sense: 6-Dimension Evaluation Prompt

Please evaluate the following video
    strictly based on whether it
    follows real-world physics laws.
This is not a test of whether it is AI
    -generated, but only judge its
    physical plausibility.
Scoring System (1 to 5 scale):
For each of the six categories below,
    assign a score from 1 to 5:
5 = Fully consistent with real-world
    physics
4 = Mostly consistent, only small
    flaws
3 = Partially consistent, with some
    clear issues
2 = Often inconsistent, unrealistic in
     multiple ways
1 = Clearly violates physical laws
- null = Not applicable (no relevant
    content in video)
For each of the following categories:
- Give a score (1 to 5) or null if the
     category does not appear at all
    in the video
- Provide a brief explanation (1 to 2
    sentences) for each score
Return your output in structured JSON
    as shown below.

---

Evaluation Criteria (Detailed):
1. Object Interaction and State
    Changes
  - Do objects respond naturally when
     they touch, collide, bounce,
     break, or are pushed/pulled?
  - Is there a realistic cause-and-
     effect in object contact?
  - Do interactions follow Newton's
     laws (e.g., action/reaction,
     inertia)?
  - Are deformations (e.g., squashing,
      bending) consistent with
     materials?
  - Do objects maintain consistent
     color, quantity, shape, and size
     throughout interactions?
2. Basic Physical Properties
  - Does gravity act in the correct
     direction and magnitude (e.g.,
     falling speed)?

  - Is there realistic motion under
     inertia and deceleration?
  - Does friction behave naturally (e.
     g., sliding slows over time)?
  - Are materials (heavy/light)
     behaving correctly based on mass
     or density?
3. Temporal and Causal Consistency
  - Do effects follow causes in the
     right order and with realistic
     timing?
  - Are there appropriate delays in
     mechanical systems, collisions,
     or triggered motions?
  - Is there a consistent flow of time
      across all elements (e.g., no
     time jumps or impossible overlaps
     )?
4. Lighting, Shadows, and Reflections
  - Are shadows cast in the correct
     direction and shape relative to
     light sources?
  - Is there consistent lighting
     across surfaces and objects?
  - Do reflective surfaces (mirrors,
     metal, water) behave correctly
     based on angles?
  - Is light intensity and falloff
     physically accurate (not flat or
     overly uniform)?
5. Fluid and Particle Behavior
  - Do water, smoke, fire, or dust
     behave according to fluid
     dynamics?
  - Is there natural turbulence,
     diffusion, and randomness?
  - Do particles (like debris, sparks,
      splashes) move in a physically
     plausible way?
  - Is there appropriate interaction
     with solid surfaces (e.g.,
     splashes bounce, smoke drifts)?
6. Local Anomalies or Violations
  - Are there any visual
     inconsistencies like:
   - sudden object movement without
      cause
   - teleportation or disappearing
      items
   - objects passing through others (
      no collision)
   - animation glitches, floating
      artifacts, or broken motion?
  - Do all parts of the scene stay
     consistent frame to frame?

```
---
Return your output in structured JSON
    as shown below.
    - object_interaction: {score: int|
      null, comment: str}
    - physical_properties: {score: int|
      null, comment: str}
    - temporal_consistency: {score: int
      |null, comment: str}
    - lighting_and_reflections: {score:
       int|null, comment: str}
    - fluids_and_particles: {score: int
      |null, comment: str}
    - local_anomalies: {score: int|null
      , comment: str}
```

## Dense Prompts Extension Prompt

```
You are an expert in Embodied AI
    tasked with generating dense and
    standardized textual descriptions
    of robot episodes. Your goal is to
     re-caption a brief, potentially
    ambiguous, episode description
    into a richer, more precise
    narrative,but no more than 150
    words. This serves two main
    purposes:
1. To clarify ambiguous action
    descriptions with precise
    technical terminology.
2. To unify the textual representation
     of similar actions across
    different datasets, ensuring
    consistency in phrasing and detail
    .

When generating the dense caption,
    adhere to the following structure
    and include these specific details
    :

1. **Overall Scene & Goal Overview:**
    * Begin with a concise summary of
       the robot's primary goal or the
        main action being performed in
        the scene, directly derived or
        inferred from the input .

2. **Environment Description:**
    * Describe the robot's operational
        environment (e.g., kitchen,
```

```
       laboratory, office, cluttered
       tabletop, outdoor path).
    * Include relevant ambient
       conditions: lighting (e.g.,
       bright, dim, natural,
       artificial), indoor/outdoor
       setting, and weather if
       applicable (e.g., sunny,
       overcast for outdoor scenes).

3. **Robot Characteristics:**
    * Specify key characteristics of
       the robot if inferable or if it
       's a common type in embodied AI
        (e.g., humanoid, mobile
       manipulator with a single arm,
       quadruped, drone).
    * Note its general appearance if
       significant (e.g., color scheme
       , predominant material like
       metallic or plastic).
    * Mention its gripper type if
       relevant to the action (e.g.,
       parallel jaw gripper, suction
       cup, multi-finger dexterous
       hand).

4. **Camera Perspective & Motion:**
    * Detail the camera's perspective (
       e.g., first-person (egocentric/
       robot's POV), third-person
       fixed, third-person dynamic/
       following, eye-in-hand,
       overhead).
    * Describe any significant camera
       movement (e.g., static, panning
        to follow action, smooth
       tracking shot, handheld jitter
       if applicable).

5. **Robot Sub-goal & Granular Action
    Breakdown:**
    * Infer and explicitly state the
       immediate sub-goal the robot is
        trying to achieve based on the
        input.
    * Provide a granular breakdown of
       the robot's actions. Use
       precise, standardized technical
        action terminology (e.g., "
       Navigate_to [location]", "
       Perceive [object]", "Grasp [
       object]", "Lift [object]", "
       Move_arm_to [target_pose]", "
       Place [object] at [location]",
```

"Release [object]", "Push [object]", "Rotate [object]").
   * Clearly identify the object(s) being manipulated and any changes in their state or location.

6. **Post-Action State:**
   * Briefly describe the state of the robot and the key manipulated object(s) *after* the described action sequence is completed.

**Constraints:**
* The final description should be dense, focusing on factual and observable details.
* Aim for a total length of approximately 100-150 words.
* Use consistent terminology for common actions and objects.
* Avoid filler words or subjective stylistic descriptions not grounded in typical embodied AI contexts.
* Only describe elements plausibly present or inferable from a typical embodied AI episode matching the 'Given Input'. Do not invent fantastical details.
* Never mention elements outside the implied scene or action.

---
**Example of Use:**

Let's say the given short episode description is: "drop the apple in the bowl"

**Given Input :drop the apple in the bowl

**Your Dense Recaption (Output based on the prompt above):**

The robot's primary goal is to place an apple into a designated bowl.
The scene is set in a brightly lit kitchen environment, specifically on a clean countertop. The robot is a single-arm manipulator, predominantly white with grey joints, equipped with a parallel jaw gripper, likely mounted on a stationary or mobile base (partially visible). The camera is a static third-person view, positioned to offer a clear side-angle of the robot's arm and the workspace, including the target bowl.
The robot's sub-goal is to successfully deposit the apple. The action sequence begins with the robot's arm, already grasping a red apple, positioned directly above a ceramic bowl. The arm executes a 'Move_arm_to [bowl_center_above]' command, followed by a 'Lower_arm' motion. The gripper then performs a 'Release [apple]' action. Finally, the arm might execute a 'Retract_arm' motion away from the bowl.
After the action, the red apple rests inside the bowl on the countertop, and the robot's gripper is open and clear of the bowl.

Given Input : {instruction}
Your Dense Recaption: