

VBench-2.0: Advancing Video Generation Benchmark Suite for Intrinsic Faithfulness

Dian Zheng^{1,3*} Ziqi Huang^{2*} Hongbo Liu¹ Kai Zou¹ Yanan He¹ Fan Zhang¹
Yuanhan Zhang² Jingwen He^{1,4} Wei-Shi Zheng^{3✉} Yu Qiao^{1✉} Ziwei Liu^{2✉}

¹Shanghai Artificial Intelligence Laboratory ²S-Lab, Nanyang Technological University

³Sun Yat-Sen University ⁴The Chinese University of Hong Kong

<https://vchitect.github.io/VBench-2.0-project/>

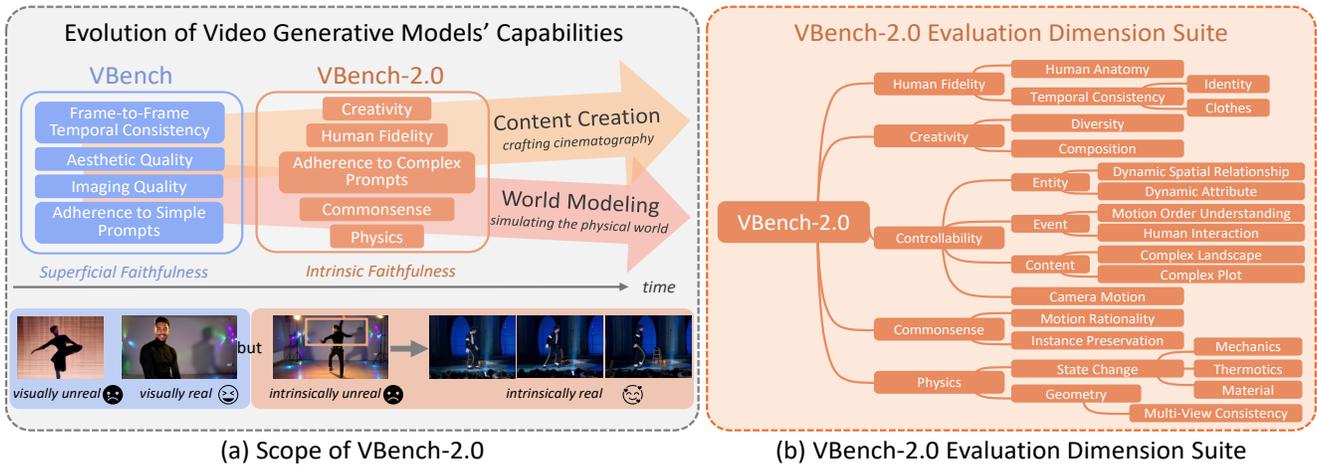


Figure 1. **Overview of VBench-2.0.** (a) **Scope of VBench-2.0.** Video generative models have progressed from achieving *superficial faithfulness* in fundamental technical aspects such as pixel fidelity and basic prompt adherence, to addressing more complex challenges associated with *intrinsic faithfulness*, including commonsense reasoning, physics-based realism, human motion, and creative composition. While VBench primarily assessed early-stage technical quality, VBench-2.0 expands the benchmarking framework to evaluate these advanced capabilities, ensuring a more comprehensive assessment of next-generation models. (b) **Evaluation Dimension of VBench-2.0.** VBench-2.0 introduces a structured evaluation suite comprising five broad categories and 18 fine-grained capability dimensions.

Abstract

Video generation has advanced significantly, evolving from producing unrealistic outputs to generating videos that appear visually convincing and temporally coherent. To evaluate these video generative models, benchmarks such as VBench have been developed to assess their faithfulness, measuring factors like per-frame aesthetics, temporal consistency, and basic prompt adherence. However, these aspects mainly represent *superficial faithfulness*, which focus on whether the video appears visually convincing rather than whether it adheres to real-world principles. While recent models perform increasingly well on these metrics, they still struggle to generate videos that are not just visually plausible but fundamentally realistic. To achieve real “world models” through video generation, the next frontier

lies in *intrinsic faithfulness* to ensure that generated videos adhere to physical laws, commonsense reasoning, anatomical correctness, and compositional integrity. Achieving this level of realism is essential for applications such as AI-assisted filmmaking and simulated world modeling.

To bridge this gap, we introduce **VBench-2.0**, a next-generation benchmark designed to automatically evaluate video generative models for their *intrinsic faithfulness*. VBench-2.0 assesses five key dimensions: Human Fidelity, Controllability, Creativity, Physics, and Commonsense, each further broken down into fine-grained capabilities. Tailored to individual dimensions, our evaluation framework integrates generalists such as state-of-the-art VLMs and LLMs, and specialists, including anomaly detection methods proposed for video generation. We conduct extensive human preference annotations to ensure evaluation

*equal contributions. ✉corresponding authors. Code is available

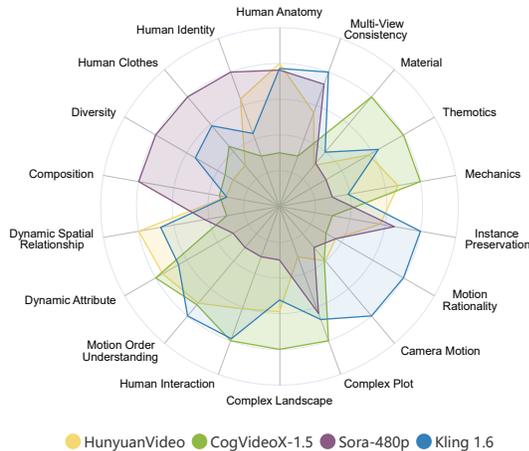


Figure 2. **VBench-2.0 Evaluation Results of SOTA Models.** The figure presents the evaluation results of four recent state-of-the-art video generation models across 18 VBench-2.0 dimensions. The results are normalized per dimension for a clearer comparison. For detailed numerical results, refer to Table 2.

alignment with human judgment. By pushing beyond superficial faithfulness toward intrinsic faithfulness, VBench-2.0 aims to set a new standard for the next generation of video generative models in pursuit of intrinsic faithfulness.

1. Introduction

Video generation aims to create realistic and temporally coherent video sequences, with a wide range of applications in video editing [4, 5, 9, 22, 30, 39, 42, 44, 58, 79, 85, 92–94], customization [25, 35, 38], image animation [23, 83], and world models [1].

Earlier video generative models [7, 29, 78, 80] primarily focused on generating short video clips of around two seconds, emphasizing fundamental capabilities like per-frame aesthetics and temporal consistency. To systematically evaluate these capabilities, benchmarks [32, 33, 46] such as VBench [32, 33] have been developed to assess aspects like per-frame aesthetics, frame-to-frame temporal smoothness, and adherence to simple text prompts, which we refer to as *superficial faithfulness*, the degree to which generated videos appear visually convincing. As video generative models continue to evolve, recent state-of-the-art models, including Sora [53], Kling [69], Gen-3 [59], HunyuanVideo [72], and Veo 2 [68], have demonstrated strong performance on these metrics, and many aspects of superficial faithfulness are now approaching saturation.

However, as video generation moves towards more advanced applications, particularly in areas that require AI models to simulate and reason about the real world [1], such as AI-driven storytelling and video-generation-based simulation, the new frontier shifts from *merely appearing real to*

being intrinsically real. We term this as **intrinsic faithfulness** - a concept that extends beyond per-frame quality and smooth motion, requiring that generated videos adhere to deeper principles such as physical laws, commonsense reasoning, anatomical correctness, and compositional integrity. Achieving this level of faithfulness is essential for applications ranging from AI-assisted filmmaking to virtual environments for embodied intelligence, ultimately paving the way for the development of true world models that can accurately represent and predict real-world dynamics.

To drive video generative models towards this next generation of capabilities, we introduce **VBench-2.0**, a benchmark suite designed to evaluate video generative models along five emerging dimensions beyond superficial faithfulness: *Human Fidelity*, *Controllability*, *Creativity*, *Physics*, and *Commonsense*. Each dimension is further broken down into isolated sub-abilities (shown in Figure 1(b)), providing a *fine-grained assessment* of the *intrinsic faithfulness* of video generative models. Given the complexity of these evaluations, we leverage *generalists* such as state-of-the-art Video Language Models (VLMs) and Large Language Models (LLMs) to perform structured reasoning and judgment. Specifically, we design two complementary evaluation methods: 1) *text description alignment* to assess abstract concept and semantic understanding, leveraging modern VLM’s strong captioning and LLM’s reasoning ability, and 2) *video-based multi-question answering* for basic visual understanding. To enhance evaluation robustness in specific domains, we incorporate *specialists*, such as human anomaly detection pipelines trained for generated videos. Furthermore, we use *evaluation safeguards* like *pre-filtering*, and *redundant questioning* to mitigate hallucinations and inconsistencies in our tailored evaluation pipelines for each dimension. Additionally, we follow VBench [32, 33] and conduct extensive human preference annotations to validate and align our automated evaluation results with human judgment.

Our evaluation provides *comprehensive insights* into the strengths and weaknesses of state-of-the-art video generative models. While recent models demonstrate emerging abilities in human anatomy, consistency, and some degree of novel creativity, they still struggle with generating complex plots, handling simple dynamic changes in objects, and remain unstable in commonsense reasoning, highlighting key open challenges in video generation towards synthesizing the world with *intrinsic faithfulness*. We provide an in-depth discussion on possible causes, potential solutions, and inherent trade-offs in Section 5.

VBench-2.0 will be fully open-sourced, being complementary to VBench [32, 33], and providing a standardized framework for evaluating future breakthroughs in video generation. While VBench remains essential for assessing *superficial faithfulness*, VBench-2.0 extends the evaluation

scope to *intrinsic faithfulness*, addressing deeper aspects of video realism. We will continually integrate newly released video generative models into VBench-2.0. By setting a higher standard for evaluation, VBench-2.0 aims to play a pivotal role in guiding the development of next-generation video generative models. Together, VBench and VBench-2.0 form a comprehensive benchmarking system, driving the field beyond superficial faithfulness towards truly intrinsically faithful video generation.

2. Related Works

Video Generative Models. With the advancements in diffusion models [2, 13, 14, 17, 27, 28, 31, 50, 63–65, 90], variational autoencoder-based compression techniques [16, 37, 56, 77, 87], and transformer architectures [15, 54], video generation has emerged as one of the most dynamic frontiers in artificial intelligence research. Prior to Sora’s breakthrough, predominant text-to-video models primarily focused on synthesizing short video clips (2-3 seconds duration) [3, 5, 7, 21, 23, 24, 29, 36, 48, 53, 57, 68, 78, 80, 86, 88, 95] through incremental improvements in visual fidelity and temporal consistency. Sora [53] pioneered the scaling paradigm in video generation by demonstrating unprecedented model capacity through large-scale training, paving the way for the development of next-generation video foundation models [1, 18, 55, 59, 61, 69–73, 86] that achieve remarkable visual quality and robust spatiotemporal coherence. They have shifted focus toward enhancing video generation adhere to deeper principles such as physical laws and commonsense reasoning that focuses more on action continuity and realistic physical-world perception [1, 53, 86], or high-quality human-centric generation with creativity potential [69, 72]. Existing benchmarks cannot systematically evaluate these new explorations, and VBench-2.0 takes the initiative to provide a comprehensive benchmark for evaluating emerging capabilities towards the goal of achieving intrinsic faithfulness through video generation.

Evaluation of Video Generative Models. Initially, video generative models primarily relied on conventional evaluation metrics such as Fréchet inception distance (FID) [26], Inception Score (IS) [60], and Fréchet video distance (FVD) [76]. However, these metrics provided limited insight into the diverse and complex capabilities of modern video generation. Recent evaluation frameworks [32, 33, 45, 46, 89] such as VBench [32, 33] introduced a more structured approach by disentangling evaluation into multiple capability dimensions, enabling more detailed and interpretable assessments. These benchmarks focus on fundamental technical attributes such as per-frame quality, temporal consistency, and basic prompt adherence. However, as models continue to improve, certain dimensions within VBench begin to saturate, necessitating broader evaluation

Table 1. **Comparison of Video Generation Benchmarks.** We compare existing video generation benchmarks based on their evaluation aspects. VBench-2.0 is the first comprehensive benchmark to assess intrinsic faithfulness in video generation, complementing VBench [32, 33]. Detailed aspects include per-frame quality (Frame Wise), temporal consistency (Temp Cons), adherence to simple prompts (Simp Pmpt), compositional creativity (Comp Crea), commonsense reasoning (Com Sense), physics-based realism (Phy), human anatomy (Human Anat), and adherence to complex prompts (Cplx Pmpt).

	Superficial Faithfulness			Intrinsic Faithfulness				
	Frame Wise	Temp Cons	Simp Pmpt	Comp Crea	Com Sense	Phy	Human Anat	Cplx Pmpt
VBench [32, 33]	✓	✓	✓					
T2V-CompBench [66]			✓	✓	✓			
PhyGenBench [49]						✓		
StoryEval [81]		✓						✓
VBench-2.0 (Ours)		✓		✓	✓	✓	✓	✓

scope that assess deeper aspects of intrinsic faithfulness in video generation. To address this, specialized benchmarks have emerged. PhyGenBench [49] evaluates a model’s understanding of physical laws through Vision-Language Models (VLMs). T2V-CompBench [66] assesses compositionality, including motion, actions, spatial relationships, and attributes. StoryEval [81] focuses on storytelling capabilities by aggregating the responses from two VLMs. Unlike prior benchmarks, which either focus on fundamental capabilities [32, 33, 46] or specific emerging domains [49, 66], VBench-2.0 introduces a comprehensive framework to systematically evaluate next-generation video generation capabilities, bridging the gap in the evolving landscape of video generation.

3. VBench-2.0 Suite for Intrinsic Faithfulness

In this section, we introduce the evaluation framework of VBench-2.0. Section 3.1 presents the five key evaluation dimensions and their respective assessment methods. Unlike VBench, which primarily evaluates *superficial faithfulness*, VBench-2.0 introduces a suite of tests to assess *intrinsic faithfulness*, focusing on deeper properties such as physics, commonsense, and creativity, human, and controllability. To ensure robust evaluation, we integrate multiple assessment methodologies, including LLM-assisted text alignment, video-based multi-question answering, and specialist models trained for anomaly detection. Our prompt suite is introduced in Section 3.2.

3.1. Evaluation Dimension Suite

VBench-2.0 evaluates video generation along five key dimensions: *Human Fidelity*, *Creativity*, *Controllability*, *Physics*, and *Commonsense*. Each dimension is further decomposed into sub-dimensions, ensuring a fine-grained assessment of a model’s capabilities. We employ a structured approach that combines *generalist* reasoning models (VLMs/LLMs) with *specialist* detectors. When using *generalist*, we adopt two evaluation schemes where each is tailored to different types of semantic understanding required

across evaluation dimensions.

Text Description Alignment. This scheme is suitable for complex or subtle scenarios, such as those involving nuanced human interactions or multi-step plots. In these cases, the composite scene is broken down and interpreted step by step by the Vision-Language Model (VLM), which generates a descriptive caption guided by system prompts that focused on specific aspects of the video (*e.g.*, prompting only about human interactions or a specific part of the plots). The correctness of the generated content is then judged by a Large Language Model (LLM), which compares the VLM-generated caption with a ground-truth reference. The reference may be the original text prompt, a predefined answer, or relevant metadata. This process is formalized as:

$$Answer = LLM(VLM(V|S_v), T | S_l), \quad (1)$$

where V is the generated video, T denotes the text prompt or crafted reference, and S_v and S_l are the system prompts for the VLM and LLM, respectively. The LLM outputs a binary judgment (“yes” or “no”) that is discretized to a score of 1 or 0 to reflect caption-reference matching. This scheme excels in semantic understanding dimensions like *Complex Plot* and *Human Interaction*, where VLMs usually struggle with high-level interpretation, but LLMs demonstrate stronger reasoning capabilities. Thus, we decouple caption generation (by VLMs) from semantic alignment (by LLMs) to improve evaluation reliability.

Video-Based Multi-Question Answering. It is designed for evaluation dimensions where one salient concept is prominent and can be directly queried through video question answering (VQA). In this approach, we construct a series of complementary and sometimes redundant questions to reduce the risk of accidental errors and ask the VLM to perform direct VQA. The formalization is as follows:

$$Answer = \sum_i^N VQA(Q^i, V | S), \quad (2)$$

where Q is a set of multiple questions. For example, in the *Dynamic Attribute* dimension focusing on color changes, we may ask: 1. *Initially, is the color of the river mostly blue?* 2. *Finally, is the color of the river mostly brown?* 3. *Does the color of the river change?* The answer to each question is binary (“yes” or “no”), and scores are either averaged or awarded only if all responses are correct, depending on the dimension’s scoring scheme. This scheme is particularly effective for surface-level visual understanding, where modern VLMs can confidently answer targeted queries without requiring high-level semantic reasoning.

3.1.1. Human Fidelity

We evaluate both the structural correctness and temporal consistency of human figures in generated videos. Structural issues commonly seen in current video generation

models, such as sudden turns or the “thousand-hand yoga” effect, are considered. For temporal consistency, we assess the consistency human identity and clothing across frames.

Human Anatomy. We assess structural realism by identifying unnatural per-frame deformations in hands, faces, and bodies. To achieve this, we train three anomaly detection models using a pre-trained ViT-based backbone [82], using a curated dataset of 150k real and generated human frames.

Human Temporal Consistency - Clothes. Clothing consistency is assessed using video-based multi-question answering to ensure outfits remain stable throughout video.

Human Temporal Consistency - Identity. Identity consistency is evaluated by measuring facial feature similarity using ArcFace [11], with face detection performed by RetinaFace [12].

3.1.2. Creativity

We evaluate creativity by analyzing a model’s ability to generate diverse outputs and complex compositions beyond real-world constraints.

Diversity. Given a text prompt, we sample 20 videos from a model, and measure inter-sample variation using style and content diversity metrics, computed from pre-trained VGG-19 [62] feature representations.

Composition. We assess species combination, single-entity actions, and multi-entity interactions using a structured video-based multi-question answering pipeline. This approach measures whether the model can generate novel and uncommon compositions.

3.1.3. Controllability

We evaluate a model’s ability to follow complex prompts and simulate dynamic changes during video generation. This dimension measures how accurately the model can render specific entities, events, content, and camera movements in response to detailed textual instructions.

Entity - Dynamic Spatial Relationship. We assess whether models accurately reposition objects in response to spatial instructions (*e.g.*, “A dog is on the left of a sofa, then the dog runs to the front of the sofa.”) using *video-based multi-question answering*.

Entity - Dynamic Attribute. We test whether models can modify attributes (*e.g.*, color, size, texture) mid-video using the *video-based multi-question answering* pipeline.

Event - Motion Order Understanding. We evaluate whether models generate several actions or motions in the specified order. We use the *text description alignment* pipeline to measure whether the generated motion sequence matches the text prompt.

Event - Human Interaction. We assess whether two humans can interact (*e.g.*, “One person hands an object to another”) based on the text prompt. We use the *text description alignment pipeline*.

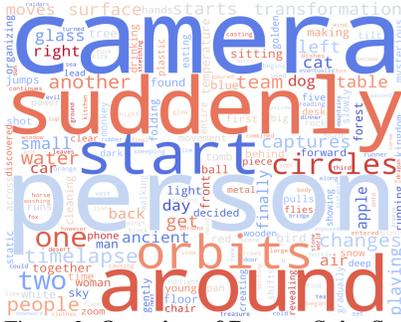


Figure 3. **Overview of Prompt Suite Statistics.** *Left:* distribution of words in the prompt suites. *Right:* number of prompts per evaluation dimension.

Content - Complex Landscape. We evaluate whether models faithfully follow long-form landscape descriptions (150+ words) that include multiple scene transitions driven by camera movements. We assess adherence using *text description alignment* with landscape-specific system prompt.

Content - Complex Plot. We test a model’s ability to construct multi-scene narratives from prompts describing multi-stage events (e.g., a four-act story with 150+ words). We evaluate plot consistency using *text description alignment*, similar to the landscape evaluation.

Camera Motion. We evaluate whether models can generate specified camera movements. We extend VBench++ [33]’s camera motion taxonomy to nine types, adding “Orbit” and “Oblique shot, airborne dolly movement.” The generated camera motion is assessed via point tracking with CoTracker-v2 [34] and carefully tailored heuristics.

3.1.4. Physics

The *Physics* dimension evaluates whether video generation models adhere to real-world physical principles. We assess two key areas: 1) *State Change*, which examines how well models simulate mechanical, thermal, and material transformations, and 2) *Geometry*, which evaluates the 3D consistency of objects and scenes across different frames. We extend prior benchmarks (e.g., PhyGenBench [49]) by increasing physics scenario difficulty and significantly improving evaluation accuracy with tailored pipelines.

State Change - Mechanics. We evaluate whether models follow basic mechanical principles such as gravity, buoyancy, and stress. This is done using a *video-based multi-question answering* pipeline. Unlike PhyGenBench, which relies on abstract physics concepts, we prompt GPT-4o to generate explicit *visual* descriptions of expected physical behavior. Notably, the GPT-4o-generated descriptions are based solely on the text prompt, *remain fixed during evaluation*, and *do not compromise evaluation reproducibility*. To ensure a focused assessment of state changes, we also apply a pre-filtering step to exclude cases where the initial state of the generated video does not align with the prompt.

State Change - Thermotics. We evaluate how well the models simulate state transitions such as vaporization, liq-

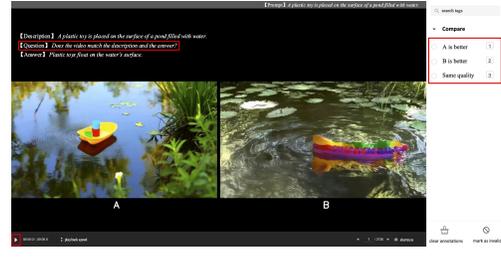
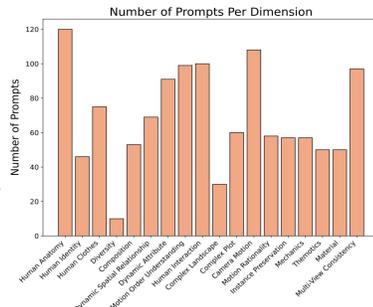


Figure 4. **Interface for Human Preference Annotation.** *Top:* Question descriptions. *Right:* Choices available to annotators. *Bottom left:* Controls for stopping and playback.

uefaction, and sublimation. To increase complexity, we introduce temperature-specific prompts (e.g., “A timelapse captures dry ice transforming at -90°C ”). Evaluation follows the same *video-based multi-question answering* approach used in the *Mechanics* dimension.

State Change - Material. We evaluate whether the models correctly depict color mixing, hardness, combustion, and solubility. Similar to *Mechanics* and *Thermotics*, we use video-based multi-question answering, testing whether generated videos properly follow material properties.

Geometry - Multi-View Consistency. Geometry is a critical aspect for 3D/4D video generation, ensuring that generated entities and scenes maintain structural consistency when viewed from different angles or as the camera moves. However, we do not have access to explicit 3D ground truth of the generated videos, making direct 3D validation infeasible. Instead, we assess multi-view consistency using two complementary metrics: 1) *Feature Matching Stability* – Measures how well objects retain their geometric consistency across frames, inspired by [41], and 2) *Camera Motion Speed* – Accounts for the effect of motion strengths on feature stability, ensuring fair cross-model comparisons. Specifically, we follow [41] to extract frame-level keypoint features using SIFT [47], efficiently match points across frames with FLANN [51], and eliminate incorrect matches using RANSAC [20]. We then estimate camera motion speed using RAFT [75], and adjust feature-matching frame intervals based on the motion strength and the video frame rate.

3.1.5. Commonsense

We assess commonsense reasoning in video generation models across two key aspects: 1) *Motion Rationality*, evaluating whether generated motions are physically plausible and correctly executed, and 2) *Instance Preservation*, ensuring that object counts remain stable throughout the video.

Motion Rationality. We evaluate whether the generated motion leads to the correct real-world consequences. A prevalent issue in video generation is the presence of fake motions that visually appear correct but lack expected effects on the environment. For instance: 1) Fake eating: A

Table 2. **VBench-2.0 Evaluation Results per Dimension.** This table presents evaluation results for four recent state-of-the-art video generation models across all 18 VBench-2.0 dimensions. A higher score indicates better performance in the corresponding dimension.

Models	Human Anatomy	Human Clothes	Human Identity	Composition	Diversity	Mechanics	Material	Thermotics	Multi-view Consistency
HunyuanVideo [72]	88.58%	82.97%	75.67%	43.96%	39.73%	76.09%	64.37%	56.52%	43.80%
CogVideoX-1.5 [86]	59.72%	87.18%	69.51%	44.70%	42.61%	80.80%	83.19%	67.13%	21.79%
Sora [53]	86.45%	98.15%	78.57%	53.65%	67.48%	62.22%	64.94%	43.36%	58.22%
Kling 1.6 [69]	86.99%	91.75%	71.95%	43.89%	53.26%	65.55%	68.00%	59.46%	64.38%
Models	Dynamic Spatial Relationship	Dynamic Attribute	Motion Order Understanding	Human Interaction	Complex Landscape	Complex Plot	Camera Motion	Motion Rationality	Instance Preservation
HunyuanVideo [72]	21.26%	22.71%	26.60%	67.67%	19.56%	10.11%	33.95%	34.48%	73.79%
CogVideoX-1.5 [86]	19.32%	24.18%	26.94%	73.00%	23.11%	12.42%	33.33%	33.91%	71.03%
Sora [53]	19.81%	8.06%	14.81%	59.00%	14.67%	11.67%	27.16%	34.48%	74.60%
Kling 1.6 [69]	20.77%	19.41%	29.29%	72.67%	18.44%	11.83%	61.73%	38.51%	76.10%

person bites into food, but the food remains unchanged. 2) Fake walking: A character moves their legs, but does not actually progress forward. and 3) Fake cutting: A knife moves through an object, but the object does not split. To address this, we adopt the *video-based multi-question answering* pipeline, specifically to ensure the motion’s impact is checked explicitly (e.g., “Did the food reduce in size after the bite?”). Redundant questioning helps filter out false positives and accidental misinterpretations.

Instance Preservation. Video generative models frequently struggle to maintain object counts due to unnatural merging, duplication, or disappearance, particularly during large-motion sequences or object interactions. We use YOLO-World [8], an open-vocabulary detection model, to detect and count entities frame by frame.

3.2. Prompt Suite

The VBench-2.0 Prompt Suite is designed to be compact yet representative. Given the increasing computational cost of video sampling, especially for longer and higher-resolution videos (e.g., HunyuanVideo and CogVideoX-1.5 taking over five minutes per sample on 8xA100 GPUs), we strategically limit the number of test cases to reduce sampling costs during evaluation, while ensuring coverage across diverse evaluation dimensions and content scenarios. Figure 3 visualizes the prompt distributions.

Tailored Prompts for Each Dimension. For each evaluation dimension in VBench-2.0, we carefully construct a suite of approximately 70 prompts, specifically tailored to probe the model’s capabilities in that dimension. Prompts are designed to systematically analyze the core ability being tested. For example, in the *Multi-View Consistency* dimension, we evaluate both object-level and scene-level prompts, ensuring that models are tested across varying spatial structures. In the *Composition* dimension, we systematically divide prompts into species combination, single-entity action, and multi-entity tasks, covering different levels of creativity and compositional reasoning. The *Physics* dimension follows a structured approach, incorporating *Mechanics* (e.g., gravity, buoyancy, stress), *Thermotics* (e.g., vaporization, freezing), and *Material* properties (e.g., color

mixing, solubility) to comprehensively assess adherence to physical laws. To further challenge model reasoning, additional constraints, such as temperature-specific prompts in *Thermotics*, require models to demonstrate a deeper understanding in how temperature affect thermotics beyond simple pattern matching.

Ensuring Disentangled Evaluation. Prompts are designed to eliminate confounding factors and ensure focused assessment in the dimension of interest. For example, in *Dynamic Spatial Relationship* and *Dynamic Attribute*, only one entity is allowed to move or changed, ensuring the test solely assesses positioning and attribute rather than taking the irrelevant ability of multi-object interactions into account.

Evaluation Robustness. To improve the robustness of evaluation, we design the actions and events in each dimension’s prompts to be explicitly recognizable by VLMs and LLMs. In the *Human Interaction* dimension, we ensure that prompts include physical contact interactions (e.g., “A person shakes hands with another”) rather than ambiguous social scenarios that are harder to verify visually (e.g., “Two people are having a picnic”). For *Motion Rationality*, prompts are designed to ensure that the counterexamples involve visually observable outcomes, such as fake eating (food remains unchanged), fake walking (not moving forward), or fake cutting (objects remain unaltered), rather than examples like lip-syncing, which are difficult to conclusively assess through visual observation alone.

By following these structured design principles, VBench-2.0 provides a *compact, diverse, and reliable* benchmark for evaluating video generation models across a diverse range of real-world and abstract scenarios. *The full details of prompt suite design for each dimension are provided in the Supplementary File.*

4. Experiments

4.1. VBench-2.0 Evaluation Results

We have used VBench-2.0 to evaluate the intrinsic faithfulness of four recent state-of-the-art video generative models so far, including Kling 1.6 [69], Sora-480p [53], HunyuanVideo [72] and CogVideoX-1.5 [86]. Additional models

Table 3. Information on Evaluated Models.

Model Name	Video Length	Per-Frame Resolution	Frame Rate (FPS)
HunyuanVideo [72]	5.3s	720×1280	24
CogVideoX-1.5 [86]	10.1s	768×1360	16
Sora-480p [53]	5.0s	480×854	30
Kling 1.6 [69]	10.0s	720×1280	24

will be incorporated as they become available. We provide an overview of the four evaluated models in Table 3, and describe the detailed setting of each model and the sampling procedures in the Supplementary File.

For each sub-ability dimension, videos were generated using the models based on the corresponding prompt suite described in Section 3.2. The evaluation method introduced in Section 3.1 is then applied to obtain numerical scores between 0 and 1, where a higher value indicate relatively stronger performance in that dimension. The evaluation results of the four video generative models are summarized in Table 2 and visualized in Figure 2.

4.2. Human Alignment of VBench-2.0

To ensure that VBench-2.0’s evaluation aligns closely with human judgment across all evaluation dimensions, we conducted human preference labeling on a large set of generated videos, following the approach of VBench [32, 33]. The annotation procedure and guidelines are similar to those used in VBench and are described in detail in the Supplementary File (along with human annotator workload and details). Multiple rounds of quality assurance were conducted to ensure annotation accuracy.

Following VBench, we computed the correlation between our evaluation results and human annotations. Figure 5 presents the correlation plot, illustrating the alignment between human judgment and VBench-2.0 evaluations in terms of model-level win ratios across video generation models in each dimension. Further details on the correlation computation are provided in the Supplementary, along with a table detailing numerical win ratios from both VBench-2.0 and human annotations for each model across all dimensions.

5. Insights and Discussions

In this section, we present key insights from VBench-2.0’s evaluation, highlighting trade-offs, model characteristics, and in-depth discussion on superficial versus intrinsic faithfulness in video generation.

5.1. Characteristics of Recent SOTA Models

From Figure 2, we can observe the relative strengths and weaknesses of each model under evaluation.

Sora-480p [53]. Sora clearly excels in *Human Fidelity* and *Creativity* dimensions compared to other SOTA models. It demonstrates a strong ability to generate human figures with reasonable anatomical consistency throughout a

video while also showing improvisational skill in producing novel and imaginative content. This makes Sora a potential tool for human-centric filmmaking and artistic exploration. However, it falls short in *Controllability*, *Physics* and *Commonsense* dimensions, indicating that the generated videos may not align well with user-provided text prompts and sometimes violate real-world principles.

Kling 1.6 [69]. Kling demonstrates relative strengths in the *Commonsense*, *Controllability* and camera-related (*Multi-View Consistency*, *Camera Motion*) dimensions. These capabilities make Kling well-suited not only for tasks that require precise camera control, but also for broader applications involving coherent, accurate, and user-guided visual storytelling or simulation. Additionally, Kling does not show significantly weaker performance in any particular area, suggesting that its training data is broad and well-rounded, making it a valuable reference for future model development.

CogVideoX-1.5 [86]. This model is relatively strong in most dimensions related to complex prompt adherence (*e.g.*, *Complex Landscape* and *Complex Plot*) and *Physics*, but shows notably poorer results in human-centric dimensions such as *Human Fidelity* and *Motion Rationality*. These outcomes suggest that CogVideoX’s training data contains limited high-quality human-related content.

HunyuanVideo [72]. Although HunyuanVideo is relatively weaker in many VBench-2.0 dimensions, it demonstrates impressive strengths in human-related aspects (*Human Fidelity* and *Motion Rationality*). This suggests that HunyuanVideo likely benefits from training data rich in high-quality, human-related content.

5.2. Key Limitations of Recent SOTA Models

Key Challenge: Generating Complex Plots. In the *Complex Plot* dimension, state-of-the-art video generation models struggle to follow detailed text descriptions involving multiple scenes, character interactions, and logical story progression. A major limitation is that current foundation video generative models typically produce single-shot videos under 10 seconds in length, insufficient for conveying coherent narratives. This highlights an important direction for future research: building models capable of functioning more like true filmmakers.

Surprising Weakness: Controllability in Simple Dynamics. Most models perform poorly in capturing *Dynamic Spatial Relationships* and *Dynamic Attributes*. Even in relatively simple cases, where an entity’s position, relationship, or attribute (*e.g.*, color) is instructed to change, and models fail in about 80% of the time. These shortcomings, despite the simplicity of the underlying semantics, are likely due to inadequate captioning granularity in video generation datasets. Existing video captioning pipelines may not be intentionally describing how object attributes or posi-

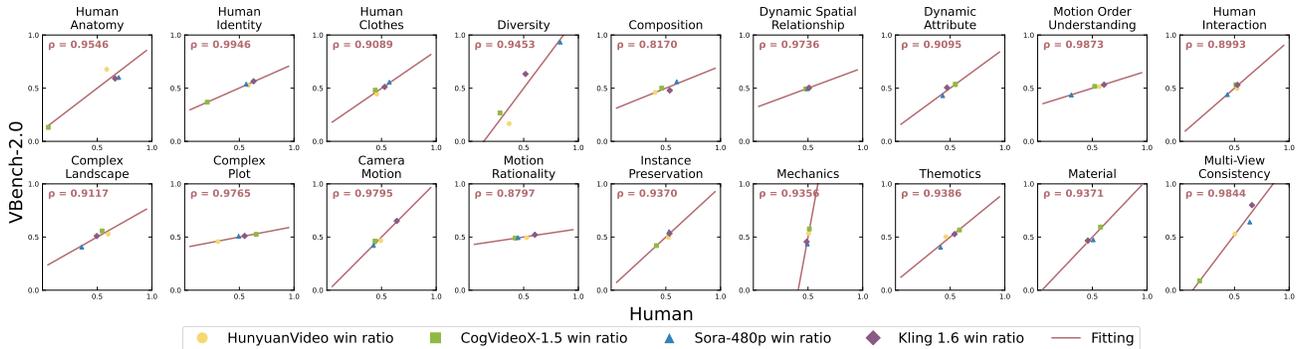


Figure 5. **Human Alignment of VBench-2.0 Evaluation.** Each plot represents the alignment verification for a specific VBench-2.0 dimension. In each plot, a dot corresponds to the human preference win ratio (horizontal axis) and the VBench-2.0 evaluation win ratio (vertical axis) for a given video generation model. A linear fit is applied to visualize the correlation, and Spearman’s correlation coefficient (ρ) is computed for each dimension. Experiments show that VBench-2.0 evaluations closely align with human judgement in all dimensions.

tions evolve over time, weakening models’ understanding of these dynamics. Enhancing these pipelines with more focused, temporally grounded instructions during video captioning could help address this gap.

5.3. The Role of Prompt Engineering

In modern video generation, the Prompt Refiner rewrites or augments input text prompts to enhance video generation quality. For implementation details, please refer to the *Supplementary File*. Below, we highlight several notable and potentially surprising observations.

Controllability vs. Creativity. Our evaluation results reveal a trade-off between *Creativity* and *Controllability*. Sora performs well in creative tasks but struggles with controllability, while the other models show the reverse pattern. This suggests that models emphasizing creativity and diversity are better at flexibly imagining novel content, but may be less capable of strictly and accurately following control signals. Alternatively, the Prompt Refiner used with the other three models may improve controllability by fine-graining the text at the expense of diversity. Sora’s internal prompt optimization appears to take a different approach. Going forward, prompt refinement strategies should consider not only precision but also the preservation of diversity, which is essential for creating open-ended content and simulating the distributions of the real world.

Physical Laws May Not Be So Challenging. Physical reasoning is often seen as a difficult frontier in video generation. Yet, except for Sora-480p, the other models guided by the external Prompt Refiner performed relatively well in the *Physics* dimension to some extent. This implies that even if models lack an inherent understanding of physical laws, a well-crafted prompt can guide them toward physically plausible outcomes. Thus, the primary challenge may not be physical reasoning itself, but rather achieving accurate text-video alignment. These insights highlight that the Prompt Refiner is more than a pre-processing tool; it is a critical component in enabling realistic simulation. Future

refiners could explore more sophisticated strategies to offset limitations in model comprehension.

No-Impact Dimensions. For dimensions that rely on model’s intrinsic visual understanding and prior knowledge, like *Human Fidelity*, *Camera Motion*, *Geometry*, and *Commonsense*, we observe no consistent performance trend from models using different Prompt Refiners. These aspects likely extend beyond the scope of logical inference or direct text-to-video mapping, limiting the Prompt Refiner’s influence. This suggests that success in these areas may depend less on prompt engineering and more on underlying data quality and model architecture.

5.4. Superficial Faithfulness vs. Intrinsic Faithfulness: Do Not Miss Out on Any Pillar

Superficial Faithfulness (e.g., cinematographic quality) often shapes the first impression viewers get from a video. As a result, models that produce aesthetically pleasing and smooth outputs are frequently perceived as “better”. However, this perception could be misleading. In practice, *Intrinsic Faithfulness*, which includes elements like storytelling, logical progression, and world simulation, is equally important for determining a model’s potential for real-world applications. For example, as shown in Figure 2 and Table 2, CogVideoX performs relatively well across many VBench-2.0 dimensions, though its visual *Quality Score* in VBench suggests room for improvement compared to models like Sora and HunyuanVideo. Conversely, HunyuanVideo produces visually impressive results, though its performance in many structure-driven dimensions in VBench-2.0 suggests opportunities for further growth in those areas. These observations highlight a common bias toward relying primarily on visual quality when judging video generation models. To address this, we encourage the community to use both VBench and VBench-2.0 together, enabling a more comprehensive and in-depth evaluation across both *Superficial Faithfulness* and *Intrinsic Faithfulness*.

6. Conclusion

While recent video generative models have achieved *superficial faithfulness*, true progress requires advancing towards *intrinsic faithfulness*, ensuring adherence to physical laws, commonsense reasoning, and structured coherence. To address this, we introduced VBench-2.0, a benchmark assessing models on five key dimensions beyond superficial faithfulness. VBench-2.0 complements VBench by expanding evaluation to deeper aspects of video generation, providing a multi-dimensional and human-aligned evaluation framework. We believe that VBench-2.0 is an important contribution to the video generation community, shaping the field into its next era.

Limitations and Future Work. We will continually add more video generative models to VBench-2.0 when they become available.

Potential Negative Societal Impacts. Although VBench-2.0 does not directly generate videos, the evaluation process inevitably involves working with generated content. As video generative models grow more powerful and capable of producing increasingly realistic scenes, we emphasize the importance of safety and ethical considerations when using these models. We encourage responsible use to mitigate potential risks associated with AI-generated media.

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. 2, 3
- [2] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 3
- [3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *CVPR*, 2023. 3
- [4] Duygu Ceylan, Chun-Hao P Huang, and Niloy J Mitra. Pix2video: Video editing using image diffusion. In *ICCV*, 2023. 2
- [5] Wenhao Chai, Xun Guo, Gaoang Wang, and Yan Lu. Stable-video: Text-driven consistency-aware diffusion video editing. *arXiv preprint arXiv:2308.09592*, 2023. 2, 3
- [6] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 19
- [7] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter1: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023. 2, 3
- [8] Tianheng Cheng, Lin Song, Yixiao Ge, Wenyu Liu, Xinggang Wang, and Ying Shan. Yolo-world: Real-time open-vocabulary object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16901–16911, 2024. 6, 18, 19
- [9] Paul Couairon, Clément Rambour, Jean-Emmanuel Haugeard, and Nicolas Thome. Videdit: Zero-shot and spatially aware text-driven video editing. *arXiv preprint arXiv:2306.08707*, 2023. 2
- [10] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, et al. Objaverse-xl: A universe of 10m+ 3d objects. In *NIPS*, 2023. 19
- [11] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *CVPR*, 2019. 4, 18
- [12] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. Retinaface: Single-shot multi-level face localisation in the wild. In *CVPR*, 2020. 4, 18
- [13] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. CogView: Mastering text-to-image generation via transformers. In *NeurIPS*, 2021. 3
- [14] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. In *NeurIPS*, 2022. 3
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3
- [16] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *CVPR*, 2021. 3
- [17] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *ICML*, 2024. 3
- [18] Weichen Fan, Chenyang Si, Junhao Song, Zhenyu Yang, Yinan He, Long Zhuo, Ziqi Huang, Ziyue Dong, Jingwen He, Dongwei Pan, et al. Vchitect-2.0: Parallel transformer for scaling up video diffusion models. *arXiv preprint arXiv:2501.08453*, 2025. 3
- [19] Guian Fang, Wenbiao Yan, Yuanfan Guo, Jianhua Han, Zuntao Jiang, Hang Xu, Shengcai Liao, and Xiaodan Liang. Humanrefiner: Benchmarking abnormal human generation and refining with coarse-to-fine pose-reversible guidance. In *European Conference on Computer Vision*, pages 201–217. Springer, 2024. 18
- [20] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 1981. 5, 14

- [21] Songwei Ge, Seungjun Nah, Guilin Liu, Tyler Poon, Andrew Tao, Bryan Catanzaro, David Jacobs, Jia-Bin Huang, Ming-Yu Liu, and Yogesh Balaji. Preserve your own correlation: A noise prior for video diffusion models. In *ICCV*, 2023. 3
- [22] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv preprint arxiv:2307.10373*, 2023. 2
- [23] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. In *ICLR*, 2024. 2, 3
- [24] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity video generation with arbitrary lengths. *arXiv preprint arXiv:2211.13221*, 2022. 3
- [25] Yingqing He, Menghan Xia, Haoxin Chen, Xiaodong Cun, Yuan Gong, Jinbo Xing, Yong Zhang, Xintao Wang, Chao Weng, Ying Shan, et al. Animate-a-story: Storytelling with retrieval-augmented video generation. *arXiv preprint arXiv:2307.06940*, 2023. 2
- [26] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 3
- [27] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 3
- [28] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 3
- [29] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. CogVideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022. 2, 3, 18
- [30] Jiahui Huang, Leonid Sigal, Kwang Moo Yi, Oliver Wang, and Joon-Young Lee. Inve: Interactive neural video editing. *arXiv preprint arXiv:2307.07663*, 2023. 2
- [31] Ziqi Huang, Kelvin C.K. Chan, Yuming Jiang, and Ziwei Liu. Collaborative diffusion for multi-modal face generation and editing. In *CVPR*, 2023. 3
- [32] Ziqi Huang, Yanan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *CVPR*, 2024. 2, 3, 7, 20
- [33] Ziqi Huang, Fan Zhang, Xiaojie Xu, Yanan He, Jiashuo Yu, Ziyue Dong, Qianli Ma, Nattapol Chanpaisit, Chenyang Si, Yuming Jiang, Yaohui Wang, Xinyuan Chen, Ying-Cong Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. Vbench++: Comprehensive and versatile benchmark suite for video generative models. *arXiv preprint arXiv:2411.13503*, 2024. 2, 3, 5, 7, 20
- [34] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. Co-tracker: It is better to track together. In *ECCV*, 2024. 5, 16
- [35] Johanna Karras, Aleksander Holynski, Ting-Chun Wang, and Ira Kemelmacher-Shlizerman. Dreampose: Fashion image-to-video synthesis via stable diffusion. *arXiv preprint arXiv:2304.06025*, 2023. 2
- [36] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *arXiv preprint arXiv:2303.13439*, 2023. 3
- [37] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3
- [38] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. *arXiv preprint arXiv:2212.04488*, 2022. 2
- [39] Yao-Chih Lee, Ji-Ze Genevieve Jang Jang, Yi-Ting Chen, Elizabeth Qiu, and Jia-Bin Huang. Shape-aware text-driven layered video editing demo. *arXiv preprint arXiv:2301.13173*, 2023. 2
- [40] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895*, 2024. 13
- [41] Xuanyi Li, Daquan Zhou, Chenxu Zhang, Shaodong Wei, Qibin Hou, and Ming-Ming Cheng. Sora generates videos with stunning geometrical consistency. *arXiv preprint arXiv:2402.17403*, 2024. 5, 14
- [42] Jun Hao Liew, Hanshu Yan, Jianfeng Zhang, Zhongcong Xu, and Jiashi Feng. Magicedit: High-fidelity and temporally coherent video editing. *arXiv preprint arXiv:2308.14749*, 2023. 2
- [43] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 18
- [44] Shaoteng Liu, Yuechen Zhang, Wenbo Li, Zhe Lin, and Jiaya Jia. Video-p2p: Video editing with cross-attention control. *arXiv preprint arXiv:2303.04761*, 2023. 2
- [45] Yuanxin Liu, Lei Li, Shuhuai Ren, Rundong Gao, Shicheng Li, Sishuo Chen, Xu Sun, and Lu Hou. Fetv: A benchmark for fine-grained evaluation of open-domain text-to-video generation. In *NeurIPS*, 2023. 3
- [46] Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu, Tiejong Zeng, Raymond Chan, and Ying Shan. Evalcrafter: Benchmarking and evaluating large video generation models. In *CVPR*, 2024. 2, 3
- [47] G Lowe. Sift-the scale invariant feature transform. *Int. J.*, 2004. 5, 14
- [48] Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao, Jingren Zhou, and Tieniu Tan. Videofusion: Decomposed diffusion models for high-quality video generation. In *CVPR*, 2023. 3
- [49] 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, 5, 13, 19

- [50] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhong-gang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 3
- [51] Marius Muja and David G Lowe. Fast approximate nearest neighbors with automatic algorithm configuration. *VISAPP (1)*, 2009. 5, 14
- [52] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 13
- [53] OpenAI. Sora. Accessed February 15, 2024 [Online] <https://sora.com/library>, 2024. 2, 3, 6, 7, 21
- [54] William Peebles and Saining Xie. Scalable diffusion models with transformers. *arXiv preprint arXiv:2212.09748*, 2022. 3
- [55] Xiangyu Peng, Zangwei Zheng, Chenhui Shen, Tom Young, Xinying Guo, Binluo Wang, Hang Xu, Hongxin Liu, Mingyan Jiang, Wenjun Li, Yuhui Wang, Anbang Ye, Gang Ren, Qianran Ma, Wanying Liang, Xiang Lian, Xiwen Wu, Yuting Zhong, Zhuangyan Li, Chaoyu Gong, Guojun Lei, Leijun Cheng, Limin Zhang, Minghao Li, Ruijie Zhang, Silan Hu, Shijie Huang, Xiaokang Wang, Yuanheng Zhao, Yuqi Wang, Ziang Wei, and Yang You. Open-sora 2.0: Training a commercial-level video generation model in \$200k. *arXiv preprint arXiv:2503.09642*, 2025. 3
- [56] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 3
- [57] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, David Yan, Dhruv Choudhary, Dingkang Wang, Geet Sethi, Guan Pang, Haoyu Ma, Ishan Misra, Ji Hou, Jialiang Wang, Kiran Jagadeesh, Kunpeng Li, Luxin Zhang, Mannat Singh, Mary Williamson, Matt Le, Matthew Yu, Mitesh Kumar Singh, Peizhao Zhang, Peter Vajda, Quentin Duval, Rohit Girdhar, Roshan Sumbaly, Sai Saketh Rambhatla, Sam Tsai, Samaneh Azadi, Samyak Datta, Sanyuan Chen, Sean Bell, Sharadh Ramaswamy, Shelly Sheynin, Siddharth Bhattacharya, Simran Motwani, Tao Xu, Tianhe Li, Tingbo Hou, Wei-Ning Hsu, Xi Yin, Xiaoliang Dai, Yaniv Taigman, Yaqiao Luo, Yen-Cheng Liu, Yi-Chiao Wu, Yue Zhao, Yuval Kirstain, Zecheng He, Zijian He, Albert Pumarola, Ali Thabet, Arsiom Sanakoyeu, Arun Mallya, Baishan Guo, Boris Araya, Breana Kerr, Carleigh Wood, Ce Liu, Cen Peng, Dmitry Vengertsev, Edgar Schonfeld, Elliot Blanchard, Felix Juefei-Xu, Fraylie Nord, Jeff Liang, John Hoffman, Jonas Kohler, Kaolin Fire, Karthik Sivakumar, Lawrence Chen, Licheng Yu, Luya Gao, Markos Georgopoulos, Rashel Moritz, Sara K. Sampson, Shikai Li, Simone Parmeggiani, Steve Fine, Tara Fowler, Vladan Petrovic, and Yuming Du. Movie gen: A cast of media foundation models, 2025. 3
- [58] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing. *arXiv preprint arXiv:2303.09535*, 2023. 2
- [59] runway. Gen-3. Accessed June 17, 2024 [Online] <https://runwayml.com/research/introducing-gen-3-alpha>, 2024. 2, 3
- [60] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and Xi Chen. Improved techniques for training gans. In *NeurIPS*, 2016. 3
- [61] Chenyang Si, Weichen Fan, Zhengyao Lv, Ziqi Huang, Yu Qiao, and Ziwei Liu. Repvideo: Rethinking cross-layer representation for video generation. *arXiv 2501.08994*, 2025. 3
- [62] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 4, 14
- [63] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 3
- [64] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021.
- [65] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021. 3
- [66] 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. *arXiv preprint arXiv:2407.14505*, 2024. 3
- [67] Xingwu Sun, Yanfeng Chen, Yiqing Huang, Ruobing Xie, Jiaqi Zhu, Kai Zhang, Shuaipeng Li, Zhen Yang, Jonny Han, Xiaobo Shu, et al. Hunyuan-large: An open-source moe model with 52 billion activated parameters by tencent. *arXiv preprint arXiv:2411.02265*, 2024. 22
- [68] Google Team. Veo2. Accessed December 18, 2024 [Online] <https://deepmind.google/technologies/veo/veo-2/>, 2025. 2, 3
- [69] Kuaishou Team. Kling. Accessed December 9, 2024 [Online] <https://klingai.kuaishou.com/>, 2024. 2, 3, 6, 7, 21
- [70] Minmax Team. Minmax. Accessed August 31, 2024 [Online] <https://hailuoai.com/>, 2023.
- [71] StepFun Team, 2025.
- [72] Tencent Team. Hunyuanvideo: A systematic framework for large video generative models, 2024. 2, 3, 6, 7, 18, 21, 22
- [73] Wan Team. Wan: Open and advanced large-scale video generative models, 2025. 3
- [74] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020. 14
- [75] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020. 5
- [76] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. FVD: A new metric for video generation. In *ICLRW*, 2019. 3
- [77] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. In *NeurIPS*, 2017. 3
- [78] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023. 2, 3

- [79] Wen Wang, kangyang Xie, Zide Liu, Hao Chen, Yue Cao, Xinlong Wang, and Chunhua Shen. Zero-shot video editing using off-the-shelf image diffusion models. *arXiv preprint arXiv:2303.17599*, 2023. 2
- [80] Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent diffusion models. *arXiv preprint arXiv:2309.15103*, 2023. 2, 3
- [81] Yiping Wang, Xuehai He, Kuan Wang, Luyao Ma, Jianwei Yang, Shuohang Wang, Simon Shaolei Du, and Yelong Shen. Is your world simulator a good story presenter? a consecutive events-based benchmark for future long video generation. *arXiv preprint arXiv:2412.16211*, 2024. 3
- [82] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9653–9663, 2022. 4, 18
- [83] Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Xintao Wang, Tien-Tsin Wong, and Ying Shan. Dynamicrafter: Animating open-domain images with video diffusion priors. *arXiv preprint arXiv:2310.12190*, 2023. 2
- [84] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024. 16
- [85] Xiangpeng Yang, Linchao Zhu, Hehe Fan, and Yi Yang. Videograin: Modulating space-time attention for multi-grained video editing. In *ICLR*, 2025. 2
- [86] 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. 3, 6, 7, 18, 21, 22
- [87] Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama, Han Zhang, Huiwen Chang, Alexander G Hauptmann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, et al. Magvit: Masked generative video transformer. In *CVPR*, 2023. 3
- [88] David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation. *arXiv preprint arXiv:2309.15818*, 2023. 3
- [89] Fan Zhang, Shulin Tian, Ziqi Huang, Yu Qiao, and Ziwei Liu. Evaluation agent: Efficient and promptable evaluation framework for visual generative models. *arXiv preprint arXiv:2412.09645*, 2024. 3
- [90] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 3
- [91] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. 13
- [92] Zicheng Zhang, Bonan Li, Xuecheng Nie, Congying Han, Tiande Guo, and Luoqi Liu. Towards consistent video editing with text-to-image diffusion models. *arXiv preprint arXiv:2305.17431*, 2023. 2
- [93] Min Zhao, Rongzhen Wang, Fan Bao, Chongxuan Li, and Jun Zhu. Controlvideo: Adding conditional control for one shot text-to-video editing. *arXiv preprint arXiv:2305.17098*, 2023.
- [94] Yuyang Zhao, Enze Xie, Lanqing Hong, Zhenguo Li, and Gim Hee Lee. Make-a-protagonist: Generic video editing with an ensemble of experts. *arXiv preprint arXiv:2305.08850*, 2023. 2
- [95] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2023. 3

VBench-2.0: Advancing Video Generation Benchmark Suite for Intrinsic Faithfulness

Supplementary Material

In this supplementary file, we provide comprehensive details of VBench-2.0. Section G details each evaluation dimension, including definitions, full evaluation pipelines, and implementation specifics, with illustrative examples. Section H introduces the prompt suite for each dimension, outlining design principles and examples. Section I describes the human annotation process, validating the alignment of VBench-2.0’s automated evaluation with human perception. Finally, Section J covers additional implementation details on video generative models’ sampling.

G. Details on Evaluation Dimension and Method Suite

G.1. Physics

State Change. For a video generative model to simulate the physical world, it is important to reason the text prompt under basic physical law of the real world (*i.e.*, gravity, buoyancy, stress, *etc.* in *Mechanics*; solidification, melting, deposition, *etc.* in *Thermotics*; hardness, solubility, combustibility, *etc.* in *Material*) We call it *physical state change ability*. The existing method, PhyGenBench [49], designs the prompt with the initial state of a physical law and generates four text results (containing specific physical concepts) for one prompt with different levels of correctness based on GPT-4o, where one of them is absolutely correct descriptions while others have different extend of wrong visual phenomena, then they use VLM [40] to select the best option through a multiple-choice format to judge the quality of the generated video by checking whether the VLM selects the best option. However, we experimentally observe that at the current stage, none of the VLM can effectively understand specific physical concepts in natural language terminologies like “microgravity”, nor can they adequately distinguish the differences between options of varying levels. To handle these problems, we first modify the system prompt for GPT-4o [52] to leave the visible phenomena only while removing the specific physical concepts in true results (*i.e.*, unlike PhyGenBench that uses terminologies to describe the physical phenomena “*The liquid’s behavior aligns with the microgravity environment, floating freely, and forming natural blobs without noticeable distortion.*”, we describe the physical phenomena’s visual results “*The liquid floating, and forming blobs*”). We name it Q_{visual}

Then we design a more robust evaluation method with two steps based on video question answering (VQA): pre-filtering and explicit visual phenomena questioning. The



Figure S6. Example for Mechanics.

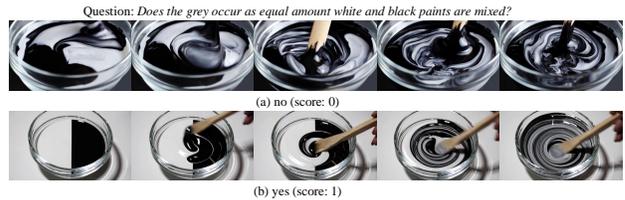


Figure S7. Example for Material.

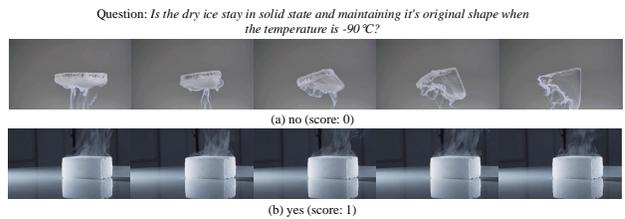


Figure S8. Example for Thermotics.

reason for pre-filtering is that the generation of correct physical phenomena depends on the model’s ability to generate the correct initial state from the given prompt, so we first use LLaVA-video-7B [91] to apply the VQA to filter the unsatisfied videos with a question Q_{filter} that describes the initial state in text prompt (*e.g.*, for the microgravity, the question is that “*Is the environment in space?*” to avoid the situation that a glass of water is thrown into the pool, and the water starts bubbling up, which aligns with the phenomena in space). The full pipeline is as follows:

$$Answer = \text{VQA}(Q_{visual}, V | \text{VQA}(Q_{filter}, V)), \quad (S3)$$

where V is the generated video, only the video matches the Q_{filter} will be considered as the valid case. We show our evaluation results of *Mechanics*, *Material*, and *Thermotics* in Fig S6, Fig S7, Fig S8.

Multi-View Consistency. *Multi-View Consistency* aims to judge the 3D consistency along the video, which is an important evaluation dimension for video generation, especially in the field of 3D/4D reconstruction. However, this dimension could not be evaluated at the pure 3D level as we do not have 3D ground truth in text-to-video generation. We evaluate it from two 2D perspectives to simulate it: the valid



Figure S9. Example for Multi-View Consistency.

number of feature matching between frames (*i.e.*, Feature matching captures the geometric consistency of a video by matching local geometric features (e.g., corners, textures) across frames, and in dynamic scenes, this process directly reflects geometric consistency through the stable alignment of object positions, shapes, and motion trajectories. We follow [41] to utilize it as part of the evaluation methods) and the moving speed of the camera. The consideration here is that one model could generate a fast camera motion with high geometric consistency is better.

For camera moving speed, we use the RAFT [74] to calculate the global flow score F_{score} between frames. Specifically, we first down-sample the videos generated by different models to the smallest one (*i.e.*, Sora with size of 480p) to avoid extensive artifacts and use different sampling intervals based on different FPS in generated videos to uniform the scale of flow score. Note that we discard examples with optical flow scores less than 5. This is because, in such cases, the video is almost stationary but the number of matched feature points becomes abnormally high, providing no value for judging multi-view consistency. We also set the maximum score to 30 as no video exceeds this score.

For feature matching, the sampling intervals $S_{interval}$ are related to two key elements: FPS and the camera moving speed (*i.e.*, the faster the camera moving, the less the point will be matched). So we balance the $S_{interval}$ as follows:

$$S_{interval} = S_{fix} / [(FPS/8) * F_{score}/10]; \quad (S4)$$

where S_{fix} is a default interval, which is set as 40. In this way, we omit the influence of FPS and camera moving speed in different methods and uniform the feature matching scale.

After obtaining the sampling interval, we first use SIFT [47] to extract local feature points from each frame in the video, then utilize FLANN [51] to efficiently match features between nearby frames, finally we use RANSAC [20] to filter the incorrect matches from the feature matching process.

Then based on some perfect cases of showing geometry consistency, we manually define the largest valid matchable point as 750, largest flow score as 10 (the 30 used before is to scale the sampling interval). We use the upper bound to rescale the valid points and flow score to [0, 1] and multiply them as the final score. The downsampling is applied before



(a) B has higher diversity (score A: 0.34 < score B: 0.70)

Figure S10. Example for Diversity.

both the flow scale step and the feature matching step to ensure fair comparison. As shown in Fig S9, we present an evaluation example of *Multi-View Consistency*.

G.2. Creativity

Diversity. *Diversity* is a fundamental ability of generative models, which measure the variety of content given the same prompt input. In video generation, this ability is also important for diverse content creations. Here we focus on the diversity of the various dimensions we proposed (*e.g.*, *Composition, Motion Order Understanding, Human Interaction, Physics, etc.*), aiming to evaluate the model’s diversity in different semantic levels. Specifically, we decouple the diversity here into style diversity and content diversity. Suppose the user requires generating 20 diversified content with a single prompt, we calculate frame-wise style diversity as follows:

$$S_{diff} = \frac{1}{N} \sum_{i < j} \sum_{l=1}^5 \|G_l^{(i)} - G_l^{(j)}\|_2^2, \quad (S5)$$

$$G_l^{(i)} = \frac{1}{CHW} F_l^{(i)T} F_l^{(i)} \quad l \in \{1_1, 2_1, 3_1, 4_1, 5_1\}, \quad (S6)$$

where F_l^i is the feature of i^{th} sample of convolution layer l in pretrained VGG19 [62]. G is the Gram matrix (popular in style transfer to express the style information), C, H, W are the channel, height, width of the style feature and N is the number of sample pairs. We calculate the average among all the pair samples as the style diversity score S_{diff} . As for the content diversity, we use the conv4.2 layer in VGG19 as the content feature and calculate the L1 loss as follows:

$$C_{diff} = \frac{1}{N} \sum_{i < j} L_1(F^{(i)} - F^{(j)}). \quad (S7)$$

The final diversity score is obtained as follows:

$$Score = \lambda S_{diff} + C_{diff}, \quad (S8)$$

where the λ is set to 1000 followed by style transfer. Finally, we set the upper bound of the Score as the highest score in all the cases (17.712) and rescale the score to [0,1].

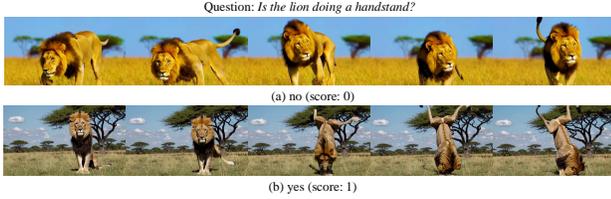


Figure S11. Example for Composition.

As shown in Fig S10, we present an evaluation example of *Diversity*.

Composition. In this paper, we evaluate it in three levels with different difficulties: species composition, single-entity action and multi-entity task. *Species combination* focuses on whether models could generate a new creature that combines the feature of several existing creatures. *Single-entity action* concentrates on whether generation models could combine a creature or object with an action in a way that is impossible to occur in real life. As for *multi-entity action*, we evaluate whether generation models can combine multi-entity to collectively perform a task.

For all of the three levels, we use a unified evaluation method: *video-based multi-question answering*. Suppose the user wants to generate a video with species composition prompt, we first decompose the text prompt into several questions Q which contains the detailed species part only, then we use the LLaVA-video-7B to perform VQA with ‘yes’ or ‘no’ answer as:

$$Score = \frac{1}{N} \sum_{i=1}^N VQA(Q^i, V), \quad (S9)$$

where V is the generated video, N is the number of questions for each prompt and the final score will fall into $[0,1]$. We additionally observe that pure VQA can not ensure robust results as current video generation models will generate separate species for the given prompt and VQA is a discrete process that could not consider the separate situations. So in this paper, we will omit the cases that show more than one creature in the video by a pre-VQA process with the question “*Is there only one creature in the video?*”.

Commonly speaking, the evaluation of composition dimension could follow the one used in basic physical law: text-description alignment. However, we experimentally observe that LLaVA-video could not describe the combined species accurately and usually misunderstand several parts of the creatures. We believe this is due to the fact that such creative topics cannot be found in real-world data, and we hope this will also provide some help for the future development of VLM. As shown in Fig S11, we present an evaluation example of *Composition*.

G.3. Controllability

Dynamic Spatial Relationship. Here, we explore the understanding of spatial relationship and the mobility capabil-

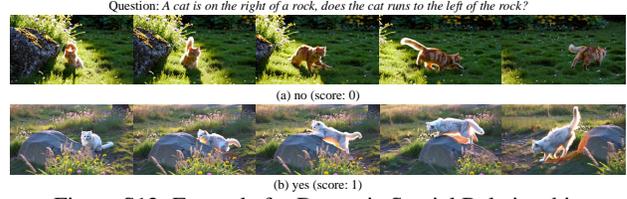


Figure S12. Example for Dynamic Spatial Relationship.

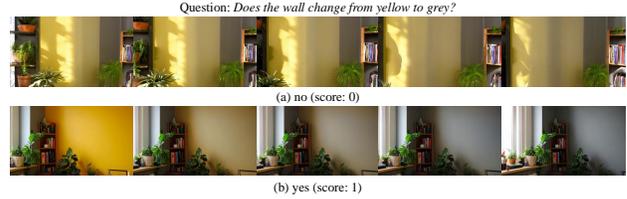


Figure S13. Example for Dynamic Attribute.

ities of video generation models. For each prompt, there is one creature and one object with an initial and a final spatial relationship (e.g., A dog is on the right of an apple, then the dog moves to the left of the apple.). We evaluate it by VQA, each prompt will be divided into two questions: the initial and the final spatial relationship and only the video matches both questions will be considered as correct. Here we do not consider whether the creature truly moves or just generates two creatures that satisfy the initial and final spatial relationship simultaneously. The reason is that the movement is too complex to judge in current video generation videos and we leave it as a future work. As shown in Fig S12, we present an evaluation example of *Dynamic Spatial Relationship*.

Dynamic Attribute. Here we explore whether the video generation models can change a creature’s or object’s attribute following text prompt. We consider four common attributes: Lightness, color, size and material. All of the attributes use a unified evaluation method, VQA. Each text prompt will be decomposed into three questions: the initial and final state of the object attribute and a justification of whether the attribute truly changed or just generate both states of the object at the same time. We consider the attribute changing here as we experimentally observe that LLaVA-video-7B can recognize the changing process of these four attributes precisely. The same as Dynamic Spatial Relationship, only all of the three questions are satisfied will be considered as correct. As shown in Fig S13, we present an evaluation example of *Dynamic Attribute*.

Motion Order Understanding. Here we explore the action order understanding ability of current video generation models. Each text prompt contains two actions, which are coordinated by a temporal adverb. We use text-description alignment to evaluate it. Specifically, we first split the text prompt into two descriptions with action only in turn, then we obtain the video caption by LLaVA-video-7B with an additional system prompt “*Return the action order in video. Here is the template: ‘1. ; 2. .’*”. The reason here is that LLaVA-video has a strong ability to generalize temporal ac-

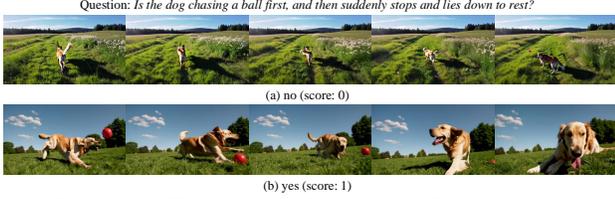


Figure S14. Example for Motion Order Understanding.

tions. After obtaining the captions, we assess whether both actions match the ground truth descriptions in turn; only when both actions match are they considered correct. This approach cleverly avoids the hallucination problem of LLM in order judgment. As shown in Fig S14, we present an evaluation example of *Motion Order Understanding*.

Human Interaction. Human interaction aims to evaluate whether the video generation models could truly generate the interaction mentioned in the text prompt. We adopt *text-description alignment* pipeline to evaluate it, which will first utilize VLM to capture the different levels of video captions (*i.e.*, up to the difficulty of dimensions) and utilize LLM to judge the consistency between captions and text prompt. Specifically, for each video, we first use the LLaVA-video-7B to obtain two levels of video captions: dense caption T_{dense} with the system prompt “Describe the video in detail.” and action interaction caption T_{action} with the system prompt “Describe the human interaction in the video, following the template as [a person xx to another person.]”. The insight here is that LLaVA-video can effectively summarize the action interactions into a standard format, which helps eliminate the influence of background and details when feeding the input into the LLM for text matching. However, due to the misleading nature of our input format, LLaVA-video transforms all instances into this standard format regardless of how many people appear in the video, which leads to misinterpretations, such as treating situations with only one person as interactions between two individuals. This is why we need an extensive dense video caption, which is used to detect the videos that only contain one person. Specifically, we design a system prompt S_{num} as shown below to prompt Qwen-2.5-7B-Instruct with the human number judgment criteria and obtain the final result.

Qwen system prompt S_{num} for human number judgment

You are Qwen, created by Alibaba Cloud. You are a helpful assistant and a brilliant person number judge.

You need to judge whether the description contains more than one person. Return yes or no only.

The formula is as follows:

$$\begin{aligned} \text{Response} = & \text{LLM}(T_{dense}|S_{num}) \\ & \text{and } \text{LLM}(T_{action}, T_{prompt}|S_{align}), \end{aligned} \quad (\text{S10})$$

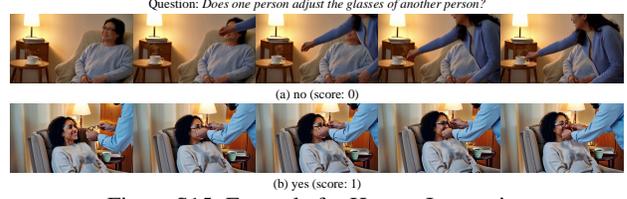


Figure S15. Example for Human Interaction.

where LLM is the text alignment judge (answer yes or no only) and we use Qwen2.5-7B-Instruct [84] in this paper by default and S_{align} is the system prompt to prompt the LLM with the ability of text alignment. Note that only the video that matches both of the questions will be considered as correct.

Commonly speaking, using a detection model could also handle the problem of human numbers. However, we experimentally observe that current video generation models usually generate videos with part of people (*e.g.*, for “a person is helping another person comb their hair”, one person may only show their hand.). The detection model could not handle this situation while this could be involved in the caption by LLaVA-video and can be detected by Qwen. As shown in Fig S15, we present an evaluation example of *Human Interaction*

Camera Motion. In films, camera movement is a fundamental operation, and therefore it is also very important for video generation models. VBench-2.0 considers nine basic camera movements (*i.e.*, pan left/right, tilt up/down, zoom in/out, static, orbit, and oblique aerial view) to systematically evaluate the ability of current and future video generation models. We use tracking technology (CoTracker2 [34]) to trace the edge regions, and determine the camera’s motion by analyzing the positions of points on the top, bottom, left, and right edges at the start and end moments. Specifically, we use uniform sampling of grid points to set the tracking points in the first frame with grid size 10. Besides the “orbit” type of movement, for the other eight types, we only use the position changes of the central points in the outermost edge regions of each side at the start and end moments of the video to determine the camera motion type. For example, if the camera movement is pan left, all of the tracking points will move to the right, so it is only necessary to determine whether the point at the end moment is on the right of those at the start moment. As for “orbit”, due to the uncertainty in the camera movement trajectory of the current video model, we make a judgment every 20 frames (*i.e.*, in each column of the grid, the first and last points move in opposite directions, with minimal movement in the y-direction and significant movement in the x-direction.), and if it appears at least once in the video, it is considered correct. As illustrated in Fig S17, we present an example of camera zoom-in motion.

Complex Plot. To meet the demands of future movie-

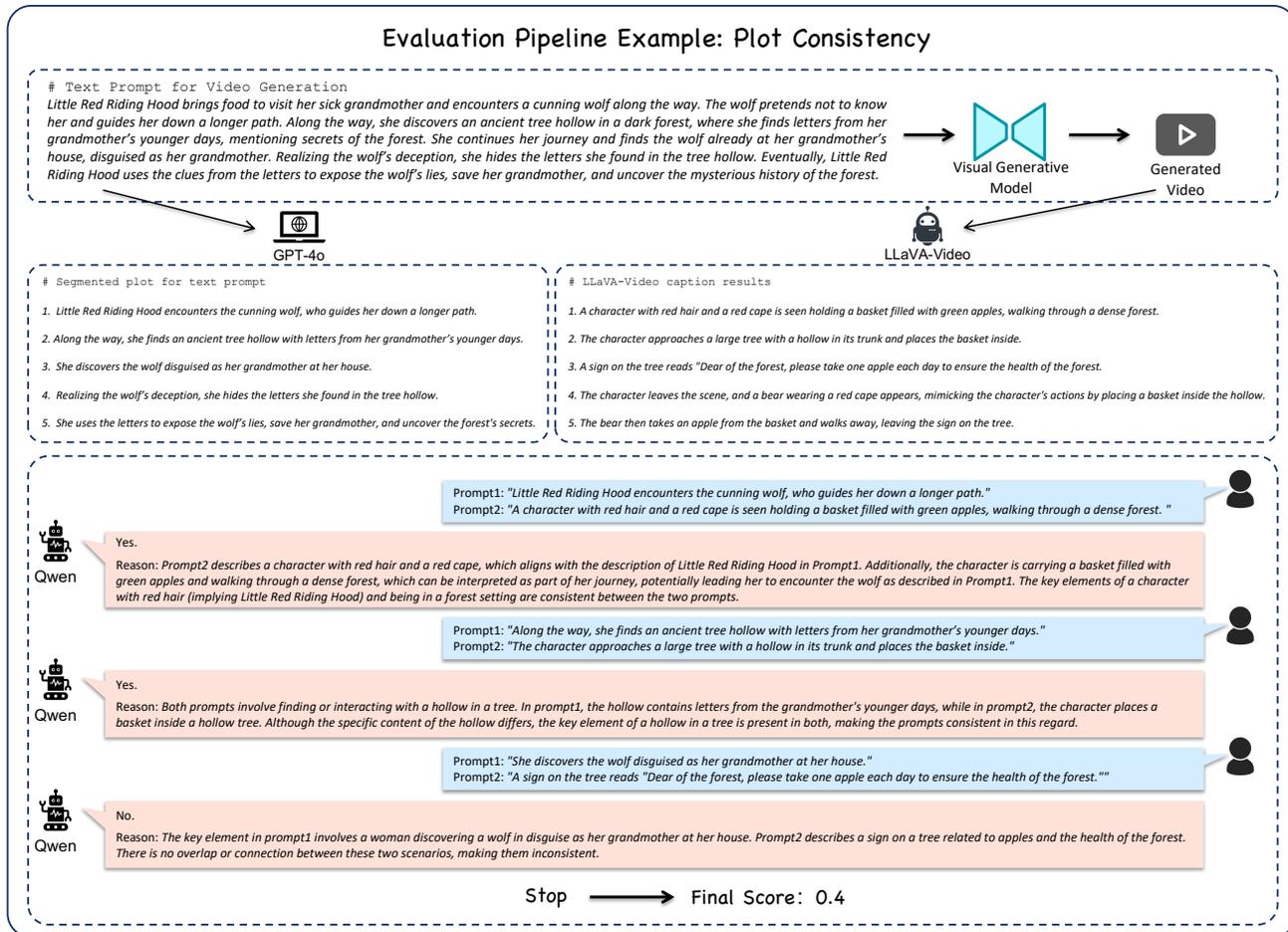


Figure S16. Example for Complex Plot.

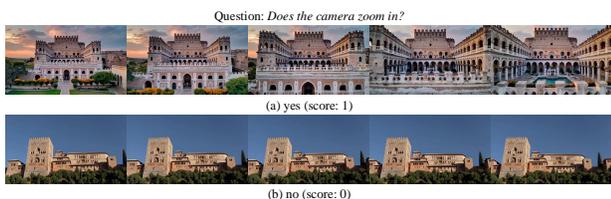


Figure S17. Example of Dynamic Camera Motion.

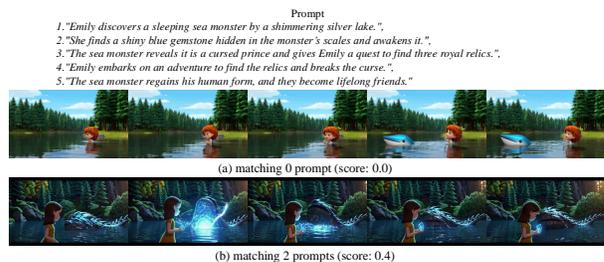


Figure S18. Example for Plot Consistency.

level video generation, the ability to handle multi-plot storytelling is a crucial capability. As shown in Fig S18, we present an evaluation example of a complex plot. We utilize *text description alignment* to evaluate it and the formula is similar to Eq.S10. Specifically, we first split and summarize the long text prompt into 4 or 5 plot descriptions as ground truth. Then we feed the video into LLaVA-video-7B to obtain 4 or 5 video captions with the system prompt "Return the plot in video. Here is the template: [1. ; 2. ; 3. ; 4. .]". Finally, we sequentially match the corresponding text and captions in the order of the plot. If the LLM (Qwen2.5-7B-Instruct) determines that a plot element appears in the caption, we proceed to evaluate the next plot element; otherwise, the evaluation stops. We discretize the

plot elements to calculate the accurate score, simulating the effect of LLM scoring. For better understanding, we show a complete example in Fig S16

Complex Landscape. Long-duration close-ups of scenery are another essential element in films. Similarly, we have designed five landscape descriptions for each text prompt, with the transitions between the scenes achieved through camera movements. The evaluation is completely the same as *Plot Consistency*. As shown in Fig S19, we present an evaluation example of a *Complex Landscape*.

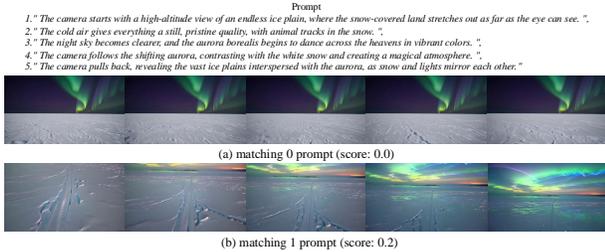


Figure S19. Example for Complex Landscape.

G.4. Human Fidelity

Human Anatomy. This section focuses on detecting potential anomalies in the appearance and structure of the human body in videos, evaluated at two levels of granularity. At the coarse level, we address global anomalies such as extra or missing limbs, body penetration, unreasonable poses, etc. At the fine-grained level, we specifically assess anomalies in hands and faces, including abnormal finger counts, facial distortions, etc. Our methodology leverages a pre-trained ViT-base model [82] to train three anomaly detection models, which will detect the human in each frame and judge the anomaly score of each human, targeting the human body, hands, and faces, respectively. The classification token (`cls_token`) is passed through an MLP to produce binary outputs. The dataset consists of two components:

1. *Real samples*: We collected approximately 1,000 real motion videos from the web and used the YOLO-World [8] detection model to extract patches of humans, hands, and faces as positive samples.
2. *Generated samples*: We generated around 1,000 videos related to human motion using CogVideo [29, 86] and HunyuanVideo [72]. These were processed with YOLO-World to extract image patches, followed by meticulous manual annotations for training. Additionally, negative samples were sampled from the HumanRefiner [19] dataset, focusing on patches related to humans, hands, and faces.

As a result, the number of training samples for humans, hands, and faces is approximately 80K, 30K, and 40K, respectively. During training, we used a batch size of 128, a learning rate of 1×10^{-4} , and the AdamW optimizer for 30 epochs. To mitigate data imbalance, we employed focal loss [43]. During inference, we first applied the YOLO-World model to detect all human instances in the input video. For each detected human region, we further detected hands and faces, extracting corresponding patches as inputs. A human instance is flagged as abnormal if any of the three models predict an anomaly. The final score is computed as

$$Score = \frac{c_{normal}}{c_{normal} + c_{abnormal}},$$

where c_{normal} and $c_{abnormal}$ represent the counts of normal and abnormal humans in the video, respectively. In practice,



Figure S20. Example for Human Anatomy.



Figure S21. Example for Human Identity.

we empirically set the YOLO-World detection threshold to 0.1 and the anomaly detection thresholds for humans faces, and hands to 0.45, 0.30, and 0.32, respectively. As shown in Fig S20, we present an evaluation example of *Human Anatomy*.

Human Identity. Here we evaluate the *Human Identity* consistency ability of video generation models. We focus on the scene with only one person to obtain the best evaluation accuracy. Firstly we utilize RetinaFace [12] to detect the human face in each frame of the video and then use Arcface [11] to extract the face feature. We select the first frame of the video as the anchor by default, and all subsequent frames are compared to the first frame to calculate similarity. However, time periods where there are multiple people or no one present are not taken into consideration. **Note** that after a scene change, whether it is still the same person is taken into consideration, which represents a more challenging scenario. As shown in Fig S21, we present an evaluation example of *Human Identity*.

Human Clothes. *Human Clothes* evaluates the human temporal consistency in another perspective. We use *video-based multi-question answering* to handle it. Specifically, we design three unified questions for all videos to exclude all abnormal situations (i.e., “*Is there only one person in the video throughout?*”, “*Is the person in the video the same throughout?*”, and “*Does the clothes of the person in the video (color, texture) remain consistent throughout?*”). Only the video that matches all the questions will be considered as correct.

Note that we do not use traditional methods such as feature similarity calculation to handle this problem based on two considerations: 1) Current models often generate characters with their clothing partially obscured by objects or with changes in visible body parts (e.g., the upper body transforming into the lower body). In such cases, traditional algorithms are unable to address these cross-temporal judgment issues. 2) LLaVA-video-7B has strong clothing color



Figure S22. Example for Human Clothes.

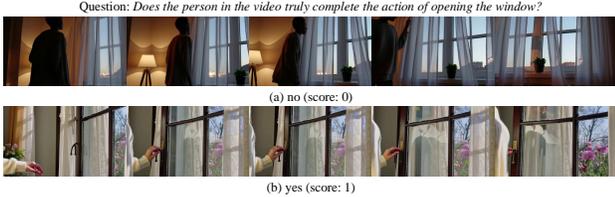


Figure S23. Example for Motion Rationality.

perception and a certain degree of cross-temporal memory capability. Therefore, we use it to evaluate this dimension. As shown in Fig S22, we present an evaluation example of *Human Clothes*.

G.5. Commonsense

Motion Rationality. *Motion Rationality* aims to evaluate whether the generated video truly acts the motion. For example, when inputting a prompt about eating hamburger, the current video generation models may generate a person pretending to eat a hamburger, but in reality, the hamburger shows no bite marks. This issue includes fake eating, fake drinking, fake cutting, fake walking, and so on, which is an urgent problem that needs to be detected and solved but cannot be done by current evaluation metrics. To evaluate it, we adopt the *video-based multi-question answering*, which is similar to Eqn. S9. Specifically, we generate several questions for a text prompt by GPT-4o and we will check whether the GPT-4o generates questions unsuitable for the given prompt. We show a question case below.

Generated questions for “A person is eating hamburger”

- 1.Does the person appear to be eating a hamburger? (yes or no)
 - 2.Is the person’s mouth in contact with the hamburger? (yes or no)
 - 3.After eating, is the hamburger visibly reduced or divided? (yes or no)
 - 4.Is the hamburger still there after eating? (yes or no)
-

The principle behind our problem design is to eliminate as many potential situations in video generation that could lead to misjudgments by VLM (LLaVA-video-7B) as possible, and to enhance the accuracy of the answers through a certain amount of redundancy (*i.e.*, only one video matches all the questions will be considered as correct). As shown in Fig S23, we present an evaluation example of *Motion Rationality*.

Instance Preservation. *Instance Preservation* aims to evaluate whether video generation models could preserve the



Figure S24. Example for Instance Preservation.

number of instances in different situations (*e.g.*, positional movement, rapid motion, multi-object collisions, and interactions). We utilize open-vocabulary detection model, YOLO-World [8] to cover more complex and diverse situations and calculate frame-level object numbers. In practice, we empirically set the YOLO-World detection threshold to 0.28 for all the classes. As shown in Fig S24, we present an evaluation example of *Instance Preservation*.

H. More Details on Prompt Suite

For each dimension in VBench-2.0, we carefully design around 70 prompts as the test cases and try to consider all the aspects that need to be evaluated as thoroughly as possible. For example, for “Multi-View Consistency”, we provide object-level and scene-level prompts; for “Composition”, we divide it into three aspects: species combination, single-entity action and multi-entity task, *etc.* We detail the prompt suite for each dimension below.

State Change. We basically follow PhyGenBench [49] to divide this dimension into three main aspects and twelve subcategories: *color mixture, flame reaction, hardness, solubility, combustibility, acidity and alkalinity, oxidative dehydration, redox* for **Material**; *gravity, buoyancy, stress, atmosphere pressure, solid pressure, elasticity, friction, surface tension* for **Mechanics**; *solidification, melting, liquefaction, boiling, deposition, sublimation* for **Thermal**. Additionally, we upgrade the difficulty of Thermal by adding specific temperature constraints based on the melting, boiling, and freezing points of different materials, exploring the higher-level semantic capabilities of video generation models or prompt optimizers.

Multi-View Consistency. This dimension is specifically designed for 3D/4D reconstruction, generation in the future. We design our prompts based on the widely used 3D object datasets ShapeNet [6], Objaverse [10].

Diversity. We carefully choose 10 prompts from various evaluation dimensions in VBench-2.0 to fully judge the diversity of existing video generation models in different semantic levels.

Composition. As mentioned above, we divide *Composition* into 3 aspects: species combination, single-entity action and multi-entity task. For all the aspects, the entities and actions we selected are common in real life, avoiding the situation where the model fails to generate results due to the training set not having encountered them.

Dynamic Spatial Relationship. We design each prompt with one common animal and one common object and only the animal could move. Additionally, we design the corresponding directions based on the capabilities of animals. For instance, cats and dogs are incapable of flying, and giraffes cannot physically be positioned beneath certain objects due to their height and anatomical constraints.

Dynamic Attribute. We consider four aspects in this dimension: size, color, lightness and material. We choose common entities with the potential ability to attribute change. For example, ant changes from small to big; leaves change from yellow to red; sky changes from dark to light; and car change from wooden to metal.

Motion Order Understanding. The prompt template of this dimension is “A [subject] is [motion A], then they start [motion B]”. We select human and common animals as the subjects and design motions for each subject which owns the possibility of successful transitions between motions.

Human Interaction. The most critical aspect of this dimension is to ensure that the designed prompts include interactions involving physical contact. For example, a scenario where one person is making a phone call to another cannot determine how the interaction is established. We cover all commonly used interactions in daily life (e.g., “shake hands”, “hand something”, “lift something together”).

Camera Motion. Since camera movement is frequently employed in films, we design prompts that are likely to appear in cinematic contexts, such as gardens (scenes), Mount Fuji (landmarks), and tables (objects).

Complex Plot. We meticulously designed text prompts across various themes to systematically evaluate the story generation capabilities of video generation models. For example, “4x100 relay race”, “fairy tale storyline”, “hero’s adventure theme” and “tomb raiding theme”. Each prompt contains at least 150 words and four or five small plots to describe a complete story, targeting a true movie generation requirement.

Complex Landscape. The design of this dimension is similar to *Complex Plot*, but the focus here is on landscape description. We design several beautiful landscape close-up such as “sea of flowers”, “tropical rainforest” and “snowy mountain”, etc. Each prompt also contains at least 150 words and five small landscape descriptions (connected by camera movement).

Human Anatomy. We divide this dimension into three aspects (action with large motions such as playing ball or running; human interaction and motion order understanding) to thoroughly evaluate the capabilities of different video generation models in this dimension. For action with large motion, we further design several harder cases, which involves multi-people in one prompt (the human abnormal will occur more frequently).

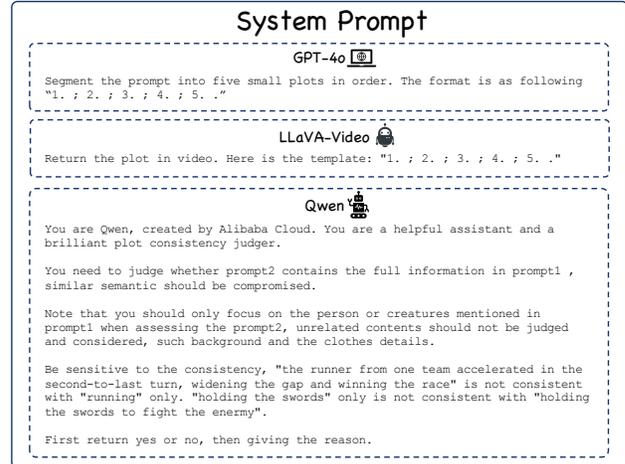


Figure S25. VBench-2.0 System Prompt for different backend models.

Human Consistency. Human consistency contains two levels of dimensions: identity and clothes. Both of them use the prompts from *Human Anatomy* apart from *Human Interaction* part as we focus on single-person consistency.

Motion Rationality. We consider all of the situations (e.g., fake eating, drinking, walking, running, cutting, and fake human object interaction such as turning off the TV) that may occur motion unrealistic and design them into text prompts.

Instance Preservation. We consider all of the situations that may occur abnormal entity division or fusion. For example, action with high motion, multi-entity interaction, object dropping or collision. We select common creatures or objects from daily life as the main subjects and for each subject, we carefully design actions and scenarios that conform to the physical laws of the real world. As we also contain color descriptions, we also select colors that are commonly associated with each subject in real life.

I. Human Preference Annotation Details

Following the approach of VBench [32, 33], we conduct large-scale human preference labeling on generated videos to validate the alignment between VBench-2.0’s evaluation and human perception across all evaluation dimensions. The collected human annotations also serve as a valuable resource for future research on fine-tuning generation and evaluation models to better reflect human judgments.

I.1. Data Preparation

There are two types of annotation formats in VBench-2.0.

The first type follows VBench and the human annotators are tasked to select a preferred video from two generated videos based on specific criteria. Given a text prompt p_i and four video generation models $\{A, B, C, D\}$,

Table S4. **Human Alignment of VBench-2.0 Evaluation Methods.** For each evaluation dimension and each video generative model, we report “VBench-2.0 Win Ratios (left) / Human Win Ratios (right)”. The results demonstrate that our evaluation metrics closely align with human perception across all dimensions.

Models	Human Anatomy	Human Clothes	Human Identity	Composition	Diversity	Mechanics	Material	Thermotics	Multi-View Consistency
HunyuanVideo [72]	67.73% / 58.73%	44.49% / 45.67%	52.60% / 58.19%	46.02% / 40.25%	16.67% / 36.67%	53.30% / 50.64%	45.70% / 45.49%	50.13% / 46.08%	52.83% / 50.39%
CogVideoX-1.5 [86]	13.10% / 5.28%	48.12% / 44.28%	36.81% / 20.62%	50.00% / 46.23%	26.67% / 28.33%	57.57% / 51.20%	59.43% / 57.49%	56.69% / 58.51%	8.77% / 18.07%
Sora [53]	60.05% / 69.71%	55.61% / 57.15%	53.88% / 56.25%	56.08% / 59.96%	93.33% / 83.33%	43.47% / 49.39%	47.38% / 50.79%	40.37% / 41.19%	64.10% / 63.94%
Kling 1.6 [69]	59.12% / 66.36%	51.18% / 52.86%	56.48% / 63.05%	47.90% / 53.56%	63.33% / 51.67%	45.56% / 48.78%	46.65% / 45.90%	52.74% / 54.02%	80.00% / 65.98%
Correlation	95.46%	90.89%	99.46%	81.70%	94.53%	93.56%	93.71%	93.86%	98.44%
Models	Dynamic Spatial Relationship	Dynamic Attribute	Motion Order Understanding	Human Interaction	Complex Landscape	Complex Plot	Camera Motion	Motion Rationality	Instance Preservation
HunyuanVideo [72]	50.64% / 51.37%	52.75% / 55.07%	51.46% / 56.35%	49.72% / 52.45%	52.78% / 60.04%	45.74% / 30.46%	46.60% / 49.28%	49.43% / 52.68%	49.61% / 52.36%
CogVideoX-1.5 [86]	49.36% / 48.07%	53.72% / 54.70%	51.68% / 52.13%	53.28% / 51.34%	55.74% / 54.46%	52.50% / 65.37%	46.19% / 44.08%	49.04% / 42.24%	41.91% / 41.32%
Sora [53]	49.68% / 49.60%	42.98% / 43.16%	43.60% / 30.81%	43.94% / 43.49%	40.56% / 35.87%	50.74% / 49.35%	42.08% / 42.64%	49.43% / 44.83%	55.17% / 52.94%
Kling 1.6 [69]	50.32% / 50.97%	50.55% / 47.07%	53.25% / 60.73%	53.06% / 52.73%	50.93% / 49.63%	51.02% / 54.81%	65.12% / 63.99%	52.11% / 60.25%	53.31% / 53.37%
Correlation	97.36%	90.95%	98.73%	89.93%	91.17%	97.65%	97.95%	87.97%	93.70%

we generate a set of videos, forming a “group” $G_{i,j} = \{V_{i,A,j}, V_{i,B,j}, V_{i,C,j}, V_{i,D,j}\}$. For each prompt p_i , we sample five such groups $\{G_{i,0}, G_{i,1}, G_{i,2}, G_{i,3}, G_{i,4}\}$ and construct pairwise comparisons: (V_A, V_B) , (V_A, V_C) , (V_A, V_D) , (V_B, V_C) , (V_B, V_D) , (V_C, V_D) . Human annotators are asked to select their preferred video for each pair. To ensure unbiased annotations, the video order within each pair is randomized.

The second type might involve two groups of videos, and the dimension in evaluation is related to the content distribution in these two groups of videos. In this type, the pairwise construction is the same as the first type while each V_i contains 20 videos generated by single prompt p_i .

I.2. Labeling Instructions.

The annotation process follows VBench but incorporates refinements in interface design and evaluation methodology. Since VBench-2.0 introduces multiple dimensions with detailed multi-question evaluations and long text prompts (some exceeding 150 words), we enhance the interface for improved readability and efficiency. Instead of displaying extensive text in video titles, we sequentially list all key annotation instructions directly within the interface. This structured layout allows annotators to efficiently reference the necessary details while conducting evaluations.

I.3. Win Ratio.

Given human annotations, we calculate the win ratio for each model. During pairwise comparisons, if a model’s video is preferred, it scores 1, while the other model scores 0. In case of a tie, both models score 0.5. The win ratio for each model is computed as the total score divided by the number of pairwise comparisons. We show the result in Table S4

I.4. Quality Assurance

The annotation process maintains the rigorous quality control measures established in VBench while further refining the review criteria. We randomly sample 20% of the annotated pairs for verification, with a required success rate of 95%

To quantify the annotation effort across all 18 evaluation dimensions, we account for the cumulative time spent on documentation drafting, initial trials, formal annotation, and re-annotation. In total, the process required 15.75 hours of individual effort and 284 hours across 18 annotators, reflecting the scale of the task and our commitment to ensuring high-quality human preference annotations. Most evaluation dimensions go through 2 rounds of trial labeling, official labeling, and post-labeling verification

J. Video Generation Models in Evaluation

To assess our benchmark against recent advancements, we utilize four models for comparison, with more to be included as they become open-sourced. We introduce the four models in detail and a unified prompt refiner to pre-process the text prompt as follows.

Prompt Refiner. As video generation models continue to improve their understanding of text and support longer inputs, many methods have begun adopting prompt rewriting techniques to refine text input. To ensure fair comparisons and higher quality video in this paper, we employ a modified version of the original Prompt Refiner from CogVideoX to uniformly process all input prompts. We tested the following four models and found that except for Sora, our Prompt Refiner consistently result in higher-quality videos with better alignment to the input text. We hypothesize that Sora may have an embedded Prompt Refiner, which could explain its performance. Therefore, we applied our Prompt Refiner to **all dimensions** of the other three models.

Kling 1.6 [69]. We adopt the standard API version of Kling 1.6. For each prompt, we sample 241 continuous frames of size 720×1280 at 24 frames per second (FPS). For the classifier-free guidance (cfg) scale, we used the official default value of 0.5.

Sora-480p [53]. Sora represents one of the earliest large-scale video generation models to emerge. For each prompt, we directly perform sampling on the official website, containing 150 continuous frames of size 480×854 at 30 FPS (*i.e.*, the low-resolution version of Sora).

HunyuanVideo [72]. HunyuanVideo is a powerful open-sourced video generation model. It is equipped with a pre-trained Multimodal Large Language Model (MLLM) [67] to better understand text prompts, supporting a large maximum token length and incorporating a dual-branch diffusion transformer generative architecture, which produces visually convincing results. We use the best version of HunyuanVideo with no classifier-free guidance. For each prompt, we sample 129 continuous frames of size 720×1280 at 24 FPS. All of the hyper-parameters are followed the default value in the official inference code. The initial random seed is set to 42 during sampling.

CogVideoX-1.5 [86]. CogVideoX-1.5 is a transformer-based video generation model, which is equipped with a 3D causal VAE to compress the spatial and temporal dimensions and an expert transformer generative architecture to achieve better text-video alignment. We use the best version of it (5B) to evaluate. For each prompt, we sample 161 continuous frames of size 768×1360 at 16 FPS. The initial random seed is also set to 42 for a fair comparison.