TOWARDS WORLD SIMULATOR: CRAFTING PHYSI CAL COMMONSENSE-BASED BENCHMARK FOR VIDEO GENERATION

Anonymous authors

Paper under double-blind review

Abstract

Text-to-video (T2V) models like Sora have made significant strides in visualizing complex prompts, which is increasingly viewed as a promising path towards constructing the universal world simulator. Cognitive psychologists believe that the foundation for achieving this goal is the ability to understand intuitive physics. However, the capacity of these models to accurately represent intuitive physics remains largely unexplored. To bridge this gap, we introduce PhyGen-Bench, a comprehensive Physics Generation Benchmark designed to evaluate physical commonsense correctness in T2V generation. PhyGenBench comprises 160 carefully crafted prompts across 27 distinct physical laws, spanning four fundamental domains, which could comprehensively assess models' understanding of physical commonsense. Alongside PhyGenBench, we propose a novel evaluation framework called *PhyGenEval*. This framework employs a hierarchical evaluation structure utilizing appropriate advanced vision-language models and large language models to assess physical commonsense. Through Phy-GenBench and PhyGenEval, we can conduct large-scale automated assessments of T2V models' understanding of physical commonsense, which aligns closely with human feedback. Our evaluation results and in-depth analysis demonstrate that current models struggle to generate videos that comply with physical commonsense. Moreover, simply scaling up models or employing prompt engineering techniques is insufficient to fully address the challenges presented by Phy-GenBench (e.g., dynamic physical phenomenons). We hope this study will inspire the community to prioritize the learning of physical commonsense in these models beyond entertainment applications. We release the data and codes at https://github.com/PhyGenBench/PhyGenBench

034 035

037

006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

033

1 INTRODUCTION

038 Text-to-video (T2V) models such as Sora have made significant strides in visualizing complex ideas 039 and scenes based on textual input (Yang et al., 2024; Wang et al., 2023). These advancements are 040 increasingly viewed as a promising path towards constructing universal simulators of the physical 041 world, which holds immense promise for video generation (Zhu et al., 2024), autonomous driving 042 (Gao et al., 2024), and the development of embodied agents (Mazzaglia et al., 2024). Cognitive 043 psychology posits that intuitive physics, which is demonstrated even by human infants (Wood et al., 044 2024; Battaglia et al., 2013), is essential for achieving this goal. Intuitive physics emphasizes rendered scenes should be visually and interactively natural to humans, rather than adhere to strict physical accuracy. Consequently, on the path towards developing a world simulator (Xiang et al., 046 2024), video generation should first be capable of accurately reproducing simple yet fundamental 047 physical phenomenons. However, even state-of-the-art models trained on vast resources (Tan et al., 048 2024) encounter difficulties in correctly generating seemingly trivial physical phenomenons, as depicted in Figure 1, the model fails to understand that the stone should sink in water. This clear pitfall shows a substantial gap between current video generation models' and human's understanding of 051 basic physics. It reveals how far these models are from being true world simulators. 052

Given this context, it becomes important to assess the extent to which current T2V models can capture intuitive physics in their generated outputs. This requires the development of comprehensive

056

059

060 061

062 063

064

065 066

067 068

069

071

073

074

075

076 077



Figure 1: Samples of videos generated by Kling or Gen-3 in *PhyGenBench* with 4 different aspects. The results show that current T2V models struggle to generate videos that align with physical commonsense (e.g., the lack of a plane's reflection in water in the first video of the second row).

078 evaluation frameworks that beyond traditional metrics. While numerous Text-to-Video (T2V) eval-079 uation benchmarks have emerged (Sun et al., 2024; Huang et al., 2024), they primarily focus on 080 various qualities of generated videos (e.g., motion smoothness, background consistency) or spatial 081 relationships, failing to address the critical issue of whether the generated videos adhere to fun-082 damental physical laws. Although some studies have explored the alignment of generated videos 083 with dynamic motions naturalness (Bansal et al., 2024), their benchmarks fail to succinctly capture fundamental physical laws or propose sufficiently robust evaluation methods. Therefore, the devel-084 opment of benchmarks and evaluation methodologies specifically tailored to assess intuitive physics 085 in generated videos remains a critical yet largely unexplored frontier.

- 087 There are two challenges impeding the evaluation of physical commonsense in T2V models. On 088 one hand, there is a lack of benchmarks focused on evaluating physical commonsense. This requires selecting semantically simple physical phenomenons that exhibit clear physical phenomena, allow-089 ing for accurate assessment by either humans or machines. On the other hand, there is a lack of 090 corresponding evaluation metrics. Traditional metrics like FVD (Unterthiner et al., 2018) exhibit 091 limitations in detecting implausible motions (Brooks et al., 2022) and necessitate reference videos, 092 which are often challenging to procure for novel scenes. Recent studies have used video-based VLMs for comprehensive video evaluation (He et al., 2024b; Sun et al., 2024). However, they often 094 struggle to correctly assess physical commonsense. This limitation stems from their inadequate un-095 derstanding of physical laws (Jassim et al., 2023) and the fact that these methods are not specifically 096 designed to evaluate physical laws.
- To address these challenges, we propose PhyGenBench and PhyGenEval to automate the evaluation 098 of physical commonsense understanding capability from T2V models. *PhyGenBench* is designed to evaluate physical commonsense based on fundamental physical laws in text-to-video generation. 100 Inspired by (Halliday et al., 2013), we categorize physical commonsense in the world into four main 101 areas: mechanics, optics, thermal, and material properties. Then principle physical laws and eas-102 ily observable physical phenomenons are identified for each category, resulting in comprehensive 103 27 physical laws and 160 validated prompts in the proposed benchmark. Through brainstorming, 104 we construct prompts that easily reflect physical laws using sources like textbooks (Harjono et al., 105 2020). This process results in a comprehensive but simple set of prompts reflecting physical commonsense, which are sufficiently clear for evaluation. As shown in Figure 1, the correctness of 106 physical commonsense in PhyGenBench can be observed through clear phenomena (e.g., plane 107 should have reflections in water) On the other hand, benefiting from the simple yet clear physical

108 phenomena in *PhyGenBench* prompts, we can propose *PhyGenEval*, which is a novel video eval-109 uation framework for assessing physical commonsense correctness in PhyGenBench. PhyGenEval 110 first uses GPT-40¹ to analyze physical laws in text, addressing the poor understanding of physi-111 cal common sense in video-based VLMs. Moreover, considering that previous evaluation metrics 112 did not specifically target physical correctness, we propose a three-tier hierarchical evaluation strategy for this aspect, transitioning from image-based to comprehensive video analysis: single image, 113 multiple images, and full video stages. Each stage employs distinct VLMs along with custom in-114 structions generated by GPT-40 to form judgments. By combining PhyGenBench and PhyGenEval, 115 we can efficiently evaluate different T2V models' understanding of physical commonsense at scale, 116 producing results highly consistent with human feedback. 117

The contributions of our work are three-fold. i): We proposed PhyGenBench, which compasses a 118 wide range of clear physical phenomenons and explicit physical laws. This benchmark can com-119 prehensively measure whether T2V models understand intuitive physics and indirectly assess their 120 gap from world simulator capabilities ii): Along with the benchmark, we propose an automated 121 evaluation framework - PhyGenEval, which overcomes the challenges of assessing the correctness 122 of physical commonsense with other metrics and demonstrates high consistency with human feed-123 back on PhyGenBench, enabling users to conduct large-scale automated testing of various T2V 124 models. iii): We conduct extensive evaluations of popular T2V models, even the best-performing 125 model, Gen-3, scores only 0.51. This indicates that current models are still far from functioning as 126 world simulators. Based on our evaluation results, we conduct an in-depth analysis and discover 127 that addressing issues such as dynamics is still challenging through prompt engineering or simply 128 scaling up model. We hope this work inspires the community to focus on the learning of physical 129 commonsense in T2V models, rather than merely using them as tools for entertainment.

130 131

132 133

134

2 RELATED WORK

2.1 BENCHMARKS FOR TEXT-TO-VIDEO GENERATION

135 The rapid advancement of text-to-video (T2V) generation models has necessitated various bench-136 marks for accurate assessment. Traditional works in video generation, such as FVD (Unterthiner 137 et al., 2018), rely on datasets like UCF-101 (Soomro, 2012) and Kinetics-400 (Kay et al., 2017), 138 which are limited in scope. Recent benchmarks, including VBench (Huang et al., 2024) and Eval-139 Crafter (Liu et al., 2024c), aim to comprehensively evaluate general video quality across multi-140 ple dimensions. In contrast, some studies focus on fine-grained evaluation of text-to-video (T2V) 141 models from specific aspects. For instance, T2V-CompBench (Sun et al., 2024) assesses compo-142 sitional generation capabilities, while DEVIL (Liao et al., 2024) evaluates dynamic characteristics of generated videos. Although some research like VideoPhy (Bansal et al., 2024) efforts address 143 the dynamic motions naturalness of video generation, their benchmarks fail to succinctly capture 144 fundamental physical laws. Consequently, most existing works overlook this crucial aspect, which 145 forms the foundation for realizing a world simulator. To address this gap, we introduce PhyGen-146 Bench, a benchmark designed to comprehensively measure T2V models' understanding of physical 147 commonsense.

148 149 150

2.2 EVALUATION METRICS FOR TEXT-TO-VIDEO GENERATION

151 Conventional approaches to video quality assessment often employ metrics such as FVD (Un-152 terthiner et al., 2018) and IS (Salimans et al., 2016). However, the detection of unrealistic motions 153 is difficult for them (Brooks et al., 2022), and FVD requires a reference video that is hard to obtain 154 for novel scenes, making it challenging to evaluate the correctness of physical commonsense. Re-155 cent studies have explored the use of advanced vision-language models (VLMs) as evaluators. For 156 instance, VideoScore (He et al., 2024b) leverages human feedback to train models for video quality 157 assessment, while T2V-CompBench (Sun et al., 2024) utilizes powerful models like LLaVA (Liu et al., 2024a) to evaluate the correctness of spatial relationships. Although a few works demon-158 strate improved alignment with human judgments, they fall short in generalizing to assessments of 159 physical commonsense correctness. To address this limitation, we introduce PhyGenEval, a novel 160

¹⁶¹

¹The version is gpt4o-0806



(b) The construction pipeline of the PhyGenBench

Figure 2: (a) is the overview of the proposed *PhyGenBench*. (b) is the *PhyGenBench* data pipeline, which covers four physics categories. We select key physical laws and manually craft initial prompts that reflect the corresponding physical phenomena. GPT-40 adds details and enhances diversity by varying objects. After manual review, we obtain 160 T2V prompts.

method designed to evaluate physical commonsense correctness on *PhyGenBench*. We validate the efficacy of our approach through comprehensive human correlation studies.

3 PHYGENBENCH

Inspired by (Swartz, 1985), we first define the following terms: "*Physical Commonsense:*" Basic intuitive understanding of how physical objects and actions behave in everyday life; "*Physical Laws:*" Universal scientific principles that describe consistent behaviors in nature; "*Physical Phenomenon:*" Observable events or processes caused by the interaction of physical laws. The purpose of *PhyGenBench* is to evaluate whether T2V models understand physical commonsense, while each prompt in *PhyGenBench* presents a clear physical phenomenon and an underlying physical law.

Overview. As illustrated in Figure 2 (a), *PhyGenBench* encompasses four major categories of physical commonsense: "*Mechanics*", "*Optics*", "*Thermal*", and "*Material Properties*". It incorporates 27 physical phenomena with intrinsic physical laws reflected by the corresponding designed 160 prompts:

204
1. "*Mechanics*" covers 7 common mechanical laws: gravity, buoyancy, solid pressure, atmospheric
205 pressure, elasticity, friction, and surface tension, with 40 validated prompts. For example, we use
"A piece of iron is gently placed on the surface of the water in a tank filled with water" to test T2V
208 model's understanding of Buoyancy, where the iron should sink due to its higher density compared
209 to water.

209
210 2. "Optics" categorizes 6 aspects based on light phenomena: reflection, refraction, scattering, dispersion, interference & diffraction, and straight-line propagation, yielding 50 prompts. A prompt like "a kite soaring above a smooth and tranquil pond" is used to test reflection generation capability.

3. "*Thermal*" considers 6 phase transitions: Solidification, Melting, Liquefaction, Boiling, deposition, Sublimation, comprising 30 prompts. Inspired by ChronoMagicBench (Yuan et al., 2024), the vaporization (boiling) process is evaluated by the prompt "a timelapse capturing the transformation of water as the temperature rapidly rises above 100°C".

4. "*Material Properties*" includes 5 physical properties (color, hardness, solubility, combustibility, and flame reaction) and 3 chemical properties (acidity, redox potential, and dehydrating properties), resulting in 40 prompts. We reflect material properties, e.g., "*hardness*", through the prompts with expected phenomena, e.g., "*an egg being hurled with significant force towards a rock*", where the egg should break while the rock remains intact.

Multiple physical laws could be included in a single prompt, which may bring confusion to the evaluation of physical common sense in video generation, even for human annotators. To avoid this, we carefully curate prompts to ensure a one-to-one correspondence for each physical phenomenon it reflects, with clear physical law inside. By incorporating physical laws from four distinct physical categories, *PhyGenBench* offers a thorough assessment of current T2V models' understanding of physical commonsense.

227

228 **Benchmark Construction.** As shown in Figure 2 (b), we develop a comprehensive approach to 229 create *PhyGenBench*. The methodology encompasses five steps: 1) Conceptualization: Following 230 (Halliday et al., 2013), We first identify key physical commonsense from four major categories of 231 physics. For each category, we select specific physical laws from textbooks (Harjono et al., 2020), 232 which can be widely recognized and can be easily demonstrated through clear, observable phys-233 ical phenomenon. 2) Prompt Engineering: For each physical law, we manually craft the initial T2V prompts to clearly depict the underlying physical phenomenon 3) Prompt Augmentation: To 234 enhance the model's video generation capabilities, we augment the initial T2V prompts by adding 235 additional details, such as more precise descriptions of objects and actions (Yang et al., 2024). This 236 augmentation process is carefully designed to avoid revealing the expected physical phenomenon. 237 4) Diversity Enhancement: Following T2V-CompBench (Sun et al., 2024), we employ GPT-40 238 to perform object substitution on the augmented prompts. This step increases the diversity of the 239 benchmark. 5) Quality Control: We conduct a thorough review of the prompts and their associated 240 physical laws to ensure accuracy and relevance. Specifically, we ensure that the T2V prompts and 241 corresponding physical laws are clear and accurate. We then randomly use the current T2V model to 242 check if the prompts are simple enough for the model to generate semantically accurate videos. This 243 methodology yields a robust and comprehensive benchmark for assessing T2V models' comprehension of physical commonsense, providing a valuable tool for advancing research in this domain. For 244 245 more detailed information about the dataset, please refer to the Appendix B

246 247

248 249

4 PhyGenEval

PhyGenEval aims to assess whether the physical phenomena in the generated videos conform to 250 the corresponding physical laws. To obtain a clear judgment, we decompose the evaluation into 251 semantic alignment (SA) and physical commonsense alignment (PCA). While SA evaluates whether 252 the semantic meaning inferred by the generated video and the input prompt are matched, PCA 253 measures whether the evaluated physical laws are grounded in the videos. For example, for the 254 scene "an egg collides with a stone", SA requires a video containing the egg, the stone, and the 255 collision action. PCA necessitates a video for the whole physical motions in which the egg hits a 256 stone and then breaks, while the stone remains intact. Following (He et al., 2024b), we convert both 257 SA and PCA to a four-point scale, as well as the human ratings.

258 259

260

4.1 SEMANTIC ALIGNMENT EVALUATION

261 Directly asking the Vision-Language Model(VLM) to align the semantic meaning between videos 262 and input prompts are difficult, as prompts usually are mixed with semantic entities and physical phe-263 nomena, and the intermediate outcomes are subtly implied by the videos. For example, in a prompt 264 like "A timelapse captures the transformation of soup as the temperature rises above 100° C", a pos-265 sible video generation would appear like "The video shows a soup, but there is no transformation of 266 the soup". To address the challenge, we first employ GPT-40 to extract object and action from the 267 original text prompt, we then utilize GPT-40 to sequentially determine the presence of extracted objects in the video and verify the occurrence of specified actions. This decomposition provides more 268 fine-grained captures and prevents the model from confusing semantic and physical correctness dur-269 ing evaluation. Experimental results demonstrate that our automated evaluation method aligns more



Figure 3: An overview of the proposed *PhyGenEval*. *PhyGenEval* is divided into three parts: Key Physical Phenomena Detection, Physics Order Verification, and Overall Naturalness Evaluation. Each part uses an appropriate VLM in combination with physical-based customized questions generated by GPT-40. The final score is the combined result of the three parts. For the example in the figure, the three-stage scores are 0, 1 (only $Question_1$ is correct), and 0. The final score is calculated as 0 according to Overall Score calculation in Section 4.2. 292

closely with human judgment and outperforms previous methods (He et al., 2024b; Sun et al., 2024) in *PhyGenBench* (Appendix C.1).

4.2 PHYSICAL COMMONSENSE EVALUATION

299 To evaluate physical correctness in the video, we evaluated multiple common evaluation metrics comparing human assessments². Experimental results in Table 1 demonstrate that these methods 300 struggle to generalize to the assessment of physical commonsense correctness on PhyGenBench, 301 e.g., VideoScore (He et al., 2024b) has only a spearman correlation of 0.19 on PhyGenBench, which 302 is most correlated with human assessments except PhyGenEval. We attribute it to the main factor: 303 Directly using video-based VLMs fails to comprehend the embedded physical commonsense (Jassim 304 et al., 2023), as current methods are not designed with physical commonsense as a foundation. To 305 fully understand the physical commonsense in the video, there are three key factors need to solve: 306 i): Physical processes typically exhibit clear key phenomena depicted by the input prompt (e.g., 307 "the egg breaks upon hitting the rock."). It is necessary to identify these key physical phenomena 308 and detect their presence in videos. ii): Physical processes are characterized by causality, manifested in the correct sequence of critical events(e.g., "The egg touchs the rock first, then breaks."). 309 The correct sequence order validates the correctness of physical processes. iii): Physical processes 310 need to possess overall naturalness, which represents the realistic of the overall process. To address 311 these factors, we design a progressive strategy that starts with key physical phenomena, then moves 312 through the sequence of several key phenomena, and finally evaluates the overall naturalness of 313 the entire video. This hierarchical and refined approach reduces the difficulty compared to exist-314 ing methods that directly uses VLMs to evaluate physical commonsense, enabling PhyGenEval to 315 achieve results closely aligned with human judgements.

316

289

290

291

293

295 296 297

298

317 **Key Physical Phenomena Detection.** This stage aims to detect whether the key physical phe-318 nomena occur in the video. Here we define the key phenomena as an observable and distinctive 319 occurrence (e.g., a specific frame) within a physical process that can directly reveal the correspond-320 ing physical law, like deformations or color changes. For each input prompt in *PhyGenBench*, we 321 craft a retrieval prompt p_r and a set of physics-related questions Q, where the retrieval prompt is

³²² 323

²Annotators are asked to score the correctness of physical commonsense in the video. Details refer to Section 5 and Appendix D.1

used to locate the key phenomena frame, and physical-related questions are utilized to check whether
 the expected physics phenomena are present in the keyframe.

As illustrated in Figure 3 (a), we first obtained both Q and P_r by prompting GPT-40 with the input T2V prompt and corresponding physical law. Following (Hessel et al., 2021), a keyframe I_i from the video based on the retrieval prompt, where I_i is the *i*-th frame in the video. By using the keyframe, we define a confidence score of the key phenomena in the video:.

$$S_{key} = \sum_{q \in Q} \max_{i-2 \le j \le i+2} \left(VLM(I_j, q) + VLM(I_j, p_r) \right),$$

where $VLM(I_j, q)$ reflects the presence of physical phenomena in I_j for each related question qfrom Q. $VLM(I_j, p_r)$ checks whether I_j matches the retrieval prompt, which ensures key phenomena occur at the correct frame. Since videos may contain semantic errors, it's also important for determining if key physical phenomena occur (e.g., an egg shouldn't break in mid-air before hitting a rock). We consider adjacent 5 frames near the keyframe to enhance the robustness. For example, the egg may not be cracked just when it first contacts the stone. We instantiate VLM-based evaluator VLM(·) with VQAScore (Lin et al., 2024), which has been shown promising evaluation results on visual question-answering.

342 **Physics Order Verification.** In this stage, we verify whether key physical phenomena occur in the 343 correct order. The correct physical sequence is an ordered series of events in a physical process that reflects causality, which represents the necessary prerequisites and temporal order of key physical 344 phenomena. As an example, the egg should first touch the stone and then crack. Considering current 345 models in *PhyGenBench* generally maintain outcome consistency (Huang et al., 2024) (e.g., the egg 346 would not reassemble itself after it is broken). we approach this direction by investigating the order 347 correctness from the keyframes (Figure 3 (b)), e.g., the keyframe of the egg hits the stone should be 348 ahead of the keyframe of the broken egg. 349

Similar to the single image evaluation, we prompt GPT-40 to generate a retrieval prompt p_r and three physical-related questions (q_1, q_2, q_3) . p_r is used to locate the keyframe (e.g., the moment the egg slightly touches the stone.). While q_1 , q_2 , and q_3 are questions to check the order correctness from the first frame to the keyframe, from the keyframe to the last frame, and from the first frame to the last frame, respectively. Similarly, we first use CLIPScore to locate the key frame I_i , then the order correctness scores of S_{before} and S_{after} are defined as:

331 332 333

341

358 359

$$S_{\text{before}} = \max_{i-2 \le j \le i} (\text{VLM}(I_0, I_j, q_1) + \text{VLM}(I_j, p_r))$$

$$S_{\text{after}} = \max_{i \le j \le i+2} (\text{VLM}(I_j, I_{-1}, q_2) + \text{VLM}(I_j, p_r))$$

360 q_3 assesses the overall physical sequence coherence of the video. The score of answering q_3 is 361 defined as by $S_{all} = VLM(I_0, I_{i-2:i+2}, I_{-1}, q_3)$, which evaluates the overall sequence (similar to 362 the input video but using manually selected key frames). Here we employ GPT-40 or LLaVA-363 Interleave (Li et al., 2024) as the VLM-based evaluator VLM(·), as they demonstrate exceptional 364 multi-image comprehension capabilities. The overall score of whole physical order evaluation can 365 be formulated as $S_{order} = S_{before} + S_{after} + S_{all}$

365 366

375 376

Overall Naturalness Evaluation. This stage aims to evaluate the overall naturalness of the 367 video. we define naturalness as the dynamic progression that aligns with real-world physical phe-368 nomenons (Liao et al., 2024). For each prompt in PhyGenBench, we obtain a naturalness evaluation 369 standard, denoted as g_{spec} , which is used to assess the naturalness for video. As shown in Figure 370 3 (c), we first refer to DEVIL (Liao et al., 2024) to establish a general evaluation standard: g_{aen} , 371 applicable to all T2V prompts. Besides, we use each input T2V prompt p, the corresponding physical law l, and general evaluation standard g_{gen} to guide GPT-40 in generating a detailed evaluation 372 standard: g_{spec} , for the given prompt. Finally, we require the VLM to score based on p, l, g_{spec} , and 373 the corresponding video denoted by $I_{0:-1}$. Formally, we define the overall naturalness score as: 374

$$S_{natural} = VLM(I_{0:-1}, p, l, g_{spec})$$

We implement the VLM-based evaluator $VLM(\cdot)$ using InternVideo2 (Wang et al., 2024) and GPT-40, both of which have demonstrated promising results in video understanding.

Metric	Mec	hanics	Op	tics	The	rmal	Mat	erial	Overall		
Metric	$\tau(\uparrow)$	$\rho(\uparrow)$									
DEVIL (Liao et al., 2024)	0.15	0.16	0.03	0.03	0.10	0.11	0.27	0.29	0.17	0.18	
VideoPhy (Bansal et al., 2024)	0.00	-0.03	-0.15	-0.14	0.08	0.08	0.13	0.14	0.03	0.04	
VideoScore (He et al., 2024b)	0.18	0.20	0.07	0.08	0.14	0.15	0.14	0.15	0.17	0.19	
PhyGenEval	0.72	0.75	0.76	0.77	0.73	0.75	0.81	0.84	0.78	0.81	

378 Table 1: PCA correlation results with proposed PhyGenEval in video generation. PhyGenEval 379 is significantly closer to human feedback on *PhyGenBench* compared to other metrics.

We first discretize $S_{\rm key}, S_{\rm order},$ and $S_{\rm natural}$ from the three stages into a four-point **Overall Score.** scale, then take their average and apply floor rounding as the final score. For robust purposes, we evaluate S_{order} with both GPT40 and LLaVA-Interleave and $S_{natural}$ with both GPT40 and Intern-Video2. The final score is calculated as the ensemble of two methods. Detailed calculation protocols are provided in Appendix C.

5 EXPERIMENT

Experiments Setup. We evaluate 5 open-source models including OpenSora V1.2 (Zheng et al., 2024a), Lavie (Wang et al., 2023), CogVideoX 2b (Yang et al., 2024), CogVideoX 5b (Yang et al., 398 2024), and Vchitect2.0 (Wang et al., 2023), as well as proprietary models Kling (kli, 2024), Pika 399 (Pik, 2023), and Gen-3 (gen, 2024). We compare our proposed metric with existing metrics or 400 benchmarks: Videophy (Bansal et al., 2024), VideoScore (He et al., 2024b) and DEVIL (Liao et al., 401 2024) More Detailed information is provided in Appendix D. 402

For human evaluation, we compared the results across 8 T2V models. We randomly select 64 403 prompts from PhyGenBench and generate 64 videos for each T2V model. Therefore we need eval-404 uation 512 videos. We ask three annotators to provide semantic and physical scores for each video³. 405 Each annotator will give an integer score of 0-3 for the semantic and physical scores, and the final 406 score is the average of the three scores and rounded up. Finally, we calculate the correlation between 407 the human scores and automatic evaluation scores using Kendall's τ and Spearman's ρ , we put more 408 detailed information about human evaluation in Appendix D.1. 409

410 Human Evaluation. As shown in Table 1, current video generation evaluation metrics largely 411 overlook physical correctness. In contrast, *PhyGenEval* implements a detailed design for evaluating 412 physical correctness, demonstrating strong correlations with human judgments across all categories. Its overall correlation coefficient reaches 0.81, indicating that *PhyGenEval* serves as an effective 413 human-aligned physical commonsense correctness evaluator for PhyGenBench. We put more results 414 in Appendix D.2 415

416 We conduct several case studies to illustrate the differences between various metrics more clearly. 417 As shown in Figure 4, (a) and (f) reveal that VideoScore and DEVIL are prone to misclassifying 418 videos that have smooth and consistent motion but violate fundamental physical laws. Specifically, as for (a), when "an egg exhibits rubber-like elasticity upon impact with a rock instead of breaking," 419 these metrics incorrectly evaluate it as physically correct. VideoPhy exhibits similar limitations. 420 In (c), it incorrectly assesses "a rock floating on water instead of sinking" as physically correct. 421 Furthermore, our analysis reveals a major flaw in these three methodologies: they cannot incorporate 422 domain-specific physical commonsense. As illustrated in (e), where "the flame from burning copper 423 appears red instead of green," these metrics fail to identify the mistake. This indicates their inability 424 to incorporate domain-specific physical commonsense. In contrast, PhyGenEval demonstrates a 425 robust integration of physical commonsense and comprehensive video content analysis, resulting in 426 more accurate and physically consistent evaluations in *PhyGenBench*. 427

428 Quantitative Evaluation. We conduct extensive experiments on a wide range of popular video 429 generation models. As illustrated in Table 2, even the best-performing model, Gen-3, only attains 430 a PCA score of 0.51 on PhyGenBench. This indicates that even for prompts containing obvious 431

387 388

389

390

391

392

393 394

³Note that we ask the annotators to focus on the correctness of the physical phenomena for physical scores.



Figure 4: Different video generation evaluation metric in *PhyGenBench*. Except for the proposed *PhyGenEval*, the current methods cannot reasonably assess the correctness of physical commonsense in videos from *PhyGenBench*.

Table 2: Evaluation results of PCA with the proposed *PhyGenEval* in videos generated by several models. The results reveal that all models score very low in physical commonsense accuracy. The scores are normalized to a range of 0-1.

Model	Size	Mechanics (↑)	$Optics(\uparrow)$	$Thermal(\uparrow)$	$Material(\uparrow)$	Average (↑)	$Human(\uparrow)$
CogVideoX (Yang et al., 2024)	2B	0.38	0.43	0.34	0.39	0.39	0.31
CogVideoX (Yang et al., 2024)	5B	0.39	0.55	0.40	0.42	0.45	0.37
Open-Sora V1.2 (Zheng et al., 2024a)	1.1B	0.43	0.50	0.44	0.37	0.44	0.35
Lavie (Wang et al., 2023)	860M	0.30	0.44	0.38	0.32	0.36	0.30
Vchitect 2.0 (Wang et al., 2023)	2B	0.41	0.56	0.44	0.37	0.45	0.36
Pika (Pik, 2023)	-	0.35	0.56	0.43	0.39	0.44	0.36
Gen-3 (gen, 2024)	-	0.45	0.57	0.49	0.51	0.51	0.48
Kling (kli, 2024)	-	0.45	0.58	0.50	0.40	0.49	0.44

physical commonsense, current T2V models struggle to generate videos that comply with intuitive
 physics. It indirectly reflects that these models are still far from achieving the world simulator.

Furthermore, we identify the following key observations: 1): Across various categories of physical commonsense, all models consistently demonstrate superior performance in the domain of optics compared to other areas. Notably, Vchitect2.0 and CogVideoX-5b achieve a PCA score in the optics domain comparable to that of closed-source models. We posit that this superior performance in the optics domain can be attributed to the abundant and explicit representation of optical knowledge in pre-training datasets, thereby enhancing the model's comprehension in this area. 2): Kling and Gen-3 exhibit significantly higher performance compared to other models. Specifically, Gen-3 demon-strates a robust understanding of material properties, achieving a score of 0.51, which substantially surpasses other models. Kling performs particularly well in thermal, attaining the highest score of 0.50 in this domain. 3): Among open-source models, Vchitect2.0 and CogVideoX 5b perform comparatively well, both exceeding the performance level of Pika. In contrast, Lavie consistently exhibits lower physical correctness across all categories.

Qualitative Evaluation. The different video cases for 4 physical commonsense categories are illustrated in Figure 5. Our main observations are as follows: In mechanics, the models struggle to generate simple physically accurate phenomenons. As shown in Figure 5, all models fail to depict the glass ball sinking in water. As for (b), instead showing it floating on the surface, OpenSora and Gen-3 even produce videos where the ball is suspended. Additionally, the models do not capture special physical phenomenonss, such as the state of water in zero gravity, as seen in (a). In optics, the models perform relatively better. (c) and (d) show the models handling reflections of balloons in wa-ter and colorful bubbles, though OpenSora and CogVideoX still produce reflections with noticeable distortions in (d). In thermal, the models fail to generate accurate videos of phase transitions. For the melting phenomenon in (e), most models show incorrect results, with CogVideoX even producing



Figure 5: Qualitative comparisons of four categories. Current models perform relatively well in generating optical phenomenons but are weaker in mechanics, thermal, and material properties.

a video where the ice cream increases in size. Similar errors appear in the sublimation process in (f), with only Gen-3 showing partial understanding. Regarding material properties, (g) shows all models failing to recognize that an egg should break when hitting a rock, with Kling displaying the egg bouncing like a rubber ball. For simple chemical reactions, such as the black bread experiment in (h), none of the models demonstrate an accurate understanding of the expected reaction.

Ablation Study. We conduct a detailed robustness analysis of the design elements in PhyGenEval, including the role of each level in the three-tier evaluation framework and the impact
of the two-stage strategy proposed in overall naturalness evaluation. Experimental results show that
the key designs of *PhyGenEval* are essential. Detailed results are provided in Appendix D.3.

6 DISCUSSION

To explore potential solutions for the challenges posed by PhyGenBench, We focus on widely used 521 and proven-effective methods such as scaling laws (Kaplan et al., 2020), prompt engineering (Fu 522 et al., 2024), and some method like Venhancer (He et al., 2024a) aimed to improve general video 523 quality (Huang et al., 2024). Through quantitative and qualitative analysis, we find: 1) Scaling up 524 models can solve some issues but still fails to handle dynamic physical phenomenons, which we 525 believe requires extensive training on synthetic data. 2) Prompt engineering like (Fu et al., 2024) 526 only solves a few simple issues (e.g., flame color), highlighting the difficulty and significance of 527 PhyGenBench. 3) While some methods improve general video quality, they do not enhance the 528 model's understanding of physical commonsense. More detailed results are in Appendix E.

529 530

518

519

530 7 CONCLUSION

In this paper, we explore the gap between current T2V models' understanding of physical commonsense and their role as world simulators. To achieve this, we introduce *PhyGenBench* and *PhyGenEval*. *PhyGenBench* is a benchmark specifically designed to assess models' understanding of physical commonsense, featuring various physical laws and simple, clear physical phenomenons. Alongside *PhyGenBench*, we propose a novel three-tier hierarchical evaluation framework called *PhyGenEval* to automate the evaluation process. Experimental and analytical results show that current T2V models struggle to generate videos that align with physical commonsense, highlighting a significant gap from world simulation. Moreover, simply scaling up models or applying prompt engineering fails to address issues in *PhyGenBench*, such as those involving dynamics.

540	REFERENCES
541	

- 542 Pika, 2023. URL https://www.pika.art/.
- 543 544 Gen-3, 2024. URL https://runwayml.com/blog/introducing-gen-3-alpha/.
- 545 546 Kling, 2024. URL https://kling.kuaishou.com/.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
 report. *arXiv preprint arXiv:2303.08774*, 2023.
- Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chenfanfu Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. Videophy: Evaluating physical common-sense for video generation. *arXiv preprint arXiv:2406.03520*, 2024.
- Peter W Battaglia, Jessica B Hamrick, and Joshua B Tenenbaum. Simulation as an engine of physical
 scene understanding. *Proceedings of the National Academy of Sciences*, 110(45):18327–18332,
 2013.
- Daniel M Bear, Elias Wang, Damian Mrowca, Felix J Binder, Hsiao-Yu Fish Tung, RT Pramod,
 Cameron Holdaway, Sirui Tao, Kevin Smith, Fan-Yun Sun, et al. Physion: Evaluating physical
 prediction from vision in humans and machines. *arXiv preprint arXiv:2106.08261*, 2021.
- Tim Brooks, Janne Hellsten, Miika Aittala, Ting-Chun Wang, Timo Aila, Jaakko Lehtinen, Ming Yu Liu, Alexei Efros, and Tero Karras. Generating long videos of dynamic scenes. Advances in
 Neural Information Processing Systems, 35:31769–31781, 2022.
- Xingyu Fu, Muyu He, Yujie Lu, William Yang Wang, and Dan Roth. Commonsense-t2i challenge: Can text-to-image generation models understand commonsense? *arXiv preprint* arXiv:2406.07546, 2024.
- Shenyuan Gao, Jiazhi Yang, Li Chen, Kashyap Chitta, Yihang Qiu, Andreas Geiger, Jun Zhang, and Hongyang Li. Vista: A generalizable driving world model with high fidelity and versatile controllability. *arXiv preprint arXiv:2405.17398*, 2024.
- David Halliday, Robert Resnick, and Jearl Walker. *Fundamentals of physics*. John Wiley & Sons, 2013.
- Ahmad Harjono, Gunawan Gunawan, Rabiatul Adawiyah, and Lovy Herayanti. An interactive ebook for physics to improve students' conceptual mastery. *International Journal of Emerging Technologies in Learning (iJET)*, 15(5):40–49, 2020.
- Jingwen He, Tianfan Xue, Dongyang Liu, Xinqi Lin, Peng Gao, Dahua Lin, Yu Qiao, Wanli Ouyang, and Ziwei Liu. Venhancer: Generative space-time enhancement for video generation. *arXiv* preprint arXiv:2407.07667, 2024a.
- 581 Xuan He, Dongfu Jiang, Ge Zhang, Max Ku, Achint Soni, Sherman Siu, Haonan Chen, Abhranil
 582 Chandra, Ziyan Jiang, Aaran Arulraj, et al. Mantisscore: Building automatic metrics to simulate
 583 fine-grained human feedback for video generation. *arXiv preprint arXiv:2406.15252*, 2024b.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- Ziqi Huang, Yinan 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 21807–21818, 2024.
- Serwan Jassim, Mario Holubar, Annika Richter, Cornelius Wolff, Xenia Ohmer, and Elia Bruni.
 Grasp: A novel benchmark for evaluating language grounding and situated physics understanding in multimodal language models. *arXiv preprint arXiv:2311.09048*, 2023.

635

- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
 models. *arXiv preprint arXiv:2001.08361*, 2020.
- Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijaya narasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action
 video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- 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.
- Zhuowan Li, Xingrui Wang, Elias Stengel-Eskin, Adam Kortylewski, Wufei Ma, Benjamin Van Durme, and Alan L. Yuille. Super-clevr: A virtual benchmark to diagnose domain robustness in visual reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14963–14973, June 2023.
- Mingxiang Liao, Hannan Lu, Xinyu Zhang, Fang Wan, Tianyu Wang, Yuzhong Zhao, Wangmeng
 Zuo, Qixiang Ye, and Jingdong Wang. Evaluation of text-to-video generation models: A dynamics
 perspective. *arXiv preprint arXiv:2407.01094*, 2024.
- ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁶
 ⁶¹⁶
 ⁶¹⁷
 ⁶¹⁷
 ⁶¹⁸
 ⁶¹⁸
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹⁹
 ⁶¹¹
 ⁶¹¹
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹²
 ⁶¹³
 ⁶¹⁴
 ⁶¹⁵
 ⁶¹⁵
 ⁶¹⁴
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances
 in neural information processing systems, 36, 2024a.
- Shaowei Liu, Zhongzheng Ren, Saurabh Gupta, and Shenlong Wang. Physgen: Rigid-body physicsgrounded image-to-video generation. In *European Conference on Computer Vision ECCV*, 2024b.
- Yaofang Liu, Xiaodong Cun, Xuebo Liu, Xintao Wang, Yong Zhang, Haoxin Chen, Yang Liu, Tieyong Zeng, Raymond Chan, and Ying Shan. Evalcrafter: Benchmarking and evaluating large video generation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22139–22149, 2024c.
- Pietro Mazzaglia, Tim Verbelen, Bart Dhoedt, Aaron Courville, and Sai Rajeswar. Multimodal foundation world models for generalist embodied agents. *arXiv preprint arXiv:2406.18043*, 2024.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gem ini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016.
- K Soomro. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- Kaiyue Sun, Kaiyi Huang, Xian Liu, Yue Wu, Zihan Xu, Zhenguo Li, and Xihui Liu. T2vcompbench: A comprehensive benchmark for compositional text-to-video generation. *arXiv preprint arXiv:2407.14505*, 2024.
- 642 Norman Swartz. *The concept of physical law*. Cambridge University Press, 1985.
- Zhiyu Tan, Xiaomeng Yang, Luozheng Qin, and Hao Li. Vidgen-1m: A large-scale dataset for text-to-video generation. *arXiv preprint arXiv:2408.02629*, 2024.
- Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018.

648	Yaohui Wang Xinyuan Chen Xin Ma Shangchen Zhou Zigi Huang Yi Wang Ceyuan Yang Yinan
649	He. Jiashuo Yu, Peiging Yang, et al. Lavie: High-quality video generation with cascaded latent
650	diffusion models. arXiv preprint arXiv:2309.15103, 2023.

- Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Jilan Xu, Zun Wang, et al. Internvideo2: Scaling video foundation models for multimodal video understanding. arXiv preprint arXiv:2403.15377, 2024.
- Justin N Wood, Tomer D Ullman, Brian W Wood, Elizabeth S Spelke, and Samantha MW Wood. Object permanence in newborn chicks is robust against opposing evidence. arXiv preprint arXiv:2402.14641, 2024.
- Jiannan Xiang, Guangyi Liu, Yi Gu, Qiyue Gao, Yuting Ning, Yuheng Zha, Zeyu Feng, Tianhua Tao, Shibo Hao, Yemin Shi, et al. Pandora: Towards general world model with natural language actions and video states. arXiv preprint arXiv:2406.09455, 2024.
- 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.
- Shenghai Yuan, Jinfa Huang, Yongqi Xu, Yaoyang Liu, Shaofeng Zhang, Yujun Shi, Ruijie Zhu, Xinhua Cheng, Jiebo Luo, and Li Yuan. Chronomagic-bench: A benchmark for metamorphic evaluation of text-to-time-lapse video generation. arXiv preprint arXiv:2406.18522, 2024.
- Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all, March 2024a. URL https://github.com/hpcaitech/Open-Sora.
- Zhicheng Zheng, Xin Yan, Zhenfang Chen, Jingzhou Wang, Qin Zhi Eddie Lim, Joshua B Tenen-baum, and Chuang Gan. Contphy: Continuum physical concept learning and reasoning from videos. arXiv preprint arXiv:2402.06119, 2024b.
- Zheng Zhu, Xiaofeng Wang, Wangbo Zhao, Chen Min, Nianchen Deng, Min Dou, Yuqi Wang, Botian Shi, Kai Wang, Chi Zhang, et al. Is sora a world simulator? a comprehensive survey on general world models and beyond. arXiv preprint arXiv:2405.03520, 2024.

702 A RELATED WORK

A.1 BENCHMARKS FOR PHYSICAL UNDERSTANDING

706 Recent benchmarks such as SuperCLEVR-Physics Li et al. (2023), ContPhy Zheng et al. (2024b), 707 and Physion Bear et al. (2021) have significantly advanced the evaluation of AI's physical reason-708 ing, focusing on understanding and predictive tasks. SuperCLEVR-Physics emphasizes reasoning 709 about dynamic properties like velocity and collisions in 4D scenes, ContPhy expands the evaluation 710 scope to include diverse physical properties, such as mass and density, within continuum settings. It underscores the limitations of existing AI models in handling soft-body dynamics. And Physion 711 evaluates models' ability to predict physical phenomena like collisions and motion while bench-712 marking against human behavior. However, these works primarily target understanding or prediction 713 rather than generative capabilities. In contrast, our work introduces PhyGenBench, a comprehen-714 sive benchmark designed to evaluate whether Text-to-Video (T2V) models can generate videos that 715 adhere to physical commonsense. Unlike existing benchmarks, our work highlights the generative 716 challenges in intuitive physics, revealing critical gaps in current T2V models and underscoring the 717 need to advance physical commonsense for applications beyond entertainment.

718 719 720

721 722

704

705

B PHYGENBENCH DETAILS

B.1 DETAILED OVERVIEW

723 A fine-grained analysis of the dataset is essential for 724 a comprehensive understanding of the benchmark. As 725 shown in Table 3, PhyGenBench covers 4 major domains 726 in physics, encompassing 27 representative physical laws, 727 which enables it to provide a more comprehensive and 728 fine-grained evaluation of models' physical capabilities. 729 We generated 1280 videos by evaluating 8 advanced mod-730 els. Additionally, our captions encompass totally 165 731 unique objects and 42 unique actions with an average length of 18.75 words. 732

733 734

735

B.2 DIFFERENCE BETWEEN VIDEOPHY AND OURS

736 VIDEOPHY Bansal et al. (2024) comprises 688 curated simple prompts that focus on interactions between three 737 types of physical materials: solid-solid, solid-fluid, and 738 fluid-fluid, but lack annotations of physical laws. The 739 dataset is designed to evaluate a model's understanding 740 of physical commonsense, featuring a limited range of 741 physical phenomenons such as rigid body interactions, 742 fluid dynamics, and contact forces. We are better suited 743 than Videophy for evaluating physical commonsense due 744 to two significant differences.

Table 3: Details of PhyGenBench

Statistic	Number
Physical Laws	27
Domains	4
Optics	50
Mechanics	40
Thermal	30
Material Properties	40
Total Captions	160
Total T2V Models	8
Total Generated Videos	1280
Unique Objects	165
Unique Actions	42
Average Length of Caption	18.75

745 First As shown in Figure 2, PhyGenBench includes 160 carefully crafted prompts across 27 distinct 746 physical laws, spanning four fundamental domains, which comprehensively assess a model's un-747 derstanding of physical commonsense. While Videophy primarily focuses on interactions between 748 solid-fluid, solid-solid, and fluid-fluid, limiting its coverage and overlooking common physical laws 749 such as phase transitions and basic material properties. What' more, Videophy lacks annotations of 750 physical laws making it hard for VLM model to evaluate. Second, as shown in Table 4, the average 751 SA score of *PhyGenBench* (0.80) significantly outperforms that of Videophy (0.63). This indicates 752 that PhyGenBench prompts are well-suited and easy for T2V models to generate high-quality, well-753 aligned videos, which benefits evaluation of physical correctness. In contrast, as shown in Figure 6, We find that prompts from Videophy pose challenges for T2V models in generating text-aligned 754 and high-quality videos for two main reasons: 1. The prompts lack detail and specificity. For in-755 stance, "A tissue blots a tear from an eye" is overly simplistic (without augmentation). Modern T2V

models, such as CogVideo5B Yang et al. (2024), are typically trained with longer and more descriptive captions, which enhance their ability to comprehend and generate content based on prompts.
The scenes are often complex and unrealistic. For example, "The wristwatch knob winds the inner spring tightly" describes a process involving intricate internal mechanisms that are not visible externally. As a result, it is exceedingly difficult for T2V models to generate such scenes accurately.



Figure 6: **Samples of videos generated by Kling, Vchitect, and Cogvideo5b in Videophy.** All T2V models struggle to achieve proper text alignment and produce high-quality videos, making it meaningless to evaluate physical correctness in Videophy.

Table 4: **Comparison of SA results for video generation between Videophy and PhyGenBench.** We randomly select 64 prompts from both Videophy and *PhyGenBench*, use different T2V models to generate videos, and then ask annotators to score based on our cretiera in Figure 10. The results show that *PhyGenBench* 's SA scores significantly outperform Videophy.

Model	Size	Videophy (↑)	PhyGenBench (\uparrow)
CogVideoX (Yang et al., 2024)	5B	0.48	0.78
Vchitect 2.0	2B	0.63	0.84
Kling	-	0.77	0.89
Average	-	0.63	0.80

C PHYGENEVAL DETAILS

C.1 SEMANTIC ALIGNMENT DETAILS

To reduce the complexity for VLM models to evaluate sementic correctness of generated videos between prompts, we adopt a two-stage strategy. Initially, we employ GPT-40 to extract objects and actions from the original text prompt. Subsequently, we employ GPT-40 to determine whether the extracted objects are present in the video and to verify the occurrence of specified actions. For each video, GPT-40 first assesses the presence of the objects mentioned in the prompt (e.g., "egg") within the video frames. This evaluation is performed according to Question 1 (Q1), where GPT-40 assigns a score from 0 to 2 based on the completeness of object presence: a score of 2 is given if all the objects are present, 1 if some of the objects are missing, and 0 if none of the objects appear in the video. After determining object presence, GPT-40 moves on to Question 2 (Q2) to check if the specified action (e.g., "*pour out*") is performed in the video. It assigns a binary score (0 or 1) depending on whether the action is present (1) or absent (0). Finally, these scores are combined to form the overall semantic alignment score. we put more details about other metric baselines in Appendix D.1.

810 C.2 Physical Commonsense alignment details

In this section, we use the same notation as in Section 4.2 and provide a more detailed description of the calculation and design of the method.

814

Key Phenomena Detection. We categorize the T2V prompts into monotonic processes (eg. 815 "melting with increasing temperature") and non-monotonic processes (eg. "an egg hitting a rock") 816 based on the physical phenomena they represent. For prompt with monotonic processes, we only 817 consider using the "Last Frame" as the retrieval prompt, resulting in a single question. We can di-818 rectly calculate $VLM(Img_i, Q)$, where the score for the corresponding video of this prompt ranges 819 from 0 to 1. For prompt with non-monotonic processes, we consider both the intermediate key 820 frames and the Last Frame, resulting in two questions. For the intermediate key frames, we calcu-821 late VLM(Img_i, Q) + VLM($\operatorname{Img}_i, P_r$), which ranges from 0-2. Consequently, the score range for 822 videos corresponding to this prompt is 0 to 3. 823

For specific calculatation, we need to calculate $VLM(I_j, p_r)$ and $VLM(I_j, q)$, where Img_j is the *j*-th frame in the video. For $VLM(I_j, p_r)$, the calculation involves assessing the matching degree between the key frame and the retrieval prompt, which can be directly obtained using the original calculation method in (Lin et al., 2024). For $VLM(I_j, q)$, we follow the computation approach from ChronoMagicBench (Yuan et al., 2024), we derive $VLM(I_j, q)$ by determining the ratio of the VQAScore for the affirmative statement to the combined VQAScores for both the affirmative and negative statements. We perform the calculations of $VLM(I_j, p_r)$ and $VLM(I_j, q)$ for each key frame within the specified range to obtain the physical correctness score for the problem.

831

832 Key Sequence Verification. For this stage, which we've primarily introduced in Section 833 4, we focus on key calculation points. The score calculation formula for q_1 is S_{before} = 834 $\max_{i=2 \le j \le i} (VLM(I_0, I_j, q_1) + VLM(I_j, p_r))$. Here, $VLM(I_j, p_r)$ determines if the retrieved key frame satisfies the retrieval prompt, as the physical phenomenon should occur in the keyframe pri-835 marily located in Key Phenomena Detection, which is crucial for Key Sequence Verification (e.g.the 836 expected physical phenomenon of egg cracking should occur in the keyframe when the egg hits 837 the stone, rather than other frames when the egg is in the air or else). $VLM(I_0, I_i, q_1)$ assesses 838 the correctness of the Key Sequence order in the video. Notably, we calculate $VLM(I_i, p_r)$ using 839 VQAScore, yielding a decimal between 0 and 1, while $VLM(I_0, I_j, q_1)$ employs VLM (GPT-4V or 840 LLaVA-Interleave) for question-answering, scoring 1 or 0 based on the model's Yes or No response. 841

Overall Naturalness Evaluation. Here we mainly explain how to get the score of this part based on the evaluation results under the two-stage strategy described in Section 4. Specifically, we ask the video-based VLM to select the most appropriate option for the video according to the detailed scoring criteria generated by the LLM, and then we map the options to scores (Completely Fantastical to Almost Realistic corresponds to 0-3 points)

Overall Score. We detail the discretization and calculation process of the scores here. In the stage of key phenomena detection, we categorize the prompts into monotonic and non-monotonic processes based on the physical phenomena they represent. For monotonic processes, the score range is 0-1, which we directly discretize by averaging into integer values from 0-3. Specifically, for non-monotonic processes with a score range of 0-3, we discretize the scores to [1, 1.5, 2.25]. This is because no points should be awarded if the physical phenomena are incorrect (VLM(I_j, p_r) = 1 and VLM(I_j, q) = 0), even with accurate retrieval. (e.g., The egg hits the stone and does not break)

In the stage of key sequence verification, we have three multi-image problems. One point is awarded for each correct answer, resulting in a final integer score from 0-3. Similar to the stage, of key phenomena detection we need to consider both the accuracy of key frame retrieval and the physical question answering. Therefore, we design the following: for Q_1 , when max_{i-2≤j≤i} (VLM(I₀, I_j, q₁) + VLM(I_j, p_r)) and VLM(I_j, p_r) > 0.5, the question is considered correct. The process for q_2 is similar. For q_3 , it is marked correct when VLM(I₀, I_{i-2:i+2}, I₋₁, q₃).

In the stage of overall naturalness evaluation, as we require video-based direct option selection,
 choosing Completely Fantastical, Clearly Unrealistic, Slightly Unrealistic, and Almost Realistic is
 scored as 0, 1, 2, and 3 points respectively. Finally, we average all scores and round down to obtain
 the final score.

Model	Duration (s)	FPS	Resolution
Open-Sora 1.2 (Zheng et al., 2024a)	4	24	1280×720
CogVideoX 2b	6	8	720×480
CogVideoX 5b	6	8	640×360
Lavie	4	8	512×320
Vchitect2.0	5	8	768×432
Pika (Pik, 2023)	3	24	1280×720
Gen-3 (gen, 2024)	11	24	1280×768
Kling (kli, 2024)	5	30	1280×720

Table 5: Details about evaluation models. The table shows duration, FPS, and resolution for each model.

For the ensamble operation, in order to reduce the bias caused by using a single VLM at a certain stage, we ensemble the results of PhyGenEval using open source models or closed source models. Specifically, we average the two results and round them down.

D EXPERIMENT

883 884 885

886

878

879

880

882

866 867 868

D.1 EXPERIMENTS SETUP

887 T2V model Implementation details. Open-Sora 1.2 (Zheng et al., 2024a) is an open-source project with the goal of reproducing Sora. CogVideoX 2b Yang et al. (2024) and CogVideoX 5b are 889 large-scale diffusion transformer models for text-to-video generation, incorporating a 3D Variational Autoencoder (VAE) for efficient video compression and an expert transformer with Expert Adaptive 890 LayerNorm to improve text-video alignment. LaVie Wang et al. (2023) is a cascaded video latent 891 diffusion model. Vchitect2.0 Wang et al. (2023), developed by the Shanghai AI Lab, is an advanced 892 video generation model featuring a Parallel Transformer architecture to scale up video diffusion 893 models and empower video creation. 894

895

896 **Evaluation Metrics details.** We compare our proposed PhyGenEval with some evaluation met-897 rics from previous methods like VideoPhy (Bansal et al., 2024) and VideoScore (He et al., 2024b). VideoPhy fine-tunes a VLM with the VIDEOPHY dataset proposed by themselves, which includes human feed back about the semantic alignment and dynamic motion correctness about videos. 899 VideoScore is trained on the VIDEOFEEDBACK dataset proposed by themselves, Initialized from 900 the Mantis model. VideoScore provides automatic assessments of video quality based on human 901 scoring criteria. To compare with *PhyGenEval* on SA and PCA, We only choose the text alignment 902 and fact consistency criteria. Specifically, for the semantic alignment evaluation, we compare the 903 Grid-LLaVA method proposed by T2V-CompBench, which extends the LLaVA (Liu et al., 2024a) 904 model to handle multi-frame inputs by sampling 6 frames uniformly from a video to create an im-905 age grid. For the physical commonsense alignment evaluation, we also compare with DEVIL (Liao 906 et al., 2024), which uses Gemini 1.5 Pro (Reid et al., 2024) to assess the overall naturalness of videos 907 and applies the same scoring standard prompt to all videos. 908

Furthermore, to evaluate the effectiveness of our *PhyGenEval* designs, we conduct a large amount of ablation studies and pue more details in Appendix D.3.

911

912 Human evaluation details. Here, we provide a detailed explanation of the human evaluation
913 described in Section 5. Specifically, we require annotators to score based on the standards outlined
914 in Figure 10, covering both semantic alignment and physical commonsense alignment. For example,
915 as for the video shown in Figure 10, The egg bounces off the rock like a rubber ball, completely
916 violating physical laws like dynamics, the annotator gives a score of 0 for physical commonsense
917 alignment. However, since the video fully includes the egg, the rock, and the collision action, the annotator gives a score of 3 for semantic alignment.

918 Table 6: SA correlation results with proposed PhyGenEval in video generation. A higher score 919 indicates better performance for a category. Bold stands for the best score, 920

Metric	Mech	Mechanics		Optics		Thermal		Material		Overall	
Wittite		$\rho(\uparrow)$	$\tau(\uparrow)$	$\rho(\uparrow)$	$\tau(\uparrow)$	$\rho(\uparrow)$	$\tau(\uparrow)$	$\rho(\uparrow)$	$\tau(\uparrow)$	$\rho(\uparrow)$	
VideoPhy (Bansal et al., 2024)	0.20	0.25	0.03	0.03	0.20	0.24	0.18	0.22	0.13	0.17	
VideoScore (He et al., 2024b)	0.14	0.16	-0.13	-0.14	0.23	0.02	0.02	0.02	0.05	0.05	
Grid-LLaVA (Sun et al., 2024)	0.39	0.43	0.45	0.49	0.30	0.33	0.22	0.26	0.35	0.39	
PhyGenEval (Grid-LLaVA)	0.35	0.38	0.46	0.48	0.41	0.44	0.42	0.45	0.42	0.44	
PhyGenEval	0.48	0.52	0.64	0.67	0.46	0.49	0.47	0.50	0.53	0.56	

Table 7: SA evaluation results with proposed PhyGenEval in video generation. Both machine and human evaluations indicate that most models achieve good semantic scores on PhyGenBench. This suggests that the scenarios in *PhyGenBench* are simple enough to clearly reflect physical phenomena.

Model	Size	$Mechanics(\uparrow)$	$Optics(\uparrow)$	$Thermal(\uparrow)$	Material(↑)	Average(↑)	$Human(\uparrow)$
CogVideoX (Yang et al., 2024)	2B	0.63	0.67	0.61	0.63	0.64	0.64
CogVideoX (Yang et al., 2024)	5B	0.78	0.88	0.78	0.64	0.78	0.78
Open-Sora V1.2 (Zheng et al., 2024a)	1.1B	0.73	0.85	0.82	0.73	0.79	0.70
Lavie (Wang et al., 2023)	860M	0.47	0.63	0.73	0.53	0.58	0.55
Vchitect 2.0 (Wang et al., 2023)	2B	0.92	0.89	0.77	0.74	0.84	0.84
Pika (Pik, 2023)	-	0.63	0.81	0.73	0.69	0.72	0.65
Gen-3 (gen, 2024)	-	0.84	0.93	0.82	0.78	0.85	0.86
Kling (kli, 2024)	-	0.88	0.91	0.87	0.74	0.85	0.89

QUANTITATIVE EVALUATION D.2

943 Comparison result about semantic alignment. Here we design a new baseline PhyGenEval 944 (Grid-LLaVA) to illustrate the superiority of the method, which uses the two-stage strategy pro-945 posed in PhyGenEval from Appendix C.1, but replaces the VLM with Grid-LLaVA proposed in 946 T2V-CompBench (Sun et al., 2024). As shown in Table 6, PhyGenEval achieves the highest cor-947 relation scores across all categories, demonstrating its effectiveness as a human-aligned semantic 948 commonsense correctness evaluator for PhyGenBench. Compared to other methods, PhyGenEval 949 consistently outperforms previous baselines like VideoPhy, VideoScore, and Grid-LLaVA. Specifically, *PhyGenEval* obtains an overall Kendall's τ of 0.53 and a Spearman's ρ of 0.56, surpassing 950 the Grid-LLaVA (τ : 0.35, ρ : 0.39). The results clearly show that our *PhyGenEval* design provides 951 a more accurate and reliable semantic commonsense evaluation in PhyGenBench. 952

953 954

955

956

957

958

929

930

931

942

Quantitative result about semantic alignment. As shown in Table 7, nearly all models achieve relatively high SA scores, whether evaluated by machines or humans. This suggests that the scenarios in PhyGenBench are relatively straightforward, making it easier to assess physical commonsense. Among all the models, Kling achieved the highest SA score, with a human evaluation score of 0.89, reflecting its strong instruction understanding and video generation capabilities.

- 959 D.3 ABLATION STUDY
- 960

961 The Component in *PhyGenEval* on physical commonsense alignment evaluation. We con-962 duct a series of ablation studies to demonstrate the necessity of our method design by examining 963 its correlation with human evaluation results, similar to those described in Section 5. Specifically, we compare: 1) The effectiveness of two-stage evaluation method proposed in Section 4.2 2) The 964 effect of the various stages of PhyGenEval, as proposed in Section 4.2; 3) Performance differences 965 when using various VLMs and their ensembles in PhyGenEval, as outlined in Section 4.2. 4) The 966 larger open models in PhyGenEval. Notice that PhyGenEval for physical commonsense alignment 967 evaluation consists of three stages: key phenomena Detection, key sequence verification, and overall 968 naturalness evaluation. And We denote them as PhyGenEval-S, PhyGenEval-M, and PhyGenEval-969 V based on the VLM they used. 970

1) We demonstrate that employing a two-stage strategy, as outlined in Section 4.2, yields superior 971 results when assessing the physical commonsense correctness of the entire video compared to one972 stage strategy. Specifically, the one-stage strategy refers to not using LLM to rewrite the scoring 973 template, but instead applying a single scoring template for all prompts' corresponding videos, al-974 lowing the VLM to score them. This method is proposed in DEVIL (Liao et al., 2024). To verify the 975 superiority of the two-stage strategy, we use InternVideo2 and GPT-40 as VLMs and perform both 976 the one-stage and two-stage strategies. We label these as PhyGenEval-V(Intern) and PhyGenEval-V(GPT-40), respectively. As shown in Table 8, the evaluation results produced by the two-stage 977 strategy are more consistent with human judgments for both InternVideo2 and GPT-40. We attribute 978 this improvement to the incorporation of LLM (GPT-40) for better comprehension of physical com-979 monsense text, which effectively reduces the complexity of the task for VLMs in evaluating the 980 physical correctness of videos. 981

2) *PhyGenEval* for physical commonsense alignment evaluation consists of three stages. We investigate the contribution of each stage to the final performance. Table 9 presents results using one or two stages (employing ensemble strategies when multiple VLMs are applicable). We find that optimal performance is achieved only when all three stages are used concurrently, demonstrating the rationale behind *PhyGenEval*'s design.

3) Given potential biases in single models and the costs associated with closed-source models, we offer two *PhyGenEval* computation methods: using GPT-40 or alternative open-source models (LLaVA-Interleave (Li et al., 2024) and InternVideo2 (Wang et al., 2024)). Table 10 shows that even using only small scale open-source models achieves a high correlation coefficient of 0.66. Notably, ensembling both methods yields the best results. Considering *PhyGenBench*'s relatively small size, we find this computational cost acceptable. Therefore we recommend users ensemble these methods.

994 4) We explore the performance of larger open-source models. Specifically, we replace LLaVA-995 Interleave used in the Physics Order Verification stage of PhyGenEval (Open-S) with InternVL-Pro (78B), denoted as PhyGenEval (Open-L). Additionally, we ensemble PhyGenEval (Open-L) with 996 PhyGenEval (Open-S), denoted as PhyGenEval (Open-Ensemble). Results show that compared to 997 smaller open-source models, the overall alignment coefficient improves from 0.66 to 0.72, indicating 998 that the method remains reproducible even when using exclusively open-source models. We believe 999 that as open-source models continue to advance, they can achieve even better performance within 1000 PhyGenEval. 1001

The Component in *PhyGenEval* on semantic alignment evaluation. we also perform necessary ablation experiments to validate the necessity of our SA evaluation design. Specifically, we compare: 1) VLM Model Selection: We leverage GPT-40 (Achiam et al., 2023) as a more robust VLM model for SA evaluation. 2) Effectiveness of our two-stage evaluation method proposed in Appendix C.1

1) As shown in Table 6, using GPT-40 in *PhyGenEval* is much better than using LLaVA, which achieve a higher Kendall's τ of 0.53 compared to 0.42, and a higher Spearman's ρ of 0.56 versus 0.44. This indicates a stronger alignment between GPT-40's evaluations and human annotations compared to open-source vlm models like Grid-LLaVA (Sun et al., 2024), justifying its selection as the preferred VLM model in the SA evaluation design. Since *PhyGenBench* includes a limited number of prompts, we believe that the cost of using GPT-40 is acceptable relative to the improvement in performance.

2) To validate the effectiveness of the two-stage strategy, we compare it with the method in T2V-1014 CompBench (Sun et al., 2024), which directly uses Grid-LLaVA to apply the same scoring standard 1015 prompt for semantic alignment evaluation across all videos. For fairness, we also use Grid-LLaVA 1016 but implement the two-stage strategy proposed in Appendix C.1. As shown in Table 6, PhyGenEval-1017 Grid-LLaVA outperforms Grid-LLaVA, achieving a higher Kendall's τ score of 0.42 compared to 1018 0.35, and a higher Spearman's ρ score of 0.44 versus 0.39. This result demonstrates the effective-1019 ness of our Two Stage Evaluation Method. By decomposing the evaluation into object detection 1020 and action detection, we effectively reduces the complexity of the task for VLMs in evaluating the 1021 sementic correctness of videos.

1022 1023

E DISCUSSION

Table 8: Comparison of PCA correlation results of the two-stage strategy for the video stage in *PhyGenEval*

Metric	Mech	Mechanics		Optics		Thermal		Material		erall
	$\tau(\uparrow)$	$\rho(\uparrow)$	$\tau(\uparrow)$	$\rho(\uparrow)$	$\tau(\uparrow)$	$\rho(\uparrow)$	$\tau(\uparrow)$	$\rho(\uparrow)$	$\tau(\uparrow)$	$\rho(\uparrow)$
			One S	tage Stra	itegy					
PhyGenEval-V(Intern) PhyGenEval-V(GPT)	$-0.03 \\ 0.39$	$-0.04 \\ 0.41$	$-0.20 \\ 0.11$	$-0.21 \\ 0.12$	$-0.26 \\ 0.19$	$-0.27 \\ 0.20$	$\begin{array}{c} 0.06 \\ 0.36 \end{array}$	$\begin{array}{c} 0.06 \\ 0.39 \end{array}$	$-0.10 \\ 0.19$	$-0.11 \\ 0.21$
			Two S	tage Stra	ntegy					
PhyGenEval-V(Intern) PhyGenEval-V(GPT)	$\begin{array}{c} 0.01 \\ 0.47 \end{array}$	$\begin{array}{c} 0.01 \\ 0.51 \end{array}$	$\begin{array}{c} 0.06 \\ 0.50 \end{array}$	$0.06 \\ 0.53$	$\begin{array}{c} 0.08\\ 0.46\end{array}$	$0.08 \\ 0.49$	$\begin{array}{c} 0.10\\ 0.53\end{array}$	$\begin{array}{c} 0.11 \\ 0.58 \end{array}$	$0.07 \\ 0.53$	$0.08 \\ 0.58$

Table 9: Comparison of PCA correlation results using each stage in PhyGenEval

	Metric	Mech	Mechanics		Optics		Thermal		Material		erall
		$\tau(\uparrow)$	$\rho(\uparrow)$								
	PhyGenEval-S	0.50	0.54	0.43	0.45	0.50	0.54	0.72	0.77	0.56	0.61
	PhyGenEval-M	0.46	0.49	0.49	0.53	0.55	0.59	0.53	0.57	0.55	0.60
	PhyGenEval-V	0.26	0.30	0.44	0.47	0.33	0.35	0.48	0.52	0.42	0.46
	PhyGenEval-SM	0.58	0.61	0.47	0.50	0.58	0.62	0.66	0.70	0.60	0.64
	PhyGenEval-SV	0.56	0.59	0.41	0.43	0.58	0.60	0.70	0.74	0.59	0.62
	PhyGenEval-MV	0.50	0.53	0.50	0.53	0.53	0.57	0.60	0.64	0.57	0.61
	PhyGenEval	0.72	0.75	0.76	0.77	0.73	0.75	0.81	0.84	0.78	0.81

Table 10: Comparison of PCA correlation results using different models such as GPT-40 or open-sourced models in *PhyGenEval*

Metric	Mech	Mechanics		Optics		Thermal		Material		erall
	$\tau(\uparrow)$	$\rho(\uparrow)$								
PhyGenEval (Open-S)	0.54	0.57	0.59	0.62	0.55	0.58	0.65	0.69	0.62	0.66
PhyGenEval (Open-L)	0.56	0.59	0.61	0.63	0.59	0.61	0.67	0.71	0.65	0.69
PhyGenEval (Open-Ensemble)	0.58	0.62	0.63	0.65	0.62	0.64	0.70	0.73	0.67	0.72
PhyGenEval (GPT40)	0.59	0.63	0.53	0.57	0.64	0.68	0.73	0.77	0.66	0.71
PhyGenEval (Ensemble)	0.72	0.75	0.76	0.77	0.73	0.75	0.81	0.84	0.78	0.81



Figure 7: Visualization of some PhyGenEval error cases.

Error case analysis. As shown in Figure 7, we visualize some error cases where both Phy-GenEval and competing methods like DEVIL fail to correctly identify the physical realism of the videos. These error cases are often caused by **confusing but iconic physical phenomena** in the videos that do not align with the correct progression of physical processes (e.g., in the erroneous case of the "burnt bread" experiment, black coloration appears but does not align with the expected phenomenon), leading to misjudgments. However, even in these cases, PhyGenEval remains closer to human ratings compared to other methods.

1087

1126

The Impact of Scaling on Physical Commonsense in Video Generation. Scaling laws have 1088 been extensively validated in video generation models (Kaplan et al., 2020). We investigate their 1089 efficacy in addressing the challenges of physical commonsense presented in PhyGenBench. As 1090 shown in Table 2, CogVideo 5B demonstrates improvements over CogVideo 2B, albeit with limited 1091 progress in the Mechanics category. Our qualitative analysis, illustrated in Figure 8, reveals signif-1092 icant advancements in static scenes with CogVideo 5B. It accurately captures complex phenomena 1093 such as colorful bubbles resulting from interference and diffraction, and oxidation-induced rusting 1094 of iron. In thermal, despite imperfections, CogVideo 5B generates more realistic boiling simula-1095 tions compared to its predecessor. However, both models struggle with simple motion dynamics, exemplified by their inability to accurately depict a bouncing football. We posit that while scaling up enhances the model's capacity to generate videos that align with physical commonsense for individual objects, it may be insufficient for physical phenomenons involving dynamic physical laws. 1099 Addressing these challenges likely requires extensive training on carefully curated synthetic data, as suggested by (Liu et al., 2024b). This approach could potentially bridge the gap in the model's 1100 grasp of fundamental physical laws. 1101



Figure 8: The qualitative comparison of CogVideoX 2B and CogVideoX 5B. The result shows that simply scaling up can solve some issues, but dynamic physical phenomenons involving the design of motion patterns remain challenging.

Rewriting prompt. We aim to explore whether GPT-augmented prompts can address the *Phy-GenBench* challenges. Specifically, we rewrite the original prompts using GPT, adding expected physical outcomes and processes. For example, after "*A bottle of juice is slowly poured out in the space station, releasing the liquid into the surrounding area*", we add "*The liquid forms floating globules, spreading out and moving randomly through the air.*" in the end.

As shown in Table 11, we use CogVideoX 5b and Kling as representative models for open-source and closed-source systems, respectively, to conduct tests. The results indicate that prompt rewriting

1134 Table 11: Evaluation results of PCA using the proposed *PhyGenEval* after rewriting prompts 1135 . The results indicate that although using rewritten prompts leads to some improvement, it is still 1136 insufficient to address the challenges highlighted by PhyGenBench.

Model	Size	Mechanics (↑)	$Optics(\uparrow)$	$Thermal(\uparrow)$	$Material(\uparrow)$	Average(†)
		Before Rewri	iting Prompt	t		
CogVideoX (Yang et al., 2024)	5B	0.39	0.55	0.40	0.42	0.45
Kling	-	0.45	0.58	0.50	0.40	0.49
		After Rewri	ting Prompt			
CogVideoX (Yang et al., 2024)	5B	0.39	0.62	0.53	0.52	0.52
Kling	-	0.50	0.64	0.61	0.48	0.56

1145 1146 1147

1148

1149

Table 12: PCA evaluation results with proposed PhyGenEval in videos after VEnhancer. The results indicate that employing VEnhancer fails to enhance the model's comprehension of physical commonsense.

Model	Size	$\textbf{Mechanics}(\uparrow)$	$Optics(\uparrow)$	$\textbf{Thermal}(\uparrow)$	$Material(\uparrow)$	$Average(\uparrow)$
Vchitect 2.0	2B	0.41	0.56	0.44	0.37	0.45
Vchitect 2.0 (Venhancer)	2B	0.41	0.56	0.42	0.38	0.45

1154

1155 does help the models generate images aligned with physical laws, but it is still far from resolving 1156 the issues highlighted by PhyGenBench. Both CogVideoX 5b and Kling exhibit some growth, but 1157 even for Kling, it only achieves a score of 0.56. This demonstrates that current models still severely 1158 lack the ability to accurately render physical scenes, and this deficiency cannot be easily resolved 1159 through simple prompt rewriting. To illustrate this issue more clearly, as shown in Figure 9, our qualitative analysis shows that rewriting prompts can only address simple issues (e.g., flame color 1160 reactions), but remains ineffective for more complex physical processes (e.g., egg breaking, stone 1161 sinking). 1162

1163

1164 The robustness of PhyGenBench and PhyGenEval. VEnhancer (He et al., 2024a) is a genera-1165 tive space-time enhancement framework that improves existing videos by adding spatial details and synthetic motion in the temporal domain. After enhancement by VEnhancer, Vchitect2.0 shows sig-1166 nificant improvement on VBench, even surpassing Kling. However, VEnhancer only enhances the 1167 visual quality of videos (e.g., making them more coherent and clear) without addressing the model's 1168 poor understanding of physical commonsense. 1169

1170 As shown in Table 12, Vchitect enhanced by VEnhancer still scores similarly to the original version 1171 on PhyGenBench. We calculate a high Spearman coefficient of 0.86 between model scores on PhyGenBench before and after VEnhancer enhancement. This indicates that PhyGenEval primarily 1172 focuses on physical correctness and is robust to other factors affecting visual quality. Furthermore, 1173 it demonstrates that even if a model can generate videos with better general quality (e.g., ranking 1174 higher on VBench), it doesn't necessarily imply a better understanding of physical common sense. 1175 This highlights the distinction between PhyGenBench and benchmarks like VBench that evaluate 1176 video quality. 1177

1178 1179

1180 1181 1182

F **COMPUTAIONAL RESOURCES**

Table 13: Resource consumption of models used in PhyGenEval.

1183					
1184	Model	Batch Size	Resources	Times	Memory Utilization Per GPU
1185	GPT-40(stage2)	8	USD 1.4	5min	-
1186	GP1-40(stage3) L1 aVA-Next-Interleave-7B	8	USD 3.1 1 x A100-80GB	2min	- 20408MiB
1187	VQAScore	3	1 x A100-80GB	10min	72726MiB
	InternVideo	1	1 x A100-80GB	1 min	7766MiB

Although our method involves different stages, it remains straightforward. Table 13 provides a summary of the resource consumption, showing that the entire evaluation process can be completed quickly and at low cost.



Figure 9: The qualitative comparison of effects before and after using rewritten prompts. The results indicate that rewriting prompts addresses only a few basic issues (such as flame color reactions), while the majority of problems remain unsolved.



Figure 10: Detailed diagram of the human evaluation process. We ask the annotators to score the semantic alignment and physical commonsense alignment of the video according to the scoring criteria in the figure.