

Objects in Generated Videos Are Slower Than They Appear: Models Suffer Sub-Earth Gravity and Don't Know Galileo's Principle...for now

Anonymous CVPR submission

Paper ID ****

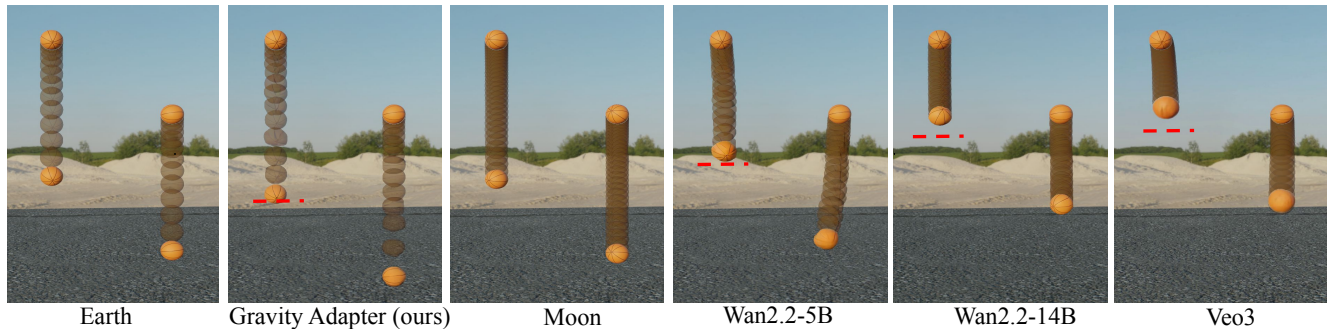


Figure 1. **Video generators produce physically implausible slow-motion falls and most fail to understand that objects fall at equal rates.** We visualize two identical balls dropped simultaneously from different heights using stroboscopic (time-lapse) composites, tracking motion until the lower ball impacts the ground. **Two failures emerge:** (1) **Galileo's principle violations:** Under Galileo's principle, both balls should fall equal distances in equal time regardless of starting height. Red dashed lines mark the expected position of the higher ball if both fell at equal rates. Earth (leftmost; simulation from Blender) shows the higher ball reaching this expected position, confirming correct physics. Moon reference ($g \approx 1.6 \text{ m/s}^2$) shows both balls falling slower but preserving equal-rate progression. In contrast, Wan 14B and Veo3 show the higher ball severely lagging behind the red line. It traveled far less distance than the lower ball despite falling for the same time, violating the fundamental principle that gravitational acceleration is universal. Even Wan 5B shows noticeable lag. (2) **Severe under-acceleration:** The spacing between successive ball positions indicates effective acceleration; wider spacing means higher acceleration. Most models exhibit compressed spacing comparable to Moon (1.6 m/s^2) or Mars (3.7 m/s^2) rather than Earth's 9.81 m/s^2 , revealing motion dramatically slower than terrestrial physics. Our Gravity Adapter (second panel), fine-tuned on Wan 5B with just 100 examples, corrects both failures by bringing the higher ball to the expected position and improving Earth-like spacing.

Abstract

001 Video generators are increasingly evaluated as potential
002 world models, which requires them to encode and under-
003 stand physical laws. We investigate their representation of
004 a fundamental law: gravity. Out-of-the-box video gener-
005 ators consistently generate objects falling at an effectively
006 slower acceleration. However, these physical tests are of-
007 ten confounded by ambiguous metric scale. We first in-
008 vestigate if observed physical errors are artifacts of these
009 ambiguities (e.g., incorrect frame rate assumptions). We
010 find that even temporal rescaling cannot correct the high-
011 variance gravity artifacts. To rigorously isolate the under-
012 lying physical representation from these confounds, we in-
013 troduce a unit-free, two-object protocol that tests the timing
014 ratio $t_1^2/t_2^2 = h_1/h_2$, a relationship independent of g , fo-
015 cal length, and scale. This relative test reveals violations of
016 Galileo's equivalence principle. We then demonstrate that
017 this physical gap can be partially mitigated with targeted

specialization. A lightweight low-rank adaptor fine-tuned
on only 100 single-ball clips raises g_{eff} from 1.81 m/s^2 to
 6.43 m/s^2 (reaching 65% of terrestrial gravity). This spe-
cialist adaptor also generalizes zero-shot to two-ball drops
and inclined planes, offering initial evidence that specific
physical laws can be corrected with minimal data.

1. Introduction

Two balls, a basketball and a tennis ball, are dropped from
the same height, side by side. Assuming no air resistance,
which hits the ground first? The answer is neither: they
strike the ground simultaneously. Galileo demonstrated this
principle over 400 years ago, and it remains one of the most
fundamental predictions of Newtonian mechanics [8, 24].
In this paper, we investigate whether state-of-the-art video
generators capture the fundamental physical law of gravity.

Galileo's principle is qualitative. A more rigorous, quan-
titative test is *how* objects fall: the time to impact is propor-

tional to the square root of the height ($t \propto \sqrt{h}$). But testing this in a generative model is not straightforward. Recent work shows that while models learn rich 3D representations [4, 9, 11, 37], this internal geometry is inherently scale-ambiguous, and video frame rates provide no absolute time reference. Any rigorous test of physical law must therefore eliminate these confounds.

Prior benchmarks audit physics through trajectory fitting [19], conservation laws [38], or broad qualitative suites [2, 17, 23, 34]. While valuable, these approaches either require camera calibration, rely on approximate heuristics, or optimize for absolute positional accuracy in a calibrated space – conflating the law of gravity with its parameterization. None directly measures whether a model’s representation of gravity quantitatively matches Earth’s $g = 9.81 \text{ m/s}^2$ in a way that is robust to the scale and time-base ambiguities inherent in generated video.

We introduce a unit-agnostic protocol based on relative timing between simultaneously falling objects. Because generated videos provide *no* guarantee that the model’s internal ‘meter’ aligns with physical units or that frame intervals correspond to absolute time, we must eliminate reliance on explicit units. Under a static pinhole camera with zero initial velocity, the ratio of impact times satisfies

$$\frac{t_1^2}{t_2^2} = \frac{h_1}{h_2},$$

where metric scale, focal length, time base, and gravitational constant all cancel—yielding a calibration-free diagnostic of physical understanding.

We find that, out of the box, most models violate Galileo’s principle that objects fall at the same rate; see Figure 1. In our single-ball drop experiments ($h \in [0.5, 4.0] \text{ m}$), impacts are consistently late across all tested video generators (Wan 5B, Wan 14B, Veo3, Cosmos 2B, Cosmos 14B). Interpreted in their native time gauge, the implied effective gravity g_{eff} ranges from 0.38 to 2.27 m/s^2 , far below Earth’s 9.81 m/s^2 ; balls therefore fall much more slowly than real-world gravity would predict, exhibiting an effective “sub-Earth” gravity. A straightforward temporal rescaling can move the *mean* g_{eff} closer to 9.81 m/s^2 , but the resulting distributions remain broad with large variance and heavy tails, indicating substantial inconsistency across scenes. Moreover, our two-object protocol reveals systematic, nonzero time differences Δt between the two balls when they traverse the same vertical distance. This relative-timing discrepancy is invariant to any global time rescaling and thus indicates a genuine violation of Galileo’s principle. Even more surprisingly, larger models (Wan 14B, Cosmos 14B) exhibit *slower* motion than their smaller counterparts (Wan 5B, Cosmos 2B), contradicting the assumption that scale alone improves physical consistency.

To repair this error, we train a lightweight LoRA adapter on just 100 single-ball sequences. Applied to Wan 5B, it

increases g_{eff} on average from 1.81 m/s^2 to 6.43 m/s^2 —reaching 65% of terrestrial gravity on average. Interestingly, this specialist generalizes zero-shot to two-ball drops and inclined planes despite never seeing these scenarios during LoRA training. In summary, objects in generated videos are slower than they appear. Video generators excel as generalists but struggle as physics engines; our adapter supplies targeted specialization with minimal data. Just as Galileo showed that the falling rate is independent of mass, we show that physical accuracy is independent of model scale; it requires targeted correction, not raw capacity.

Contributions. (1) We introduce a unit-free, two-object measurement protocol that isolates gravitational acceleration while being invariant to camera scale and frame rate, providing a rigorous diagnostic for physical consistency in video generators. (2) We quantify systematic physics violations across state-of-the-art models, revealing effective gravity at 5–20% of Earth’s value in their native time gauge and showing that larger models can be *less* physically accurate than smaller ones. (3) We show that a 100-example LoRA adapter on Wan 5B and Wan 14B partially improves these gravity metrics and shows generalization to unseen scenarios, demonstrating the feasibility of correcting specific physical laws with minimal data.

2. Related Works

Physics Evaluation in Video Generators. Recent benchmarks evaluate physics understanding through pixel-to-ground-truth matching [19, 23], human/VLM-as-a-judge [2, 13, 14, 20], or trajectory modeling [17, 38]. However, these approaches face a fundamental challenge: scale and time ambiguity. Monocular vision [10, 40] and neural rendering [1] cannot recover absolute metric scale without calibration. Video generators inherit this limitation: their internal representations are unitless, and frame rates provide no absolute time base [15]. Yet existing physics benchmarks implicitly assume calibrated measurements, conflating the physical law itself with its parameterization (e.g., testing $g = 9.81 \text{ m/s}^2$ specifically rather than the square-root scaling law). While VLMs enable broad qualitative scope [12, 22], they suffer from hallucinations and prompt shortcuts [22, 27], making them unreliable for quantitative physical laws. We address this by introducing unit-free relative measurements that isolate physical principles from metric parameterization.

Falling Objects. For free-fall specifically, PISA [19] optimizes trajectory IoU using 5,000 simulated videos with reward-based fine-tuning, while Morpheus [38] compares extracted trajectories against physics simulations. However, neither quantifies whether models understand gravity’s fundamental properties: the square-root time-height relationship ($t \propto \sqrt{h}$) or Galileo’s equivalence principle (all objects fall at equal rates). They only measure deviation

139 from ground-truth videos with implicit, unspecified gravity
140 values. Kang et al. [17] train 2D models from scratch to
141 test generalizability but similarly do not isolate the physical
142 law from scale ambiguity. Our two-ball protocol eliminates
143 these confounds by testing timing ratios ($t_1^2/t_2^2 = h_1/h_2$),
144 where scale, focal length, and gravity all cancel. This pro-
145 vides a calibration-free test of gravity.

146 **Probing Knowledge in Generative Models.** Recent
147 work investigates what visual knowledge generative mod-
148 els encode, including depth, normals, semantics, and object
149 relationships [6, 11, 37], as well as disentangled represen-
150 tations of lighting and geometry [4, 5, 9, 35]. Sarkar et
151 al. [28] revealed that diffusion models systematically fail
152 at projective geometry, mispredicting shadows and vanish-
153 ing points despite encoding 3D cues. Our work extends
154 this inquiry from static geometry to temporal dynamics, un-
155 covering systematic failures in gravitational physics. Criti-
156 cally, we find that a 100-example LoRA corrects these fail-
157 ures and generalizes zero-shot to unseen scenarios, suggest-
158 ing physical laws occupy sparse, learnable subspaces within
159 foundation models. This connects interpretability research
160 with physics-guided generation, showing that models en-
161 code latent physical knowledge that requires targeted acti-
162 vation rather than learning from scratch.

163 **Physics-Guided Generation.** Methods for improv-
164 ing physical correctness follow two paradigms: distilling
165 physics through post-training alignment [19, 39] or explic-
166 itly simulating scenes with physics engines [7, 21, 33, 36].
167 Post-training approaches like PISA require massive com-
168 pute, while simulation-based methods depend on accurate
169 scene estimation and cannot handle complex real-world sce-
170 narios. Both align with the intuitive physics hypothesis [3]
171 that physical understanding emerges from observing dy-
172 namics. However, our findings reveal that despite exposure
173 to vast video data, models learn correlational patterns with-
174 out internalizing fundamental rules. They systematically vi-
175 olate both acceleration laws and Galileo’s equivalence prin-
176 ciple. Unlike prior work, our lightweight LoRA adapter
177 achieves substantial correction and generalizes zero-shot to
178 two-ball drops and inclined planes, demonstrating that tar-
179 geted specialization can efficiently bridge the gap between
180 visual plausibility and physical correctness.

181 3. Experimental Setup

182 We investigate how well modern video generators capture
183 gravitational physics in the canonical setting of falling ob-
184 jects. Our evaluation requires controlled generation and
185 measurement; we describe our synthetic benchmark, the
186 models tested, and our tracking-based evaluation protocol.

187 3.1. Synthetic Benchmark

188 **Rendering.** Each video in our benchmark is created syn-
189 thetically using Blender. We use various standard sports

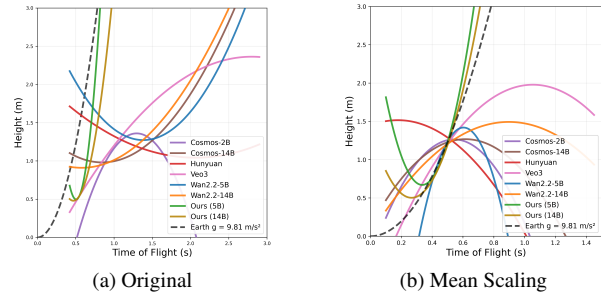


Figure 2. **Effect of time-scaling on $h-t$ relationships.** (a) We plot h versus t for all models. We repeat each test example with 4 seeds and fit polynomials through the means. The gray dashed line indicates terrestrial motion. All models systematically under-accelerate, and none obey the square root scaling law of time with height. The Gravity Adapters (green, gold) substantially improves Wan 5B and Wan 14B towards correct gravity. (b) **Mean time scaling.** We compute a Mean time scalar using a subset of random 30 samples from our dataset, which scales the effective time of the 30 samples to better match the ground truth time. The Mean time scalar, when applied to the second subset of 45 samples, brings the mean effective gravity closer to 9.81, m/s² for many models, but the variance remains high indicating that under-acceleration is not simply a frame rate artifact.

190 balls as falling objects, which should be well-represented in
191 the models’ training data. Balls are simulated to be dropped
192 from random heights in the range [0.5, 4.0] m, and the ini-
193 tial height is recorded. Videos are rendered at 1280×720
194 resolution for 2 seconds at 24 fps, yielding 48 frames per
195 sequence.

196 We generate N videos with heights uniformly sampled
197 from [0.5, 4.0] m. The dataset includes K different ball
198 types (basketball, soccer ball, tennis ball, volleyball, base-
199 ball) across M unique HDRI environments. For the adapter
200 training (Sec. 5), we use 100 single-ball sequences; all other
201 experiments use held-out test data.

202 The camera is fixed at a sufficient distance and height
203 to capture the entire trajectory of the ball and is positioned
204 perpendicular to the falling motion. The scene is lit, and
205 the background is occupied by HDRI maps. We use indoor
206 and outdoor HDRI maps with a variety of different ground
207 materials to match the setting. Backgrounds are chosen to
208 match the scale of the ball—objects in the background do
209 not appear too large or too small compared to the ball. We
210 include samples with backgrounds having objects, to pro-
211 vide more context, as well as samples with more plain back-
212 grounds that provide less context for more challenging gener-
213 ations (additional objects provide scale context).

214 The camera uses a focal length of $f = 50$, mm and is po-
215 sitioned between 1.0 and 8.0 meters from the drop zone. Its
216 height is set to half the drop height, with an additional offset
217 sampled from the range $[-0.5, 0.5]$ meters. This configura-
218 tion ensures that perspective distortion remains minimal

Table 1. **Effect of time-scaling on g_{eff} estimation across models.** (a) **Original.** We report effective gravity values computed as $g_{\text{eff}} = 2h/t^2$ (m/s^2). The ground truth is 9.81 m/s^2 . All models under-accelerate, and Gravity Adapters consistently reduce this deficit. Reported mean values are averaged over four random seeds and all test examples. Median and Range values are across all seeds and test samples. (b) **Mean time scaling.** A global scalar (MTS) is estimated from a 30-sample subset and applied to a disjoint 45-sample split. This shifts several models closer to 9.81 m/s^2 , but variance remains high.

(a) Original				(b) Mean-Time Scaled				
Model	Mean (m/s^2)	Median (m/s^2)	Range (m/s^2)	Model	MTS	Mean (m/s^2)	Median (m/s^2)	Range (m/s^2)
Cosmos 2B [30]	1.85	1.30	[0.23, 14.18]	Cosmos 2B [30]	2.43	10.34	7.24	[1.40, 69.96]
Cosmos 14B [30]	1.51	1.01	[0.24, 10.31]	Cosmos 14B [30]	2.77	10.86	7.19	[1.81, 76.04]
Hunyuan [18]	1.97	1.15	[0.23, 15.84]	Hunyuan [18]	2.63	13.85	7.06	[1.55, 104.64]
Veo3 [26]	2.27	2.08	[0.28, 6.66]	Veo3 [26]	2.12	9.39	8.5	[1.26, 29.11]
Wan 5B [31]	1.81	1.24	[0.26, 8.26]	Wan 5B [31]	2.43	10.24	7.22	[1.76, 49.15]
Wan 14B [31]	2.18	1.19	[0.27, 59.98]	Wan 14B [31]	2.56	13.78	7.78	[1.75, 321.30]
Gravity Adapter 5B [16]	6.43	6.38	[1.24, 16.64]	Gravity Adapter 5B [16]	1.28	10.27	9.74	[2.4, 26.67]
Gravity Adapter 14B [16]	5.51	5.63	[1.52, 11.67]	Gravity Adapter 14B [16]	1.36	10.00	10.10	[2.82, 19.08]

219 across the falling trajectory.

220 **Single-Ball Protocol.** We first quantify the effective grav- 255
 221 itational acceleration that models implicitly use by measur- 256
 222 ing fall time vs. height. This establishes a baseline before 257
 223 we eliminate scale ambiguity with our unit-free two-ball 258
 224 protocol. We generate 75 sequences with drop heights uni- 259
 225 formly sampled from $[0.5, 4.0]$ m. This tests whether mod- 260
 226 els produce trajectories consistent with $t = \sqrt{2h/g}$, allow- 261
 227 ing us to measure effective gravity $g_{\text{eff}} = 2h/t^2$. 262

228 **Two-Ball Protocol.** To eliminate metric scale and time- 263
 229 base ambiguity, we generate 50 sequences with two (iden- 264
 230 tical) balls dropped simultaneously from different heights 265
 231 h_1, h_2 within the same frame. To ensure both balls experi- 266
 232 ence identical perspective effects, they are positioned at the 267
 233 same distance from the camera and separated horizontally 268
 234 by ball-diameter with an offset of 0.5 meters to prevent oc- 269
 235 clusion. This guarantees that any scale ambiguity affects 270
 236 both balls equally, making the timing ratio truly unit-free. 271
 237 Height ratios h_1/h_2 range from 0.25 to 3.5. The relative 272
 238 timing should satisfy $t_1^2/t_2^2 = h_1/h_2$, which cancels grav- 273
 239 ity, scale, focal length, and frame rate—providing a unit- 274
 240 free test of physical understanding. 275

241 3.2. Models Evaluated

242 We evaluate models with Image-to-Video capability to con- 276
 243 strain the generations to fall from a known height specified 277
 244 using a conditioning image. Experiments are conducted on 278
 245 Wan 5B and Wan 14B [31], Veo3 [26], Hunyuan [18], and 279
 246 Cosmos 2B and Cosmos 14B [30]. For each model, we sup- 280
 247 ply the first frame of the video and a text prompt as initial 281
 248 conditions. We use consistent text prompts across models 282
 249 to isolate physics understanding from prompt sensitivity. 283

250 3.3. Measurement Protocol

251 Generated videos are resized and adjusted to 1280×720 , 284
 252 24 fps. We track the ball using SAM2 [25], initialized 285
 253 with Blender centroids. Impact time t is determined when: 286
 254 (1) the bottom of the ball drops below a point located one 287

255 ball-radius (in pixels) above the ground threshold y_{ground} , 256
 257 and (2) vertical velocity drops below $\varepsilon = 1$ pixels/frame. 258
 259 The ground threshold prevents false detections for hover- 260
 261 ing balls; the velocity threshold accounts for impact defor- 261
 262 mation. For single-ball drops, we compute $g_{\text{eff}} = 2h/t^2$ 262
 263 and report mean, median and range. For two-ball drops, we 263
 264 compute timing ratios t_1^2/t_2^2 and compare against the theo- 264
 265 retical h_1/h_2 . We manually verify extreme acceleration re- 265
 266 sults above 50 m/s^2 and exclude any samples that are found 266
 267 to be incorrectly evaluated because of SAM failures. 267

268 4. Results

269 4.1. Single-Ball Drops: Sub-Earth Gravity

270 Figure 2a plots h versus t for all the models. Under correct 270
 271 physics, points should lie on the curve $t^2 = 2h/g$. All the 271
 272 models deviate from this reference. 272

273 Table 1a summarizes effective gravity measurements. 273
 274 No model reproduces Earth-like acceleration. Smaller mod- 274
 275 els show less extreme under-acceleration (closer to Earth 275
 276 gravity) than larger models, though all remain far below 276
 277 $g = 9.81 \text{ m/s}^2$. Paradoxically, scaling increases capaci- 277
 278 ty but does not improve physical accuracy. In the case 278
 279 of Wan 14B, even though it reports a higher mean value, 279
 280 the median and the corresponding distribution reveal that 280
 281 the majority of generations remain substantially under- 281
 282 accelerated. The increased mean is primarily driven by a 282
 283 small number of extreme samples. Our lightweight Grav- 283
 284 ity Adapter, fine-tuned on just 100 sequences (Sec. 5), 284
 285 improves Wan 5B from $g_{\text{eff}}=1.81$ to $g_{\text{eff}}=6.43$ (range: 285
 286 $1.24 - 16.64 \text{ m/s}^2$), demonstrating that targeted specializa- 286
 287 tion can partially correct gravity. 287

288 **Does prompting help? No.** To reduce any scale ambi- 288
 289 guity for the models, we experiment with more detailed text 289
 290 prompts containing explicit height, diameter of the ball, dis- 290
 291 tance of the camera from the ball, and height of the camera 291
 above the ground. However, we see no significant improve-
 ment in any of them (Tab 2).

Is under-acceleration a frame rate artifact? No. We first

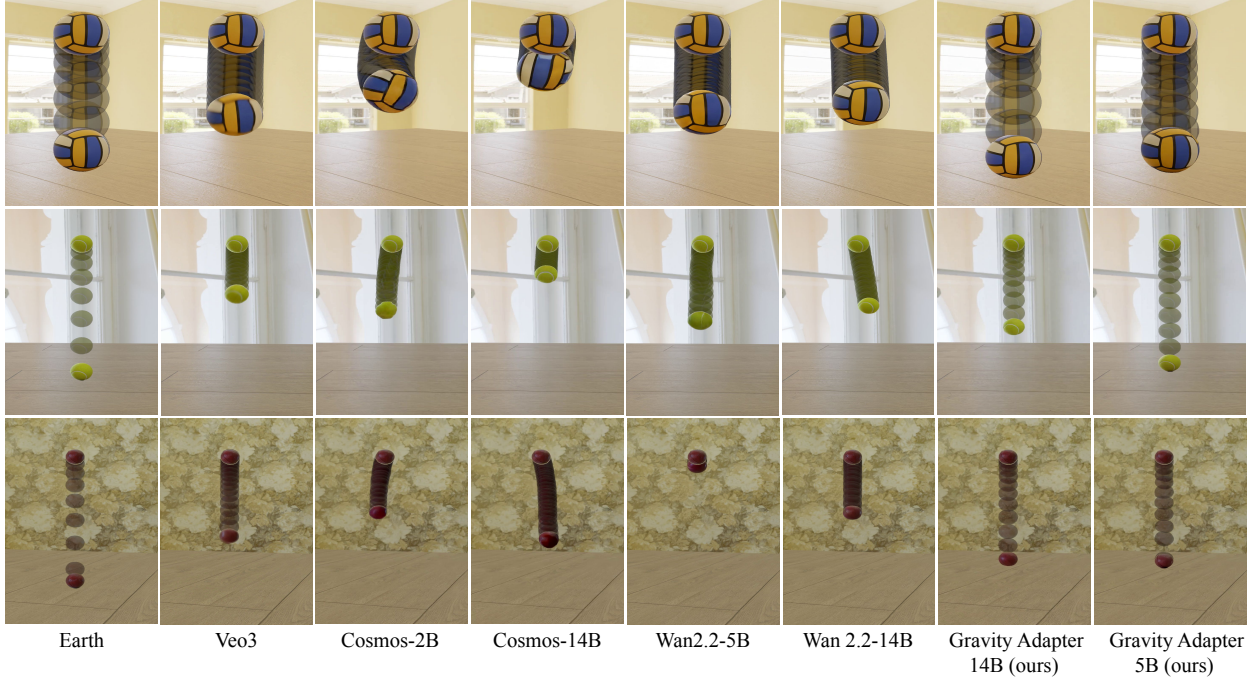


Figure 3. **Scaled Single-ball drops reveal systematic under-acceleration across all models.** Stroboscopic composites (left) visualize ball positions at equal time intervals from release. The panels show the trajectories performed by each model during the time it takes a ball falling under $9.8m/s^2$ to reach the ground, scaled by MTS (Tab. 1b). All models showcase severe under-acceleration (easily visible in the compressed spacing in the composites). The Gravity Adapters (seventh and eighth column) substantially improves Wan 5B and Wan14B toward terrestrial dynamics.

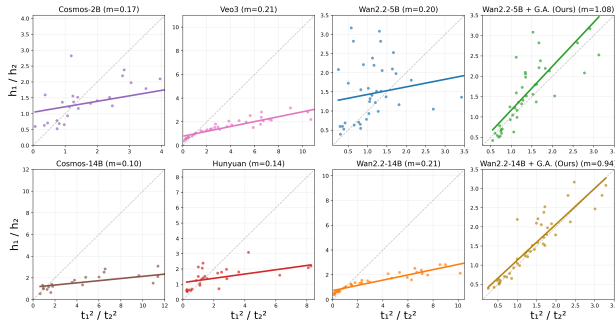


Figure 4. **Two-ball relative timing results.** We plot measured timing ratios t_1^2/t_2^2 against theoretical predictions h_1/h_2 across multiple height ratios. The gray dashed line indicates perfect agreement. All models deviate systematically, confirming that under-acceleration is not an artifact of scale estimation but reflects genuine physics and Galileo’s principle violations. We also measure the slope (m) for each of them to understand the deviation.

292 test whether temporal stretching explains the low gravity
 293 values. Using the single-ball dataset, we split 75 samples
 294 into two subsets (30 and 45 videos). From the first sub-
 295 set, we compute each model’s average time-scaling factor
 296 (MTS) as $\text{mean}(t_{eff}/t_{gt})$. We then apply this global scal-
 297 ing factor to videos in the second subset ($1/MTS \times t_{eff}$),
 298 computing scaled gravity values g_{scaled} . Tab. 1b shows the
 299 results across samples and seeds. Since the absolute scale

Table 2. **Use of expanded prompts.** Detailed prompts describing the scene with explicit parameters do not significantly affect the model’s understanding of gravity. Wan 5B increase marginally in g_{eff} . Veo 3 and Cosmos 2B show a decline.

Model	Base Prompt			Expanded Prompt		
	Mean (m/s^2)	Median (m/s^2)	Range (m/s^2)	Mean (m/s^2)	Median (m/s^2)	Range (m/s^2)
Wan 5B	1.81	1.24	[0.26, 8.26]	2.06	1.37	[0.29, 6.80]
Veo3	2.27	2.08	[0.28, 6.66]	1.63	1.24	[0.34, 4.94]
Cosmos 2B	1.85	1.30	[0.23, 14.18]	1.25	0.83	[0.27, 3.84]

is unknown, we treat the single-ball experiment as a test
 of temporal consistency. If the model simply had a slow
 clock, a linear time scalar should fix the error. As Tab. 1b
 shows, it does not. This failure necessitates a unit-free met-
 ric to understand the root cause. See Fig. 3 for qualita-
 tive examples demonstrating under-acceleration. We also
 tested other scaling baselines (see supplement): (1) per-
 scene mean-time scaling computed across multiple seeds,
 and (2) height-adjusted scaling for non-vertical trajec-
 tories. Neither approach yields substantial improvement, con-
 firming that physics errors persist even with privileged per-
 sample correction.

4.2. Two-Ball Drops: Failure of Galileo’s Principle

Prior work evaluating physics in video generators[19, 38]
 predominantly tests single-object scenarios, comparing
 generated trajectories against ground truth. While such ap-
 proaches can measure whether models follow text prompts

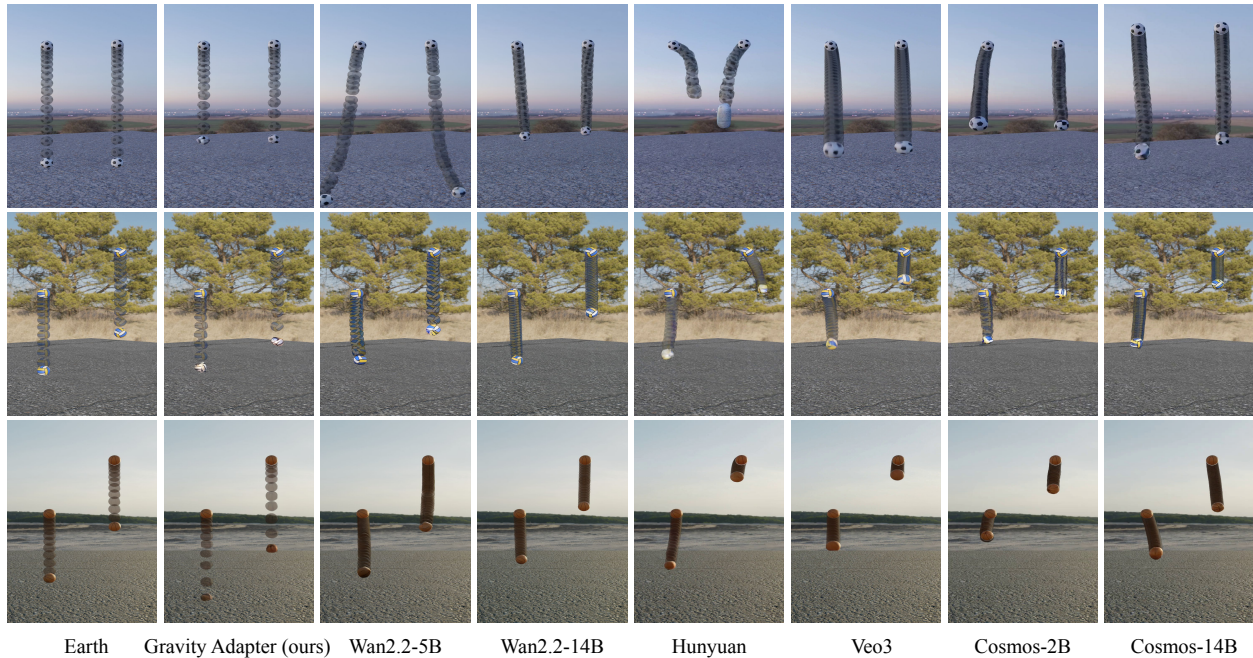


Figure 5. **Most models fail Galileo’s principle of gravitational equivalence.** We freeze at the moment the lower ball (from height h_1) impacts the ground. Under correct physics, both balls should have fallen equal distances in equal time, demonstrating that gravitational acceleration is universal (leftmost column: Ground Truth Earth). **Catastrophic failures:** Veo3 (sixth column) shows the higher ball barely moving while the lower ball lands—a complete violation of 400-year-old physics. Wan 14B, Veo3, and Cosmos 14B show similar failures with the higher ball remaining significantly elevated, suggesting these models believe in gravity depends on the starting height or object ordering. **Correction:** The Gravity Adapter (second column) finetuned on Wan 5B restores near-terrestrial acceleration *and* perfects equal-rate falling, demonstrating that this fundamental physics deficit can be corrected with only 100 training examples.

Table 3. **Time deviation statistics (Δt in frames) for the two-ball falling experiment.** The lower ball hits the ground at time t_1 after traversing a vertical distance y_1 . We then measure the time t_2 for the upper ball to traverse the *same* distance and report $\Delta t = t_2 - t_1$ (ideal constant- g motion gives $\Delta t = 0$). Positive values (e.g., Wan 14B, Veo3, Cosmos 14B) mean the upper ball is *slower*, indicating under-acceleration relative to the lower ball; negative values (Wan 5B, Cosmos 2B) mean it is *faster*, indicating over-acceleration. Interestingly, the *direction* of the violation flips with scale: smaller models tend to have negative mean Δt , while larger models have positive mean Δt . Gravity Adapters move the means close to zero (e.g., -0.95 for 5B and 0.14 for 14B), thereby shrinking the ranges, which largely cancels this scale-dependent Galileo violation.

Model	Mean (Δt)	Median (Δt)	Range (Δt)
Cosmos 2B [30]	-2.77	-2.0	[-35.0, 23.0]
Cosmos 14B [30]	4.03	5.0	[-36.0, 30.0]
Veo3 [26]	6.71	7.00	[-42.0, 28.0]
Wan 5B [31]	-4.22	-4.0	[-35.0, 30.0]
Wan 14B [31]	2.22	1.0	[-14.0, 16.0]
Gravity Adapter 5B [16]	-0.95	-1.00	[-8.0, 4.0]
Gravity Adapter 14B [16]	0.14	0.0	[-5.0, 3.0]

317 accurately, they cannot determine whether models under-
 318 stand the underlying physical principles. A model might
 319 correctly match a single falling trajectory through a pattern

recognition without understanding gravity. 320

We pose a more fundamental question: Do models un- 321
 322 derstand Galileo’s principle of gravitational equivalence –
 323 that all objects fall at the same rate regardless of mass
 324 or starting height? To test this, we generate scenes with
 325 two identical balls dropped simultaneously from different
 326 heights $h_1 < h_2$. We use identical balls to eliminate depth-
 327 perception confounds that could arise from size differences.
 328 Under correct physics, both balls experience identical ac-
 329 celeration and thus fall equal distances in equal time. At
 330 the moment when the lower ball impacts at time $t_1 = \sqrt{\frac{2h_1}{g}}$,
 331 the higher ball should have fallen the same distance $d = h_1$,
 332 placing it at height $h_2 - h_1$. **This is a zero-ambiguity test:**
 333 either both balls fall together (pass) or they don’t (fail).

Figure 5 reveals that **most models catastrophically fail.** 334
 335 Veo3, Wan 14B, and Cosmos 14B show the higher ball
 336 traveling an unequal distance compared to the lower ball,
 337 when the lower ball lands. This observation is confirmed
 338 by plotting t_1^2/t_2^2 versus h_1/h_2 across multiple height ra-
 339 tios as shown in Figure 4. This is not just a uniform timing
 340 error; it reflects a failure to respect that gravitational ac-
 341 celeration is universal. Our two-ball experiments show that
 342 the upper and lower balls often take different times to tra-
 343 verse the same distance: in some models, the higher ball is

344 effectively slower, in others faster, but in all cases the two
345 objects experience different accelerations within the same
346 scene, contradicting four centuries of physics.

347 To quantify deviations from Galileo’s principle we mea-
348 sure the temporal lag between the two falling balls in units
349 of video frames. First, we record the time at which the
350 lower ball reaches the ground, denoted as t_1 , having cov-
351 ered a vertical pixel distance y_1 . We then measure the time
352 t_2 taken by the upper ball to traverse the same distance, and
353 compute the deviation $\Delta t = t_2 - t_1$. A positive deviation
354 indicates that the upper ball requires more time to cover the
355 same distance, reflecting under-acceleration relative to the
356 lower ball; and a negative deviation implies the opposite, an
357 over-acceleration. The measured deviations for each model
358 are reported in Tab 3. We find an interesting pattern: the
359 smaller models - Wan 5B and Cosmos 2B show an opposite
360 trend compared to their larger counterparts, with the higher
361 ball accelerating faster than the slower ball. As shown in
362 Tab. 3. Our gravity adapters (Sec. 5) noticeably improve
363 both Wan-5B and Wan-14B, bringing their behavior closer
364 to the physically correct regime in which gravity acts uni-
365 versally, not separately on each object.

366 5. Gravity Adapter: Targeted Specialization

367 Video generators excel as generalists but struggle as physics
368 engines. We investigate whether these physical deficits are
369 due to a lack of model capacity or simply a lack of align-
370 ment. By training a lightweight adapter, we probe the learn-
371 ability of gravitational physics.

372 5.1. Training Protocol

373 We train a LoRA [16, 29] on Wan 5B to correct its sys-
374 tematic under-acceleration. The training set consists of 100
375 single-ball drop sequences rendered in diverse HDRI envi-
376 ronments distinct from our test benchmark. We train for
377 5,000 iterations with rank $r = 32$, learning rate 10^{-4} . Total
378 training time is approximately 6 hours on two A100 GPUs.

379 5.2. Results on Benchmark Tasks

380 **Single-Ball Drops.** Figure 2a shows that the adapter suc-
381 cessfully shifts the mean effective gravity significantly
382 closer to Earth’s standard ($1.81 \rightarrow 6.43m/s^2$). While the
383 distribution remains broad (indicating that stochastic gener-
384 ation artifacts persist), the shift demonstrates that the model
385 can be realigned to physical laws with negligible data (100
386 examples).

387 **Two-Ball Drops (Zero-Shot).** We also find that the
388 adapter shows zero-shot transfer to two-ball relative tim-
389 ing tasks without ever seeing two-object scenarios during
390 training (Figure 4). Timing ratios move substantially closer
391 to the theoretical $t_1^2/t_2^2 = h_1/h_2$ prediction, indicating the
392 adapter learned general gravitational dynamics rather than
393 memorizing single-ball trajectories.

Table 4. **Zero-shot generalization to real-world PISA bench-
mark.** The Gravity Adapter (5B), trained only on 100 syn-
thetic sequences, improves trajectory accuracy on 361 real-world
videos [19]. The marginal improvement when full finetuning (FT)
and reward optimization (ORO)[19] does not justify its enormous
training cost.

Model	L2 ↓	Chamfer ↓	IoU ↑	Params
Wan 5B (baseline) [31]	0.148	0.413	0.073	—
Wan 5B + Gravity Adapter [16]	0.127	0.346	0.086	80M
Wan 5B + FT [19] + ORO [19]	0.123	0.332	0.089	5B

Real-World Transfer: PISA Benchmark. To test
whether the adapter overfits to synthetic training data, we
evaluate zero-shot on the PISA benchmark [19]—361 real-
world videos of objects falling onto cluttered surfaces. This
is challenging: our adapter was trained only on synthetic
scenes with sports balls and clear ground planes. Table 4
shows that the adapter over the baseline across all metrics.

5.3. Comparison with Alternative Approaches

Table 5 validates our adaptor against three alternatives.

LoRA Rank Ablation We increase the LoRA rank from
32 to 64 to examine whether additional capacity improves
 g_{eff} . Despite doubling the parameter count (80M \rightarrow
160M), we observe no meaningful improvement, indicat-
ing that effective gravity cannot be enhanced through scal-
ing alone. Lower-rank variants (8 and 16) further degrade
performance.

Explicit Guidance Models We test whether providing
more information improves physical accuracy. The First-
Last Frame model [31] receives initial and final frames to
reduce depth ambiguity, but performs poorly—likely con-
strained to 5-second generations that default to slow motion.
We also evaluate trajectory matching [32] using ground-
truth 2D centroid trajectories from drop height to impact, af-
ter which we let the model freely generate. This model sim-
ilarly underperforms ($g_{eff} = 0.38 m/s^2$, range 0.04–3.22),
suggesting difficulty handling rapid motion. Our adapter
exceeds both methods without any explicit guidance.

Full Model Fine-Tuning Replicating PISA’s proto-
col [19], we fine-tune all Wan 5B parameters on our 100-
example set using their two-stage approach: 5,000 su-
pervised iterations followed by 1,000 reward-model itera-
tions [25]. Due to computational constraints, we reduce res-
olution to 480p and 10 denoising steps. Our LoRA matches
or exceeds full fine-tuning performance (Fig. 6) while train-
ing only 1% of parameters and using an order of magni-
tude less compute, suggesting physics correction may be
amenable to low-rank adaptation.

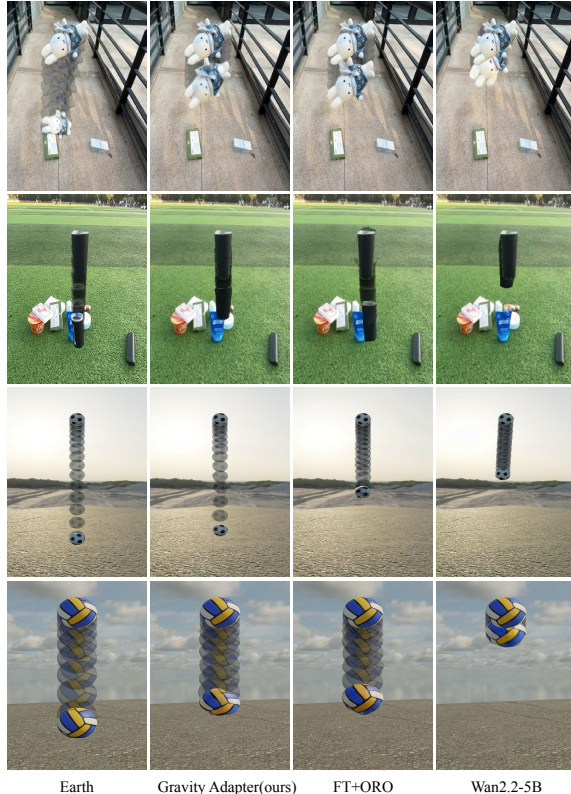


Figure 6. **Comparison with Full FT + ORO.** The Gravity Adapter (5B) shows impressive performance on both the real world dataset [19] and our benchmark with comparable and sometimes better performance at reduced computational cost.

Table 5. The Gravity Adapter (5B) exceeds full fine-tuning performance on our benchmark with less training computation and outperforms methods with additional guidance.

Method	Mean (m/s ²)	Median (m/s ²)	Range (m/s ²)
Wan 5B (baseline) [31]	1.81	1.24	[0.26, 8.26]
Baseline + FT [19] + ORO [19]	4.07	3.63	[0.89, 10.72]
Gravity Adapter (ours, rank 8) [16]	3.07	2.83	[1.24, 16.64]
Gravity Adapter (ours, rank 16) [16]	5.64	5.72	[1.51, 12.05]
Gravity Adapter (ours, rank 32) [16]	6.43	6.38	[1.24, 16.64]
Gravity Adapter (ours, rank 64) [16]	6.11	5.72	[1.91, 14.17]
FLF (first + last frame) [31]	0.58	0.22	[0.03, 23.70]
Trajectory guided [32]	0.38	0.16	[0.04, 3.22]

431 5.4. Generalization to Canonical Scenarios

432 To test whether the adapter learned general gravitational
433 principles or merely ball-dropping heuristics, we evaluate
434 on canonical physics regimes never seen during training:

435 **Inclined Planes.** We generate 12 sequences of smooth
436 cubes sliding down frictionless inclines at angles from 30°
437 to 75°. Under correct physics, acceleration should scale as
438 $g \sin(\theta)$. Figure 7 shows the Gravity Adapter substantially
439 improves motion realism compared to baseline Wan 5B,
440 with objects accelerating appropriately for their incline
441 angle. This demonstrates that the adapter internalized gravita-
442 tional principles applicable beyond vertical free fall.

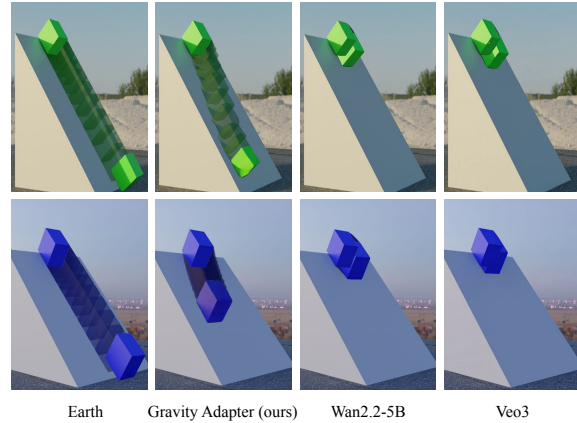


Figure 7. **Zero-shot generalization to inclined planes.** The Gravity Adapter (5B) improves motion realism for cubes sliding down inclines at various angles, despite never seeing this scenario during training. This demonstrates learning of general gravitational principles rather than ball-specific heuristics.

443 6. Discussion

444 We evaluate video generators on gravitational physics,
445 revealing systematic failures in their native time gauge.
446 Single-ball experiments show all models dramatically
447 under-accelerate ($g_{\text{eff}} \sim 1-2 \text{ m/s}^2$ vs. Earth’s 9.81 m/s^2), so
448 balls fall much more slowly than real-world gravity would
449 predict, exhibiting an effective “sub-Earth” gravity. A sim-
450 ple global temporal rescaling can move the *mean* g_{eff}
451 closer to 9.81 m/s^2 , but the resulting distributions remain broad
452 with large variance and heavy tails. Our unit-free two-ball
453 protocol eliminates scale ambiguity by testing timing ra-
454 tios ($t_1^2/t_2^2 = h_1/h_2$) and time deviations Δt for travers-
455 ing the same distance. Most models catastrophically fail
456 Galileo’s principle, with different objects effectively ex-
457 perencing different gravity within the same scene; these
458 relative-timing violations are invariant to any global time
459 rescaling. Surprisingly, larger models perform worse on
460 these metrics, contradicting naive scaling-law expectations.
461 A lightweight LoRA adapter trained on 100 synthetic ex-
462 amples substantially corrects these failures (e.g., $1.81 \rightarrow$
463 6.43 m/s^2) and moves Δt close to zero with tighter ranges,
464 and it generalizes zero-shot to real-world videos, two-ball
465 drops, and inclined planes.

466 **Limitations and Future Work.** Our adapter achieves
467 only 65% of terrestrial gravity on average with high vari-
468 ance ($1.24-16.64 \text{ m/s}^2$). We evaluate only vertical free-fall
469 and simple inclined motion, not projectile trajectories, ro-
470 tation, collisions, or friction. Future directions include: (1)
471 Broader evaluation: Apply unit-free protocols to momen-
472 tum, energy conservation, and fluid dynamics. (2) Physics-
473 informed losses: Incorporate $d = \frac{1}{2}gt^2$ and relative-timing
474 penalties directly into training objectives. (3) Multi-law
475 adapters: Train single adapters for multiple physical phe-
476 nomena.

477

References

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

- [1] Joshua Ahn, Haochen Wang, Raymond A Yeh, and Greg Shakhnarovich. Alpha invariance: On inverse scaling between distance and volume density in neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20396–20405, 2024. 2
- [2] Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chenfanfu Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. Videophy: Evaluating physical commonsense for video generation, 2024. 2
- [3] 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. 3
- [4] Anand Bhattad, Daniel McKee, Derek Hoiem, and David Forsyth. Stylegan knows normal, depth, albedo, and more. *Advances in Neural Information Processing Systems*, 36:73082–73103, 2023. 2, 3
- [5] Anand Bhattad, James Soole, and David A Forsyth. Stylitgan: Image-based relighting via latent control. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4231–4240, 2024. 3
- [6] Anand Bhattad, Konpat Preechakul, and Alexei A Efros. Visual jenga: Discovering object dependencies via counterfactual inpainting. *arXiv preprint arXiv:2503.21770*, 2025. 3
- [7] Boyuan Chen, Hanxiao Jiang, Shaowei Liu, Saurabh Gupta, Yunzhu Li, Hao Zhao, and Shenlong Wang. Physgen3d: Crafting a miniature interactive world from a single image. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 6178–6189, 2025. 3
- [8] Stillman Drake. *Galileo at work: His scientific biography*. Courier Corporation, 2003. 1
- [9] Xiaodan Du, Nicholas Kolkin, Greg Shakhnarovich, and Anand Bhattad. Generative models: What do they know? do they know things? let’s find out! *arXiv preprint arXiv:2311.17137*, 2023. 2, 3
- [10] David Eigen, Christian Puhersch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. *Advances in neural information processing systems*, 27, 2014. 2
- [11] Mohamed El Banani, Amit Raj, Kevis-Kokitsi Maninis, Abhishek Kar, Yuanzhen Li, Michael Rubinstein, Deqing Sun, Leonidas Guibas, Justin Johnson, and Varun Jampani. Probing the 3D Awareness of Visual Foundation Models. In *CVPR*, 2024. 2, 3
- [12] Quentin Garrido, Nicolas Ballas, Mahmoud Assran, Adrien Bardes, Laurent Najman, Michael Rabbat, Emmanuel Dupoux, and Yann LeCun. Intuitive physics understanding emerges from self-supervised pretraining on natural videos, 2025. 2
- [13] Jing Gu, Xian Liu, Yu Zeng, Ashwin Nagarajan, Fangrui Zhu, Daniel Hong, Yue Fan, Qianqi Yan, Kaiwen Zhou, Ming-Yu Liu, and Xin Eric Wang. ”phyworldbench”: A comprehensive evaluation of physical realism in text-to-video models, 2025. 2
- [14] Xuyang Guo, Jiayan Huo, Zhenmei Shi, Zhao Song, Jiahao Zhang, and Jiale Zhao. T2vphysbench: A first-principles benchmark for physical consistency in text-to-video generation, 2025. 2
- [15] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in neural information processing systems*, 35:8633–8646, 2022. 2
- [16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. 4, 6, 7, 8
- [17] Bingyi Kang, Yang Yue, Rui Lu, Zhijie Lin, Yang Zhao, Kaixin Wang, Gao Huang, and Jiashi Feng. How far is video generation from world model? – a physical law perspective. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025. 2, 3
- [18] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 4
- [19] Chenyu Li, Oscar Michel, Xichen Pan, Sainan Liu, Mike Roberts, and Saining Xie. Pisa experiments: Exploring physics post-training for video diffusion models by watching stuff drop. In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025. 2, 3, 5, 7, 8
- [20] Dacheng Li, Yunhao Fang, Yukang Chen, Shuo Yang, Shiyi Cao, Justin Wong, Michael Luo, Xiaolong Wang, Hongxu Yin, Joseph E. Gonzalez, Ion Stoica, Song Han, and Yao Lu. Worldmodelbench: Judging video generation models as world models, 2025. 2
- [21] Shaowei Liu, Zhongzheng Ren, Saurabh Gupta, and Shenlong Wang. Physgen: Rigid-body physics-grounded image-to-video generation. In *European Conference on Computer Vision (ECCV)*, 2024. 3
- [22] Saman Motamed, Minghao Chen, Luc Van Gool, and Iro Laina. Travl: A recipe for making video-language models better judges of physics implausibility, 2025. 2
- [23] Saman Motamed, Laura Culp, Kevin Swersky, Priyank Jaini, and Robert Geirhos. Do generative video models understand physical principles? *arXiv preprint arXiv:2501.09038*, 2025. 2
- [24] Isaac Newton, I Bernard Cohen, and Anne Whitman. *The Principia: mathematical principles of natural philosophy*. Univ of California Press, 1999. 1
- [25] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 4, 7
- [26] Google Research and DeepMind. Veo 3: Text-to-video model. Online model release, 2025. 4, 6

- 589 [27] Enes Sanli, Baris Sarper Tezcan, Aykut Erdem, and Erkut
590 Erdem. Can your model separate yolks with a water bot-
591 tle? benchmarking physical commonsense understanding in
592 video generation models, 2025. 2
- 593 [28] Ayush Sarkar, Hanlin Mai, Amitabh Mahapatra, Svetlana
594 Lazebnik, David A Forsyth, and Anand Bhattad. Shad-
595 ows don't lie and lines can't bend! generative models don't
596 know projective geometry... for now. In *Proceedings of
597 the IEEE/CVF conference on computer vision and pattern
598 recognition*, pages 28140–28149, 2024. 3
- 599 [29] ModelScope Team. Diffsynth-studio: A diffusion model
600 engine. [https://github.com/modelscope/
601 DiffSynth-Studio](https://github.com/modelscope/DiffSynth-Studio), 2025. GitHub repository, Apache-
602 2.0 license. 7
- 603 [30] NVIDIA Cosmos Team. Cosmos-predict2: World founda-
604 tion models for physical ai. [https://github.com/
605 nvidia-cosmos/cosmos-predict2](https://github.com/nvidia-cosmos/cosmos-predict2), 2025. GitHub
606 repository, Apache-2.0 license. 4, 6
- 607 [31] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao,
608 Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianx-
609 iao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jin-
610 gren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao,
611 Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang,
612 Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei
613 Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui,
614 Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang,
615 Wenmeng Zhou, Wenten Wang, Wenting Shen, Wenyuan Yu,
616 Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu
617 Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yi-
618 tong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun
619 Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi
620 Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open
621 and advanced large-scale video generative models. *arXiv
622 preprint arXiv:2503.20314*, 2025. 4, 6, 7, 8
- 623 [32] Angtian Wang, Haibin Huang, Zhiyuan Fang, Yiding Yang,
624 and Chongyang Ma. ATI: Any trajectory instruction for con-
625 trollable video generation. *arXiv preprint arXiv:2505.22944*,
626 2025. 7, 8
- 627 [33] Chen Wang, Chuhao Chen, Yiming Huang, Zhiyang Dou,
628 Yuan Liu, Jiatao Gu, and Lingjie Liu. Physctrl: Generative
629 physics for controllable and physics-grounded video genera-
630 tion. *arXiv preprint arXiv:2509.20358*, 2025. 3
- 631 [34] Thaddäus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane
632 Gu, Nick Matarese, Kevin Swersky, Been Kim, Priyank
633 Jaini, and Robert Geirhos. Video models are zero-shot learn-
634 ers and reasoners. *arXiv preprint arXiv:2509.20328*, 2025.
635 2
- 636 [35] Xiaoyan Xing, Konrad Groh, Sezer Karaoglu, Theo Gevers,
637 and Anand Bhattad. Luminet: Latent intrinsics meets dif-
638 fusion models for indoor scene relighting. In *Proceedings
639 of the Computer Vision and Pattern Recognition Conference*,
640 pages 442–452, 2025. 3
- 641 [36] Yu Yuan, Xijun Wang, Tharindu Wickremasinghe, Zeeshan
642 Nadir, Bole Ma, and Stanley Chan. Newtongen: Physics-
643 consistent and controllable text-to-video generation via neu-
644 ral newtonian dynamics. *arXiv preprint arXiv: 2509.21309*,
645 2025. 3
- [37] Guanqi Zhan, Chuanxia Zheng, Weidi Xie, and Andrew Zis-
serman. A general protocol to probe large vision models for
3d physical understanding. *Advances in Neural Information
Processing Systems*, 37:43468–43498, 2024. 2, 3
- [38] Chenyu Zhang, Daniil Cherniavskii, Antonios Tragoudaras,
Antonios Vozikis, Thijmen Nijdam, Derck W. E. Prinzhorn,
Mark Bodracska, Nicu Sebe, Andrii Zadaianchuk, and Efs-
tratis Gavves. Morpheus: Benchmarking physical reason-
ing of video generative models with real physical experi-
ments, 2025. 2, 5
- [39] Xiangdong Zhang, Jiaqi Liao, Shaofeng Zhang, Fanqing
Meng, Xiangpeng Wan, Junchi Yan, and Yu Cheng. Vide-
orepa: Learning physics for video generation through re-
lational alignment with foundation models. *arXiv preprint
arXiv:2505.23656*, 2025. 3
- [40] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G
Lowe. Unsupervised learning of depth and ego-motion from
video. In *Proceedings of the IEEE conference on computer
vision and pattern recognition*, pages 1851–1858, 2017. 2