
PhysEditWorld: A Large-Scale Dataset Toward Physics-Editable World Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent game world models can synthesize visually plausible, action-conditioned
2 rollouts. However, their interaction behaviors often remain limited to exploratory
3 or wandering trajectories, and physical dynamics are typically learned as implicit
4 correlations from data rather than as controllable variables. This limitation hinders
5 their applicability to authored game environments, where physical rules are deliber-
6 ately designed and require explicit manipulation. We introduce PhysEditWorld,
7 a multimodal dataset with physical parameters, with a primary focus on gravity
8 in this initial version. At its core, PhysEditWorld is built upon a replay paradigm
9 implemented with a UE5 replay-and-rendering pipeline. Each scenario records
10 a normalized action trace and replays the same initial state, character controller,
11 action sequence, and camera policy under multiple gravity configurations, enabling
12 controlled and attributable physical variation. PhysEditWorld contains 12 cinematic
13 UE5 scenes, over 100 hours of gameplay interactions, and more than 60 million
14 rendered rollout frames. Each sample provides synchronized multimodal signals,
15 including RGB, depth, normals, audio, action traces, camera trajectory, engine
16 states, semantic annotations, and explicit gravity labels. We further conduct initial
17 utility studies on both generative video models and world understanding models,
18 demonstrating that PhysEditWorld enables improved gravity-faithful dynamics
19 modeling, enhances consistency under physical edits, and provides a scalable
20 foundation for controllable world modeling research.

21 1 Introduction

22 Recent game world models have progressed from visual predictors to interactive generative simulators.
23 Systems such as Genie, DIAMOND, GameNGen, GameGen-X, YUMI, LingBot-World, and Matrix-
24 Game-3.0 show that large generative models can synthesize plausible gameplay trajectories, support
25 user interaction, or maintain longer-horizon world consistency [1, 2, 3, 4, 5, 6, 7]. However, these
26 models typically learn physics as an implicit regularity of the data distribution. They can imitate
27 how a game usually evolves, but they are not designed to answer a question that is central to game
28 authoring: how should the same scene evolve if a physical rule is edited?

29 Unlike natural world modeling, where physics is usually treated as a latent constant, game world
30 modeling must eventually support physical laws as editable design variables. Developers routinely
31 tune gravity scale, jump behavior, friction, drag, and wind to shape pacing, difficulty, player feel,
32 and emergent motion. A learned world model that entangles these parameters with appearance and
33 gameplay statistics may generate visually plausible clips, yet still fail as an editable game simulator.
34 In this paper, we focus on gravity as a first measurable step toward editable physics, because it is
35 widely supported across engines and produces observable changes in jump arcs, airtime, fall speed,
36 and object trajectories.

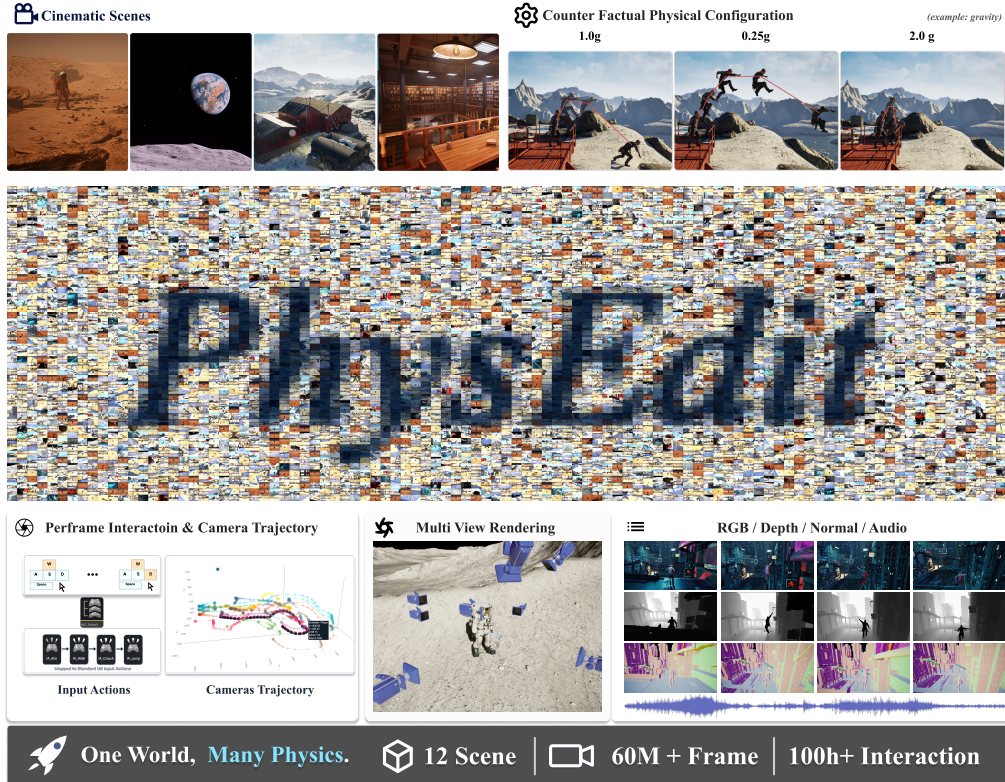


Figure 1: Overview of PhysEditWorld. The dataset contains 12 cinematic UE5 scenes, 100+ hours of gameplay interactions, and 60M+ rendered rollout frames. Recorded interactions are replayed under editable physical configurations and rendered from 8 synchronized camera views, with RGB, depth, normal, audio, action, camera, engine-state, and gravity annotations.

37 Existing datasets and benchmarks do not directly measure this capability. Game and reinforcement-
 38 learning environments such as ALE, Procgen, MineRL, and CARLA provide interactive worlds or
 39 human demonstrations, but they are primarily designed for policy learning, navigation, or gener-
 40 alization under fixed rules [8, 9, 10, 11]. World-exploration datasets such as Sekai provide large-
 41 scale first-person and drone-view videos with rich annotations and camera trajectories, but are not
 42 organized around matched physical interventions [12]. Physics reasoning and video-generation
 43 benchmarks such as PHYRE, CLEVRER, Physion, VBench, VideoPhy, PhyGenBench, Physics-
 44 IQ, and PhysInOne evaluate intuitive physics, physical plausibility, or physics-aware generation
 45 [13, 14, 15, 16, 17, 18, 19, 20]. These resources are valuable, but they do not provide matched
 46 gameplay clips in which the scene, interaction trace, and camera policy are held fixed while gravity is
 47 explicitly changed.

48 We introduce **PhysEditWorld**, a multimodal dataset for physics-editable game world modeling.
 49 PhysEditWorld contains 12 cinematic UE5 scenes and over 60M rendered frames, with matched
 50 rollouts designed for gravity-conditioned generation and evaluation. The dataset is built around
 51 *comparability*: each replay group fixes the authored scene, initial state, action trace, character
 52 controller, and camera policy, then reruns the scenario under different gravity configurations. The
 53 pipeline records synchronized first-person RGB, third-person RGB, depth, normal maps, action traces,
 54 semantic captions, engine states, and gravity annotations. Because non-gravity factors are controlled
 55 within a replay group, the resulting motion differences can be evaluated as responses to the physical
 56 intervention.

57 PhysEditWorld enables evaluation beyond visual plausibility by testing whether generated rollouts
 58 are faithful to edited gravity. We study this capability in gravity-conditioned generation, first-person
 59 world-model rollouts, and video-language gravity inference. Across representative backbones, we
 60 find that current models can maintain visual realism, but often under-express gravity-sensitive motion

61 or confuse the relative ordering of gravity levels. This suggests that editable physics remains a
62 missing capability in current game world models.

63 **Our contributions are as follows.**

- 64 • We introduce **PhysEditWorld**, a large-scale multimodal dataset for physics-editable game world
65 modeling. To the best of our knowledge, it is the first dataset organized around matched gameplay
66 rollouts with explicit gravity interventions.
- 67 • We develop a UE5 **replay-and-rendering pipeline** that automatically replays the same scene,
68 initial state, character controller, action trace, and camera policy under controlled physics configu-
69 rations.
- 70 • We conduct dataset utility studies showing that PhysEditWorld can be used to evaluate and improve
71 gravity awareness in generative video models, game world models, and video-language models.

72 **2 Related work**

73 **Game world models.** Game world models have evolved from single-game neural simulators to
74 open-domain interactive video generators. Early works such as GameGAN, Playable Video Genera-
75 tion, Playable Environments, and Promptable Game Models studied controllable neural simulation
76 and interactive video manipulation from gameplay or video data [21, 22, 23, 24]. Recent sys-
77 tems, including Genie, Oasis, DIAMOND, GameNGen, GameGen-X, GameFactory, MineWorld,
78 YUME, LingBot-World, Matrix-Game 3.0, and related interactive video world models, scale action-
79 conditioned generation, open-world exploration, long-horizon consistency, and real-time interac-
80 tion [1, 25, 2, 3, 4, 26, 27, 7, 5, 6, 28]. While these models learn rich visual and interaction dynamics,
81 physical rules are usually absorbed as implicit dataset regularities rather than exposed as editable
82 variables. PhysEditWorld addresses this gap with explicit gravity annotations and matched replays
83 under fixed scene, action, controller, and camera conditions.

84 **Video and world-model datasets.** Existing game and world-model datasets mainly emphasize
85 scale, action control, exploration coverage, or state supervision. MineRL, VPT, and MineDojo
86 provide Minecraft demonstrations, video pretraining data, simulator environments, or task suites for
87 embodied agent learning [10, 29, 30]. OGameData/GameGen-X and GF-Minecraft/GameFactory
88 target open-world or action-controllable game video generation [4, 26]; Sekai and YUME focus on
89 world exploration and interactive generation [12, 7]; and WildWorld and MultiWorld add explicit
90 state, multi-agent, or multi-view supervision [31, 32]. Related resources also study high-DoF action-
91 to-video control or physics-related gameplay failures [33, 34]. However, these datasets are not
92 organized around matched physical interventions. PhysEditWorld instead replays the same authored
93 scene, interaction trace, character controller, and camera policy under different gravity configurations,
94 making the physical edit directly comparable.

95 **Physics benchmarks and controllable generation.** Many benchmarks study physical under-
96 standing from visual data. ShapeStacks and ADEPT probe object stability and expectation vi-
97 olation, while PHYRE, I-PHYRE, IntPhys, IntPhys2, CLEVRER, CATER, and Physion eval-
98 uate intuitive physics, compositional actions, causal reasoning, intervention, and future predic-
99 tion [35, 36, 13, 37, 38, 39, 14, 40, 15]. CoPhy, ComPhy, CRIPP-VQA, ContPhy, CausalVQA, and
100 QuantiPhy further study counterfactual dynamics, hidden physical properties, causal alternatives,
101 and quantitative physical quantities [41, 42, 43, 44, 45, 46]. Video-generation benchmarks such as
102 VBench, VBench-2.0, VideoPhy, PhyGenBench, Physics-IQ, PhysInOne, WorldScore, and Newton-
103 Rewards evaluate physical plausibility, intrinsic faithfulness, world-generation quality, or Newtonian
104 motion consistency [16, 47, 17, 18, 19, 20, 48, 49]. Recent controllable generation methods condition
105 on force, torque, force fields, physical parameters, Newtonian dynamics, or learned physical pri-
106 ors [50, 51, 52, 53, 54, 55, 56]. These works provide important reasoning, evaluation, and generation
107 tools, but they do not provide interactive game rollouts where one physical rule is edited while scene
108 content and interaction are held fixed.

109 **Synthetic simulation platforms.** Simulation platforms enable controllable data generation and
110 ground-truth annotation. ALE, Procgen, MineRL, and CARLA support standardized environments
111 for games, reinforcement learning, and driving [8, 9, 10, 11]. AI2-THOR, RoboTHOR, ProcTHOR,
112 Habitat, Gibson, iGibson, TDW, Kubric, VirtualHome, and BEHAVIOR-1K provide controllable

113 3D simulation with rich sensor outputs, physical interaction, programmatic actions, or scalable
 114 embodied-AI environments [57, 58, 59, 60, 61, 62, 63, 64, 65, 66]. Robotics simulators such as
 115 SAPIEN, RL Bench, ManiSkill2, and Isaac Gym support articulated-object interaction, manipulation,
 116 and high-throughput physics simulation [67, 68, 69, 70]. UnrealCV connects Unreal Engine to
 117 computer-vision pipelines, while UnrealZoo scales UE environments for embodied AI [71, 72].
 118 PhysEditWorld builds on UE but targets a different protocol: matched replay under edited gravity.

119 3 Physics Edit Dataset

120 PhysEditWorld is a large-scale multimodal dataset for gravity-editable game world modeling. It
 121 contains 12 cinematic-quality UE5 scenes, more than 100 hours of human gameplay interactions,
 122 and over 60 million rendered rollout frames across multiple gravity configurations and synchronized
 123 camera views. All videos are rendered at 30 FPS and 1280×720 resolution. Unlike many physics-
 124 video datasets that focus on short, isolated physical events, PhysEditWorld is built from interactive
 125 game scenarios, enabling gravity effects to be studied in action-conditioned world-model rollouts.

126 Each rollout provides synchronized RGB video, depth maps, surface normals, audio when available,
 127 semantic captions, action traces, camera parameters, engine-state logs, and explicit gravity labels.
 128 The engine logs include camera trajectories, character states, object states, and relevant physical
 129 variables exported in UE5’s native world coordinate system and physical scale. The basic unit is a
 130 *matched replay group*: the scene, initial state, character controller, action trace, and camera policy are
 131 fixed, while only the gravity configuration changes. This structure makes gravity-dependent effects
 132 such as jump height, airtime, fall speed, landing timing, camera displacement, and object motion
 133 directly comparable across variants.

134 Table 1 summarizes the advantage of PhysEditWorld over prior game, world-model, and physics-
 135 oriented datasets. Physics-video and physical-reasoning benchmarks such as CLEVRER, Physion,
 136 VBench, VideoPhy, PhyGenBench, Physics-IQ, and PhysInOne provide useful tests of physical
 137 plausibility, intuitive physics, or physics-aware generation, but generally lack interactive action traces
 138 and matched physical interventions [14, 15, 16, 17, 18, 19, 20]. Game and world-modeling datasets
 139 such as ALE, Procgen, MineRL, Sekai, GameFactory/GF-Minecraft, and OGameData/GameGen-X
 140 provide action control, visual scale, or camera information to varying degrees, but typically operate
 141 under fixed physical rules [8, 9, 10, 12, 26, 4]. In contrast, PhysEditWorld combines multi-view
 142 rendering, RGB-depth-normal-audio supervision, explicit editable gravity, camera annotations, and
 143 action/control traces, making it possible to evaluate whether a model changes world dynamics
 144 consistently under the same scene, action sequence, and camera policy.

Table 1: Comparison with related datasets and data-generation resources. RGB-D-N-A denotes synchronized RGB, depth, normal, and audio. Edit phys. denotes whether physical attributes are explicitly editable; PhysEditWorld focuses on gravity in the current release. Camera ann. denotes camera parameters or trajectories. Action/control denotes recorded actions, controls, or interaction traces. ✓: yes, ✗: no, △: partial.

Resource	Multi-view	RGB-D-N-A	Edit phys.	Camera ann.	Action
ALE [8]	✗	✗	✗	✗	✓
Procgen [9]	✗	✗	✗	✗	✓
MineRL [10]	✗	✗	✗	△	✓
Sekai [12]	△	△	✗	✓	✗
GameFactory / GF-Minecraft [26]	✗	✗	✗	✗	✓
OGameData / GameGen-X [4]	△	✗	✗	✓	△
CLEVRER [14]	✗	✗	✗	✗	✗
Physion [15]	✗	△	✗	△	✗
VBench [16]	✗	✗	✗	✗	✗
PhyGenBench [18]	✗	✗	✗	✗	✗
VideoPhy [17]	✗	✗	✗	✗	✗
Physics-IQ [19]	✗	✗	✗	✗	✗
PhysInOne [20]	✓	△	✗	△	✗
PhysEditWorld	✓	✓	✓	✓	✓

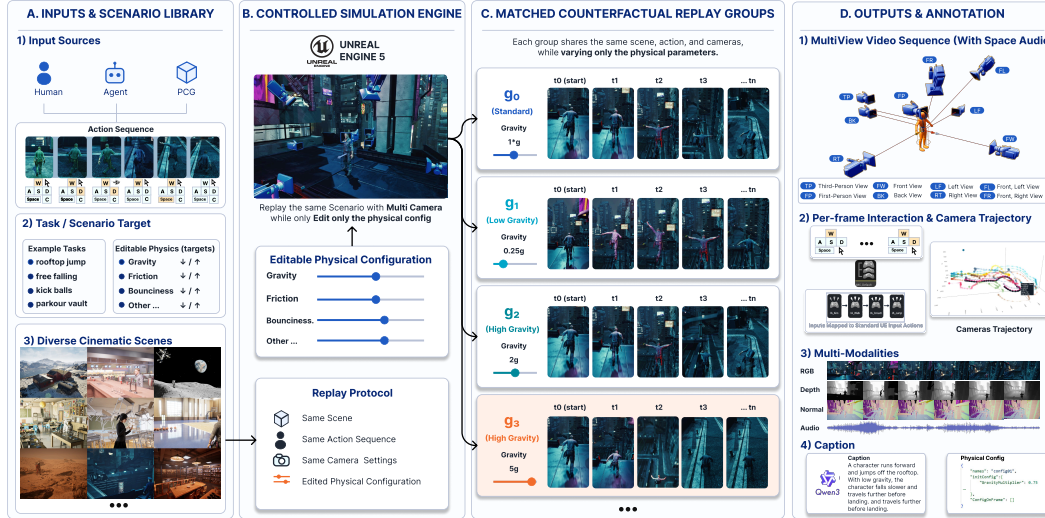


Figure 2: Overview of the PhysEditWorld data pipeline. A) Input sources and scenario library provide human-controlled, agent-generated, or scripted action sequences in diverse cinematic UE5 scenes. B) A controlled UE5 simulation engine replays the same scenario while editing only the physical configuration. C) Matched replay groups share the same scene, action sequence, character controller, and camera policy while varying gravity. D) The pipeline exports synchronized multi-view videos, spatial audio, per-frame interaction and camera trajectories, multimodal render passes, captions, and physical annotations.

145 4 Physics Data Pipeline

146 PhysEditWorld is generated by a UE5 replay-and-rendering pipeline integrated into the standard
 147 game-development workflow. Rather than reconstructing scenes in a separate simulator, the pipeline
 148 operates directly on artist-ready UE5 levels through an in-editor plug-in, converting authored game
 149 content into data-production-ready scenarios with registered scenes, validated assets, controller
 150 bindings, camera policies, and replay settings. It then builds normalized action sequences from
 151 human-controlled, agent-generated, or scripted inputs and replays them through the same UE5
 152 gameplay stack under controlled physical configurations. The final stage performs synchronized
 153 rendering, runtime logging, data organization, and post-hoc VLM annotation, producing aligned
 154 rollouts with semantic descriptions.

155 4.1 Scenario and input construction

156 The pipeline starts from artist-ready UE5 levels and converts them into replayable scenarios through
 157 the in-editor plug-in. Each scenario registers the authored level, interactable assets, character setup,
 158 collision configuration, camera anchors, and replay settings required for controlled simulation.
 159 Scenario preparation emphasizes physical visibility and replay stability: artists select events such as
 160 jumps, falls, object interactions, and height-varying traversal, and validate that the corresponding
 161 assets, controllers, spawn points, and camera anchors support reproducible replay.

162 For each prepared scenario, the pipeline constructs action sequences from human-controlled, agent-
 163 generated, or scripted inputs. Instead of recording raw device events, it records semantic Input
 164 Action sequences from UE5’s Enhanced Input System, including movement axes, jump commands,
 165 camera deltas, and key or button states. This representation preserves the intent of the interaction
 166 while avoiding hardware-specific confounders such as keyboard layout, mouse sensitivity, controller
 167 mapping, or platform-dependent event APIs.

168 4.2 Controlled UE5 simulation

169 To scale data production within standard game-development workflows, PhysEditWorld uses an
170 in-editor DataFactory plug-in rather than a separate simulator. The plug-in operates on native
171 UE5 abstractions, including authored levels, character controllers, Enhanced Input actions, camera
172 components, collision volumes, and gameplay logic. It provides unified tools for scene registration,
173 replay specification, controller binding, camera setup, physical-parameter editing, validation, and
174 batch execution.

175 During simulation, the normalized action sequence is injected back into the same UE5 gameplay
176 stack used during capture. Physical configurations are applied as explicit simulation parameters,
177 while the authored scene, controller logic, input sequence, and camera policy remain fixed. This
178 design separates interaction capture from physical intervention and allows controlled physical edits
179 without recollecting behavior or modifying the authored level.

180 4.3 Matched replay expansion

181 Let A denote the prepared UE5 scenario, S the normalized action sequence, M the character
182 controller, θ the editable physical configuration, and π_c the camera policy. A rollout is generated as

$$x = F(A, S, M, \theta, \pi_c). \quad (1)$$

183 Replay expansion keeps (A, S, M, π_c) fixed and reruns the same scenario under different values of θ .
184 This construction makes the physical edit the controlled variable, so that differences across variants
185 can be attributed to the edited simulation parameter rather than to changes in scene content, input
186 behavior, character setup, or camera specification.

187 The action sequence is treated as scene-bound because its semantics depend on local geometry,
188 obstacles, and interaction targets. The pipeline therefore expands validated scenario-action pairs
189 over compatible physical configurations and character setups, producing matched variants while
190 preserving the original authored environment and interaction intent.

191 4.4 Synchronized export and annotation

192 Each replay is exported through Movie Render Queue and runtime logging. Movie Render Queue
193 produces rendered observation streams and auxiliary render passes, while the runtime logger records
194 time-aligned interaction, camera, state, and physical-configuration metadata. All exported records are
195 indexed by a shared rollout identity and frame timeline, enabling deterministic alignment between
196 rendered observations and engine-side logs.

197 After rendering, the pipeline performs data organization, post-hoc semantic annotation, and quality
198 filtering. We annotate each rollout with a per-clip caption generated by Qwen3-VL-8B-Instruct [73];
199 the captioning model observes sampled rendered frames only and does not receive simulator metadata.
200 Replays with rendering failures, broken synchronization, severe camera clipping, unstable simulation,
201 or divergence not attributable to the intended physical edit are discarded.

202 5 Dataset Utility for Gravity-Conditioned Generation

203 We evaluate PhysEditWorld as supervision for gravity-conditioned generation. Our goal is not to
204 rank generative models comprehensively, but to test whether fine-tuning on matched gravity rollouts
205 improves a model’s response to explicit gravity conditions. This distinction matters because visually
206 plausible generation can still fail physically: the camera may remain nearly static, falling may be
207 delayed, or high-gravity rollouts may not accelerate more than low-gravity ones. PhysEditWorld
208 makes these failures measurable by holding the scene, action sequence, camera policy, and initial
209 state fixed while varying gravity.

210 5.1 Experimental Setup

211 **Data splits.** We construct training and evaluation splits at the *replay-group* level rather than the clip
212 level, so that no scene appears in both splits. The training set contains 1,530 (scene, action, gravity)
213 tuples drawn from 9 scenes; the held-out set contains 170 tuples from the remaining 3 scenes. Within

214 each replay group, all gravity variants are kept together to preserve the matched-comparison structure.
 215 We sample at most one clip per (scene, action, gravity) triplet to prevent near-duplicate frames from
 216 dominating the loss. Gravity multipliers are drawn from $\{0.05, 0.1, 0.5, 1.0, 2.0, 5.0, 20.0\} \times g_{\oplus}$.

217 **Baselines.** We evaluate two settings. For gravity-conditioned video generation, we compare zero-shot
 218 against **Wan2.2-TI2V-5B** [74] SFT version on our data. Kling 3.0 and Seedance 2.0 appear only in
 219 qualitative comparisons as they do not release training code, while Wan2.2-TI2V-5B is additionally
 220 fine-tuned via LoRA [75] on PhysEditWorld. For action-conditioned first-person world modeling,
 221 we do quality study on **LingBot-World** [5] fine-tuned on our data and compare against the frozen
 222 **Matrix-Game 3.0** [6] and **LingBot-World** [5].

223 **Training protocol.** We apply LoRA with rank 128 on all attention projection matrices of both
 224 backbones, trained for 5 epochs with AdamW (learning rate 1×10^{-4} , batch size 8) on $8 \times H100$
 225 GPUs. Inputs are 5-second clips at 30 FPS, 1280×720 , with the gravity multiplier injected as a text
 226 token in the conditioning prompt (e.g., "gravity: 0.25g"). All other hyperparameters follow
 227 each backbone’s released SFT recipe.

228 **Evaluation metric.** Since generated videos do not expose simulator states or ground-truth trajectories,
 229 we recover a per-frame camera trajectory with VGGT [76] and use its vertical-axis component as a
 230 one-dimensional fall-progress signal $q(t)$. We restrict evaluation to the pre-landing segment, detected
 231 via a simple progress threshold on smoothed $q(t)$, and discard clips with insufficient fall progress
 232 or unreliable landing detection. Because VGGT trajectories carry scale ambiguity, all speeds and
 233 accelerations are reported in normalized VGGT units, and only relative comparisons within a matched
 234 replay group are meaningful. Implementation details (smoothing window, threshold values, rejection
 235 rates) are provided.

236 For each reliable clip, we compute fall-axis speed $v(t)$ as the central-difference derivative of the
 237 smoothed $q(t)$ and fit a linear model:

$$v(t) = at + b. \quad (2)$$

238 The fitted slope a serves as a normalized acceleration proxy, and R^2 measures how well the speed
 239 curve follows a linear acceleration pattern. To test whether motion ordering follows the requested
 240 gravity ordering, we compute pairwise gravity-acceleration alignment within each matched replay
 241 group:

$$\text{Align} = \frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} [(g_i - g_j)(a_i - a_j) > 0], \quad (3)$$

242 where $\mathcal{P} = \{(i, j) \mid i < j, i, j \in \text{group}\}$ is the set of unordered pairs within a group, g_i is the
 243 requested gravity, and a_i is the fitted acceleration proxy. A perfectly gravity-faithful model achieves
 244 $\text{Align} = 1.0$; chance is 0.5.

245 5.2 Gravity-Conditioned Video Generation / WorldModel

246 Table 2 reports a representative matched case for Wan2.2-TI2V-5B before and after PhysEditWorld
 247 supervised fine-tuning. The zero-shot model is insensitive to the requested gravity: acceleration
 248 proxies are near-zero across all three settings, and the alignment of 33.3% confirms that gravity
 249 ordering is not preserved. After PhysEditWorld SFT, the acceleration proxy increases monotonically
 250 with requested gravity, alignment reaches 100%, and mean R^2 rises from 0.066 to 0.570, indicating
 251 that the model now responds to gravity as a controllable variable. Figure 3 top shows a qualitative
 252 example under a $1g$ prompt: the zero-shot model produces minimal camera motion while the scene
 253 remains nearly static, whereas the LoRA-tuned model generates strong forward self-motion toward
 254 the lunar surface with visible motion blur.

Table 2: VGGT-based gravity response analysis on a representative held-out matched replay group with three gravity settings (0.25 \times , 0.75 \times , 10.0 \times). Accel. proxy denotes the fitted slope a in $v(t) = at + b$ on the pre-landing segment, in normalized VGGT units; only *within-row*, *within-group* comparisons are meaningful. Mean R^2 averages the linear-fit quality across the three settings.

Model	Accel. @ 0.25 \times	Accel. @ 0.75 \times	Accel. @ 10.0 \times	Mean $R^2 \uparrow$	Gravity Align. \uparrow
Wan2.2-TI2V-5B (zero-shot)	0.0002	0.0046	0.0009	0.066	33.3%
Wan2.2-TI2V-5B + PhysEditWorld SFT	0.1766	0.3983	0.6214	0.570	100.0%

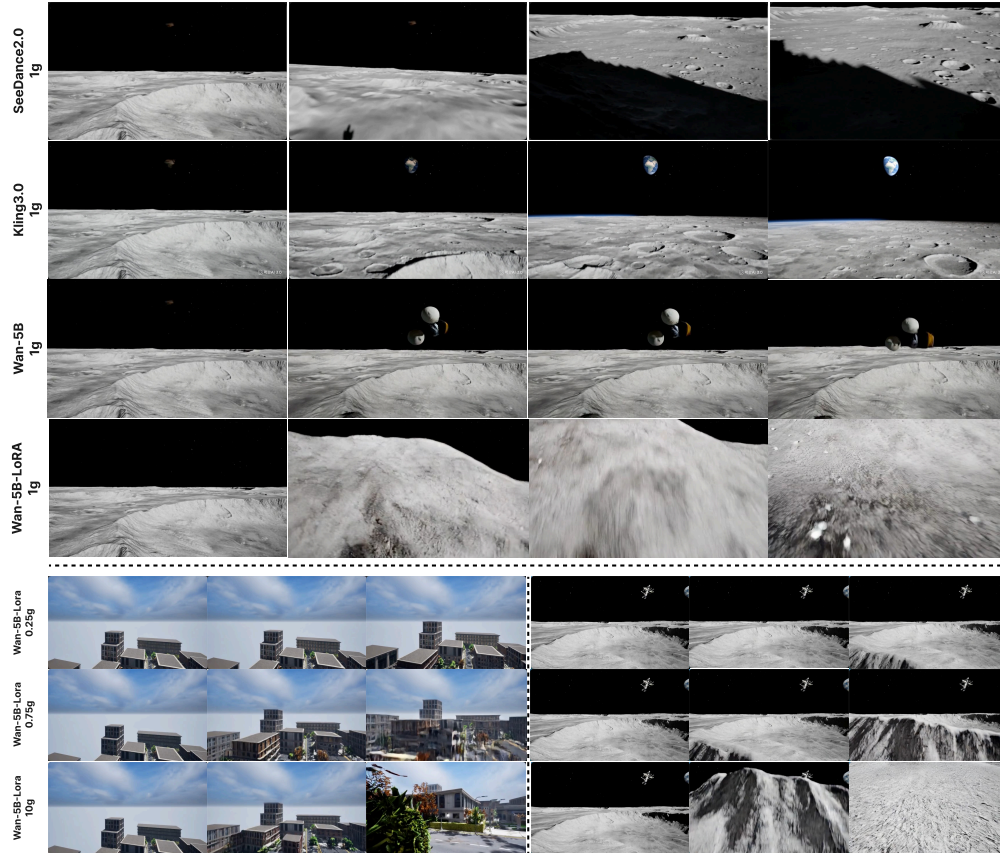


Figure 3: Qualitative results under a free-fall prompt. *Top*: SeeDance2.0, Kling3.0, Wan2.2-TI2V-5B, and +PhysEditWorld LoRA at $1g$; baselines remain near-static while the LoRA model generates strong self-motion with motion blur. *Bottom*: +PhysEditWorld LoRA at $0.25\times/0.75\times/10\times$ gravity; descent speed scales monotonically, confirming continuously controllable gravity.



Figure 4: First-person world-model case study. First frame is generated by GPT-Image2. Baseline models (Matrix-Game 3.0, zero-shot LingBotWorld) remain near the platform edge and never enter free-fall under a $1g$ prompt. After PhysEditWorld LoRA-128 tuning, LingBotWorld generates platform departure and gravity-dependent downward self-motion.

255 5.3 Action-Conditioned First-Person World Modeling

256 We evaluate whether PhysEditWorld supervision transfers to action-conditioned first-person world
257 models. We use a simple stress test: the W key is held continuously from a rooftop edge, where a
258 gravity-faithful model should produce platform departure followed by gravity-dependent downward
259 motion. As shown in Figure 4, both Matrix-Game 3.0 and zero-shot LingBotWorld remain near the
260 platform edge and never enter free-fall. After PhysEditWorld LoRA tuning, LingBotWorld generates
261 platform departure followed by downward self-motion that scales with the requested gravity. This
262 demonstrates that the gravity-response failure is not limited to text-conditioned generation, and that
263 PhysEditWorld provides effective supervision across both generation paradigms.

264 6 Dataset Utility for Gravity-Aware VLMs

265 **Experimental Setup.** We evaluate Qwen3-VL-8B-Instruct [73] on a held-out set of 170 rollouts
266 drawn from the same stratified split described in §5.1. The task requires the model to predict both the
267 gravity class (low / normal / high) and the continuous gravity multiplier from a short gameplay video
268 clip.

269 **Baseline.** We first assess the zero-shot capability of Qwen3-VL-8B-Instruct [73] without any domain-
270 specific fine-tuning, serving as a strong off-the-shelf VLM baseline.

271 **Training Protocol.** We apply LoRA SFT on the remaining 1,530 training rollouts. To discourage
272 label memorization, gravity targets are jittered by $\pm 10\%$ during training; evaluation is conducted
273 against the original held-out labels.

274 **Metrics.** We report class accuracy for the three-way gravity classification, and mean absolute error
275 (MAE), median absolute percentage error (Median APE), and within-10% rate for the continuous
276 gravity multiplier prediction.

Table 3: Gravity-aware prediction on held-out rollouts. +SFT denotes PhysEditWorld LoRA tuning.

Model	Class Acc. (%) \uparrow	MAE \downarrow	Median APE (%) \downarrow	Within 10% \uparrow
Qwen3-VL-8B[73]	24.71	5.4573	95.00	9.48
Qwen3-VL-8B + SFT	95.29	0.8305	6.22	90.59

277 As shown in Table 3, the zero-shot model achieves a class accuracy of 24.71%—below the 33.3%
278 random baseline for a three-way classification task—with a median APE of 95.00% and a within-10%
279 rate of only 9.48%, confirming that gravity magnitude is not reliably recoverable from video without
280 targeted supervision. After PhysEditWorld LoRA SFT, class accuracy improves to 95.29% and
281 median APE drops to 6.22%, with 90.59% of predictions falling within 10% of the true gravity
282 multiplier. These results demonstrate that the physics-varied rollouts in PhysEditWorld provide
283 effective supervision signal for fine-grained, gravity-aware video-language understanding.

284 7 Conclusion

285 PhysEditWorld introduces editable game-world modeling as a controlled physical-intervention prob-
286 lem. Instead of asking only whether a generated rollout looks plausible, it asks whether the same
287 authored scenario evolves consistently when a physical rule is changed. By replaying matched inter-
288 actions under explicit gravity configurations, PhysEditWorld separates visual realism from physical
289 controllability and reveals failure modes that standard video-quality metrics can overlook. The current
290 release focuses on gravity as a measurable and widely supported first step, covering effects such as
291 airtime, fall speed, jump arcs, and landing dynamics. In future releases, we plan to extend the dataset
292 and pipeline to additional editable physical attributes, such as friction, drag, restitution, wind, and
293 object-level physical parameters, further supporting controllable neural game engines that can be
294 edited as authored simulations rather than only imitated as videos.

295 References

- 296 [1] Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes,
297 Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, Yusuf Aytar, Sarah Bechtle,
298 Feryal Behbahani, Stephanie C. Y. Chan, Nicolas Heess, Lucy Gonzalez, Simon Osindero,
299 Sherjil Ozair, Scott Reed, Jingwei Zhang, Konrad Zolna, Jeff Clune, Nando de Freitas, Satinder
300 Singh, and Tim Rocktäschel. Genie: Generative interactive environments. In *Proceedings of the*
301 *41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine*
302 *Learning Research*, pages 4603–4623. PMLR, 2024.
- 303 [2] Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kanervisto, Amos Storkey, Tim Pearce, and
304 Francois Fleuret. Diffusion for world modeling: Visual details matter in Atari. In *Advances in*
305 *Neural Information Processing Systems*, 2024.
- 306 [3] Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-
307 time game engines. In *International Conference on Learning Representations*, 2025. URL
308 <https://arxiv.org/abs/2408.14837>.
- 309 [4] Haoxuan Che, Xuanhua He, Quande Liu, Cheng Jin, and Hao Chen. Gamegen-x: Interactive
310 open-world game video generation. In *International Conference on Learning Representations*,
311 2025.
- 312 [5] Robbyant Team, Zelin Gao, Qiuyu Wang, Yanhong Zeng, Jiapeng Zhu, Ka Leong Cheng,
313 Yixuan Li, Hanlin Wang, Yinghao Xu, Shuailei Ma, Yihang Chen, Jie Liu, Yansong Cheng, Yao
314 Yao, Jiayi Zhu, Yihao Meng, Kecheng Zheng, Qingyan Bai, Jingye Chen, Zehong Shen, Yue
315 Yu, Xing Zhu, Yujun Shen, and Hao Ouyang. Advancing open-source world models. *arXiv*
316 *preprint arXiv:2601.20540*, 2026.
- 317 [6] Zile Wang, Zexiang Liu, Jiaying Li, Kaichen Huang, Baixin Xu, Fei Kang, Mengyin An, Peiyu
318 Wang, Biao Jiang, Yichen Wei, Yidan Xietian, Jiangbo Pei, Liang Hu, Boyi Jiang, Hua Xue,
319 Zidong Wang, Haofeng Sun, Wei Li, Wanli Ouyang, Xianglong He, Yang Liu, Yangguang
320 Li, and Yahui Zhou. Matrix-game 3.0: Real-time and streaming interactive world model with
321 long-horizon memory. *arXiv preprint arXiv:2604.08995*, 2026.
- 322 [7] Xiaofeng Mao, Shaoheng Lin, Zhen Li, Chuanhao Li, Wenshuo Peng, Tong He, Jiangmiao
323 Pang, Mingmin Chi, Yu Qiao, and Kaipeng Zhang. Yume: An interactive world generation
324 model. *arXiv preprint arXiv:2507.17744*, 2025.
- 325 [8] Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning
326 environment: An evaluation platform for general agents. *Journal of Artificial Intelligence*
327 *Research*, 47:253–279, 2013.
- 328 [9] Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging procedural
329 generation to benchmark reinforcement learning. In *International Conference on Machine*
330 *Learning*, 2020.
- 331 [10] William H. Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codel, Manuela
332 Veloso, and Ruslan Salakhutdinov. MineRL: A large-scale dataset of Minecraft demonstrations.
333 In *International Joint Conference on Artificial Intelligence*, 2019.
- 334 [11] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun.
335 CARLA: An open urban driving simulator. In *Conference on Robot Learning*, 2017.
- 336 [12] Zhen Li, Chuanhao Li, Xiaofeng Mao, Shaoheng Lin, Ming Li, Shitian Zhao, Zhaopan Xu,
337 Xinyue Li, Yukang Feng, Jianwen Sun, et al. Sekai: A video dataset towards world exploration.
338 In *Advances in Neural Information Processing Systems Datasets and Benchmarks Track*, 2025.
- 339 [13] Anton Bakhtin, Laurens van der Maaten, Justin Johnson, Laura Gustafson, and Ross Gir-
340 shick. PHYRE: A new benchmark for physical reasoning. In *Advances in Neural Information*
341 *Processing Systems*, 2019.
- 342 [14] Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B.
343 Tenenbaum. CLEVRER: Collision events for video representation and reasoning. In *Internat-*
344 *ional Conference on Learning Representations*, 2020.

- 345 [15] Daniel M. Bear, Elias Wang, Damian Mrowca, Felix J. Binder, et al. Physion: Evaluating
346 physical prediction from vision in humans and machines. In *Advances in Neural Information
347 Processing Systems Datasets and Benchmarks Track*, 2021.
- 348 [16] Ziqi Huang, Yinan He, Jiashuo Yu, et al. VBench: Comprehensive benchmark suite for video
349 generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
350 Pattern Recognition*, 2024.
- 351 [17] Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chenfanfu
352 Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. VideoPhy: Evaluating physical
353 commonsense for video generation. In *International Conference on Learning Representations*,
354 2025.
- 355 [18] Fanqing Meng, Jiaqi Liao, Xinyu Tan, Wenqi Shao, Quanfeng Lu, Kaipeng Zhang, Yu Cheng,
356 Dianqi Li, Yu Qiao, and Ping Luo. Towards world simulator: Crafting physical commonsense-
357 based benchmark for video generation. In *International Conference on Machine Learning*,
358 2025.
- 359 [19] Saman Motamed, Laura Culp, Kevin Swersky, Priyank Jaini, and Robert Geirhos. Do gener-
360 erative video models understand physical principles? In *Proceedings of the IEEE/CVF Win-
361 ter Conference on Applications of Computer Vision*, pages 948–958, March 2026. URL
362 <https://arxiv.org/abs/2501.09038>.
- 363 [20] Siyuan Zhou, Hejun Wang, Hu Cheng, Jinxi Li, Dongsheng Wang, Junwei Jiang, Yixiao Jin,
364 Jiayue Huang, Shiwei Mao, Shangjia Liu, et al. Physinone: Visual physics learning and
365 reasoning in one suite. *arXiv preprint arXiv:2604.09415*, 2026.
- 366 [21] Seung Wook Kim, Yuhao Zhou, Jonah Philion, Antonio Torralba, and Sanja Fidler. Learning to
367 simulate dynamic environments with GameGAN. In *Proceedings of the IEEE/CVF Conference
368 on Computer Vision and Pattern Recognition*, 2020.
- 369 [22] Willi Menapace, Stephane Lathuiliere, Sergey Tulyakov, Aliaksandr Siarohin, and Elisa Ricci.
370 Playable video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision
371 and Pattern Recognition*, 2021.
- 372 [23] Willi Menapace, Stéphane Lathuilière, Aliaksandr Siarohin, Christian Theobalt, Sergey
373 Tulyakov, Vladislav Golyanik, and Elisa Ricci. Playable environments: Video manipula-
374 tion in space and time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
375 Pattern Recognition*, pages 3584–3593, 2022. URL <https://arxiv.org/abs/2203.01914>.
- 376 [24] Willi Menapace, Aliaksandr Siarohin, Stephane Lathuiliere, Panos Achlioptas, Vladislav
377 Golyanik, Sergey Tulyakov, and Elisa Ricci. Promptable game models: Text-guided game
378 simulation via masked diffusion models. *ACM Transactions on Graphics*, 2024.
- 379 [25] Decart, Julian Quevedo, Quinn McIntyre, Spruce Campbell, Xinlei Chen, and Robert Wachen.
380 Oasis: A universe in a transformer. Project page, 2024. URL <https://oasis-model.github.io/>.
381
- 382 [26] Jiwen Yu, Yiran Qin, Xintao Wang, Pengfei Wan, Di Zhang, and Xihui Liu. Gamefactory:
383 Creating new games with generative interactive videos. In *Proceedings of the IEEE/CVF
384 International Conference on Computer Vision*, 2025.
- 385 [27] Junliang Guo, Yang Ye, Tianyu He, Haoyu Wu, Yushu Jiang, Tim Pearce, and Jiang Bian.
386 Mineworld: A real-time and open-source interactive world model on Minecraft. *arXiv preprint
387 arXiv:2504.08388*, 2025.
- 388 [28] Taiye Chen, Xun Hu, Zihan Ding, and Chi Jin. Learning world models for interactive video
389 generation. *arXiv preprint arXiv:2505.21996*, 2025.
- 390 [29] Bowen Baker, Ilge Akkaya, Peter Zhokhov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon
391 Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (VPT): Learning to act by
392 watching unlabeled online videos. In *Advances in Neural Information Processing Systems*, 2022.
393 URL <https://arxiv.org/abs/2206.11795>.

- 394 [30] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew
395 Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. MineDojo: Building open-ended
396 embodied agents with internet-scale knowledge. In *Advances in Neural Information Processing
397 Systems Datasets and Benchmarks Track*, 2022. URL [https://openreview.net/forum?
398 id=rc8o_j8I8PX](https://openreview.net/forum?id=rc8o_j8I8PX).
- 399 [31] Zhen Li, Zian Meng, Shuwei Shi, Wenshuo Peng, Yuwei Wu, Bo Zheng, Chuanhao Li, and
400 Kaipeng Zhang. Wildworld: A large-scale dataset for dynamic world modeling with actions
401 and explicit state toward generative ARPG. *arXiv preprint arXiv:2603.23497*, 2026. URL
402 <https://arxiv.org/abs/2603.23497>.
- 403 [32] Haoyu Wu, Jiwen Yu, Yingtian Zou, and Xihui Liu. Multiworld: Scalable multi-agent multi-
404 view video world models. *arXiv preprint arXiv:2604.18564*, 2026.
- 405 [33] Yuang Wang, Chao Wen, Haoyu Guo, Sida Peng, Minghan Qin, Hujun Bao, Xiaowei Zhou, and
406 Ruizhen Hu. Precise action-to-video generation through visual action prompts. In *Proceedings
407 of the IEEE/CVF International Conference on Computer Vision*, 2025. URL [https://arxiv.
408 org/abs/2508.13104](https://arxiv.org/abs/2508.13104).
- 409 [34] Mohammad Reza Taesiri, Finlay Macklon, and Cor-Paul Bezemer. CLIP meets GamePhysics:
410 Towards bug identification in gameplay videos using zero-shot transfer learning. In *2022
411 IEEE/ACM 19th International Conference on Mining Software Repositories*, pages 270–281,
412 2022. doi: 10.1145/3524842.3528438. URL <https://arxiv.org/abs/2203.11096>.
- 413 [35] Oliver Groth, Fabian B. Fuchs, Ingmar Posner, and Andrea Vedaldi. ShapeStacks: Learning
414 vision-based physical intuition for generalised object stacking. In *European Conference on
415 Computer Vision*, 2018. URL <https://arxiv.org/abs/1804.08018>.
- 416 [36] Kevin A. Smith, Lingjie Mei, Shunyu Yao, Jiajun Wu, Elizabeth S. Spelke, Joshua B.
417 Tenenbaum, and Tomer D. Ullman. Modeling expectation violation in intuitive physics
418 with coarse probabilistic object representations. In *Advances in Neural Information Pro-
419 cessing Systems*, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/
420 e88f243bf341ded9b4ced444795c3f17-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/e88f243bf341ded9b4ced444795c3f17-Abstract.html).
- 421 [37] Shiqian Li, Kewen Wu, Chi Zhang, and Yixin Zhu. I-PHYRE: Interactive physical reasoning.
422 In *International Conference on Learning Representations*, 2024. URL [https://openreview.
423 net/forum?id=1bbPQShCT2](https://openreview.net/forum?id=1bbPQShCT2).
- 424 [38] Ronan Riochet, Mario Ynocente Castro, Mathieu Bernard, Adam Lerer, Rob Fergus, Veronique
425 Izard, and Emmanuel Dupoux. IntPhys 2019: A benchmark for visual intuitive physics
426 understanding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5016–
427 5025, 2022.
- 428 [39] Florian Bordes, Quentin Garrido, Justine T. Kao, Adina Williams, Michael Rabbat, and Em-
429 manuel Dupoux. IntPhys 2: Benchmarking intuitive physics understanding in complex synthetic
430 environments. *arXiv preprint arXiv:2506.09849*, 2025. URL [https://arxiv.org/abs/
431 2506.09849](https://arxiv.org/abs/2506.09849).
- 432 [40] Rohit Girdhar and Deva Ramanan. CATER: A diagnostic dataset for compositional actions and
433 TEmporal reasoning. In *International Conference on Learning Representations*, 2020. URL
434 <https://arxiv.org/abs/1910.04744>.
- 435 [41] Fabien Baradel, Natalia Neverova, Julien Mille, Greg Mori, and Christian Wolf. CoPhy: Counter-
436 factual learning of physical dynamics. In *International Conference on Learning Representations*,
437 2020. URL <https://openreview.net/forum?id=SkeyppEFvS>.
- 438 [42] Zhenfang Chen, Kexin Yi, Yunzhu Li, Mingyu Ding, Antonio Torralba, Joshua B. Tenenbaum,
439 and Chuang Gan. ComPhy: Compositional physical reasoning of objects and events from videos.
440 In *International Conference on Learning Representations*, 2022. URL [https://openreview.
441 net/forum?id=PgNEYaIc81Q](https://openreview.net/forum?id=PgNEYaIc81Q).

- 442 [43] Maitreya Patel, Tejas Gokhale, Chitta Baral, and Yezhou Yang. CRIPP-VQA: Counterfactual
443 reasoning about implicit physical properties via video question answering. In *Proceedings of*
444 *the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9856–9870,
445 Abu Dhabi, United Arab Emirates, 2022. Association for Computational Linguistics. doi: 10.
446 18653/v1/2022.emnlp-main.670. URL [https://aclanthology.org/2022.emnlp-main.](https://aclanthology.org/2022.emnlp-main.670/)
447 670/.
- 448 [44] Zhicheng Zheng, Xin Yan, Zhenfang Chen, Jingzhou Wang, Qin Zhi Eddie Lim, Joshua B.
449 Tenenbaum, and Chuang Gan. ContPhy: Continuum physical concept learning and reasoning
450 from videos. In *Proceedings of the 41st International Conference on Machine Learning*, volume
451 235 of *Proceedings of Machine Learning Research*, pages 61526–61558. PMLR, 2024. URL
452 <https://proceedings.mlr.press/v235/zheng241.html>.
- 453 [45] Aaron Foss, Chloe Evans, Sasha Mitts, Koustuv Sinha, Ammar Rizvi, and Justine T. Kao.
454 CausalVQA: A physically grounded causal reasoning benchmark for video models. *arXiv*
455 *preprint arXiv:2506.09943*, 2025.
- 456 [46] Puyin Li, Tiange Xiang, Ella Mao, Shirley Wei, Xinye Chen, Adnan Masood, Li Fei-Fei, and
457 Ehsan Adeli. QuantiPhy: A quantitative benchmark evaluating physical reasoning abilities of
458 vision-language models. *arXiv preprint arXiv:2512.19526*, 2025.
- 459 [47] Dian Zheng, Ziqi Huang, Hongbo Liu, Kai Zou, Yinan He, Fan Zhang, Lulu Gu, Yuanhan
460 Zhang, Jingwen He, Wei-Shi Zheng, Yu Qiao, and Ziwei Liu. VBench-2.0: Advancing video
461 generation benchmark suite for intrinsic faithfulness. *arXiv preprint arXiv:2503.21755*, 2025.
462 URL <https://arxiv.org/abs/2503.21755>.
- 463 [48] Haoyi Duan, Hong-Xing Yu, Sirui Chen, Li Fei-Fei, and Jiajun Wu. Worldscore: A unified
464 evaluation benchmark for world generation. In *Proceedings of the IEEE/CVF International*
465 *Conference on Computer Vision*, 2025.
- 466 [49] Minh-Quan Le, Yuanzhi Zhu, Vicky Kalogeiton, and Dimitris Samaras. What about gravity
467 in video generation? post-training newton’s laws with verifiable rewards. *arXiv preprint*
468 *arXiv:2512.00425*, 2025.
- 469 [50] Shaowei Liu, Zhongzheng Ren, Saurabh Gupta, and Shenlong Wang. PhysGen: Rigid-body
470 physics-grounded image-to-video generation. In *European Conference on Computer Vision*,
471 2024. doi: 10.1007/978-3-031-73007-8_21. URL <https://arxiv.org/abs/2409.18964>.
- 472 [51] Nate Gillman, Charles Herrmann, Michael Freeman, Daksh Aggarwal, Evan Luo, Deqing
473 Sun, and Chen Sun. Force prompting: Video generation models can learn and generalize
474 physics-based control signals. In *Advances in Neural Information Processing Systems*, 2025.
475 URL <https://arxiv.org/abs/2505.19386>.
- 476 [52] Chen Wang, Chuha Chen, Yiming Huang, Zhiyang Dou, Yuan Liu, Jiatao Gu, and Lingjie Liu.
477 PhysCtrl: Generative physics for controllable and physics-grounded video generation. *arXiv*
478 *preprint arXiv:2509.20358*, 2025.
- 479 [53] Haoze Zhang, Tianyu Huang, Zichen Wan, Xiaowei Jin, Hongzhi Zhang, Hui Li, and Wangmeng
480 Zuo. PhysChoreo: Physics-controllable video generation with part-aware semantic grounding.
481 *arXiv preprint arXiv:2511.20562*, 2025.
- 482 [54] Yu Yuan, Xijun Wang, Tharindu Wickremasinghe, Zeeshan Nadir, Bole Ma, and Stanley H.
483 Chan. NewtonGen: Physics-consistent and controllable text-to-video generation via neural
484 newtonian dynamics. In *International Conference on Learning Representations*, 2026. URL
485 <https://openreview.net/forum?id=rJ6N6sunaU>.
- 486 [55] Sriram Narayanan, Ziyu Jiang, Srinivasa Narasimhan, and Manmohan Chandraker. Phyco:
487 Learning controllable physical priors for generative motion, 2026. URL [https://arxiv.org/](https://arxiv.org/abs/2604.28169)
488 [abs/2604.28169](https://arxiv.org/abs/2604.28169).
- 489 [56] Chenyu Li, Oscar Michel, Xichen Pan, Sainan Liu, Mike Roberts, and Saining Xie. PISA
490 experiments: Exploring physics post-training for video diffusion models by watching stuff drop.
491 In *International Conference on Machine Learning*, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2503.09595)
492 [2503.09595](https://arxiv.org/abs/2503.09595).

- 493 [57] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti,
494 Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. AI2-THOR: An interactive 3d
495 environment for visual AI. *arXiv preprint arXiv:1712.05474*, 2017.
- 496 [58] Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha Kembhavi, Eric Kolve, Roozbeh Mot-
497 taghi, Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, Luca Weihs, Mark
498 Yatskar, and Ali Farhadi. RoboTHOR: An open simulation-to-real embodied AI platform. In
499 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
500 URL <https://arxiv.org/abs/2004.06799>.
- 501 [59] Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Jordi Salvador, Kiana Ehsani, Win-
502 son Han, Eric Kolve, Ali Farhadi, Aniruddha Kembhavi, and Roozbeh Mottaghi. ProcTHOR:
503 Large-scale embodied AI using procedural generation. In *Advances in Neural Information Pro-
504 cessing Systems*, 2022. URL [https://proceedings.neurips.cc/paper_files/paper/
505 2022/hash/27c546ab1e4f1d7d638e6a8dfbad9a07-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/27c546ab1e4f1d7d638e6a8dfbad9a07-Abstract-Conference.html). Out-
506 standing Paper Award.
- 507 [60] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana
508 Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for
509 embodied AI research. In *Proceedings of the IEEE/CVF International Conference on Computer
510 Vision*, 2019.
- 511 [61] Fei Xia, Amir R. Zamir, Zhiyang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gib-
512 son env: Real-world perception for embodied agents. In *Proceedings of the IEEE Conference on
513 Computer Vision and Pattern Recognition*, 2018. URL [https://openaccess.thecvf.com/
514 content_cvpr_2018/html/Xia_Gibson_Env_Real-World_CVPR_2018_paper.html](https://openaccess.thecvf.com/content_cvpr_2018/html/Xia_Gibson_Env_Real-World_CVPR_2018_paper.html).
- 515 [62] Chengshu Li, Fei Xia, Roberto Martín-Martín, Michael Lingelbach, Sanjana Srivastava, Bokui
516 Shen, Kent Elliott Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, Andrey Kurenkov, Karen
517 Liu, Hyowon Gweon, Jiajun Wu, Li Fei-Fei, and Silvio Savarese. iGibson 2.0: Object-centric
518 simulation for robot learning of everyday household tasks. In *Proceedings of the 5th Conference
519 on Robot Learning*, volume 164 of *Proceedings of Machine Learning Research*, pages 455–465.
520 PMLR, 2022. URL <https://proceedings.mlr.press/v164/li22b.html>.
- 521 [63] Chuang Gan, Jeremy Schwartz, Seth Alter, Damian Mrowca, Martin Schrimpf, James Traer,
522 Julian De Freitas, Jonas Kubilius, Abhishek Bhandwaldar, Nick Haber, et al. ThreeDWorld: A
523 platform for interactive multi-modal physical simulation. In *Advances in Neural Information
524 Processing Systems Datasets and Benchmarks Track*, 2021.
- 525 [64] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J.
526 Fleet, Daniel Gnanapragasam, Florian Golemo, Charles Herrmann, et al. Kubric: A scalable
527 dataset generator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
528 Recognition*, 2022.
- 529 [65] Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio
530 Torralba. VirtualHome: Simulating household activities via programs. In *Proceedings of the
531 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8494–8502, 2018.
532 URL <https://arxiv.org/abs/1806.07011>.
- 533 [66] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen, Sanjana Srivastava, Roberto Martín-
534 Martín, Chen Wang, Gabriel Levine, Wensi Ai, Benjamin Martinez, et al. BEHAVIOR-1k: A
535 human-centered, embodied AI benchmark with 1,000 everyday activities and realistic simulation.
536 *arXiv preprint arXiv:2403.09227*, 2024.
- 537 [67] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu,
538 Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X. Chang, Leonidas J. Guibas, and Hao Su.
539 SAPIEN: A simulated part-based interactive environment. In *Proceedings of the IEEE/CVF
540 Conference on Computer Vision and Pattern Recognition*, 2020. URL [https://arxiv.org/
541 abs/2003.08515](https://arxiv.org/abs/2003.08515).
- 542 [68] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J. Davison. RL Bench: The
543 robot learning benchmark and learning environment. *IEEE Robotics and Automation Letters*, 5
544 (2):3019–3026, 2020. URL <https://arxiv.org/abs/1909.12271>.

- 545 [69] Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang,
546 Stone Tao, Xinyue Wei, Yunchao Yao, Xiaodi Yuan, Pengwei Xie, Zhiao Huang, Rui Chen,
547 and Hao Su. ManiSkill2: A unified benchmark for generalizable manipulation skills. In
548 *International Conference on Learning Representations*, 2023. URL [https://arxiv.org/
549 abs/2302.04659](https://arxiv.org/abs/2302.04659).
- 550 [70] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles
551 Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac
552 Gym: High performance GPU-based physics simulation for robot learning. In *Advances
553 in Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. URL
554 <https://arxiv.org/abs/2108.10470>.
- 555 [71] Weichao Qiu and Alan Yuille. UnrealCV: Connecting computer vision to unreal engine. *arXiv
556 preprint arXiv:1609.01326*, 2016. URL <https://arxiv.org/abs/1609.01326>.
- 557 [72] Fangwei Zhong, Kui Wu, Churan Wang, Hao Chen, Hai Ci, Zhoujun Li, and Yizhou Wang.
558 UnrealZoo: Enriching photo-realistic virtual worlds for embodied AI. In *Proceedings of the
559 IEEE/CVF International Conference on Computer Vision*, 2025. URL [https://arxiv.org/
560 abs/2412.20977](https://arxiv.org/abs/2412.20977). Highlight.
- 561 [73] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao
562 Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie
563 Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei Li, Kaixin
564 Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang Liu, Dayiheng
565 Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen Meng, Xuancheng
566 Ren, Xingzhang Ren, Siboz Song, Yuchong Sun, Jun Tang, Jianhong Tu, Jianqiang Wan, Peng
567 Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao Xie, Yiheng Xu, Haiyang Xu, Jin
568 Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang, Bowen Yu, Fei Zhang, Hang Zhang,
569 Xi Zhang, Bo Zheng, Humen Zhong, Jingren Zhou, Fan Zhou, Jing Zhou, Yuanzhi Zhu, and
570 Ke Zhu. Qwen3-VL technical report, 2025.
- 571 [74] Wan Team, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Fei Wu Yu,
572 Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, et al.
573 Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*,
574 2025.
- 575 [75] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,
576 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International
577 Conference on Learning Representations*, 2022.
- 578 [76] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David
579 Novotny. VGGT: Visual geometry grounded transformer. In *Proceedings of the IEEE/CVF
580 Conference on Computer Vision and Pattern Recognition*, 2025.

581 A Asset Library, Character Setup, and Capture Configuration

582 A.1 Asset Library

583 PhysEditWorld is built from a manually curated UE5 asset library designed to cover diverse gravity-
584 sensitive interaction scenarios. The current asset collection contains 12 representative scene environ-
585 ments, including polar research bases, indoor sports courts, lunar and Martian surfaces, cyberpunk city
586 blocks, classrooms, laboratories, forest valleys, research camps, bowling alleys, theaters, and modern
587 urban streets. These assets span indoor, outdoor, planetary, urban, sports, and object-interaction
588 settings, allowing the dataset to expose different forms of gravity-dependent behavior such as free
589 fall, jump arcs, landing timing, object displacement, and contact response. Artists manually inspect
590 each scene for visual fidelity, collision quality, interaction suitability, and physical stability before it
591 is included in the replay pipeline.

592 In addition to scene assets, PhysEditWorld includes three character assets with compatible animation
593 and retargeting setups. The character system uses a UE5 animation stack based on motion matching,
594 pose search, chooser-driven animation selection, blend stacks, orientation warping, and animation
595 warping. To support low-gravity scenarios, we implement a low-gravity animation tuner that adjusts
596 airborne-state detection and animation playback rate according to the gravity multiplier. This prevents
597 low-gravity clips from using visually implausible Earth-gravity animation timing. The camera system
598 supports socket-based attachment to skeletal bones, such as head, spine, or pelvis sockets, so first-
599 person and actor-following views can inherit smooth character motion during jumping, falling, and
600 landing. External character meshes are connected to the animation library through IK Rig and IK
601 Retargeter assets, enabling different characters to share a consistent motion-control interface during
602 counterfactual replay.

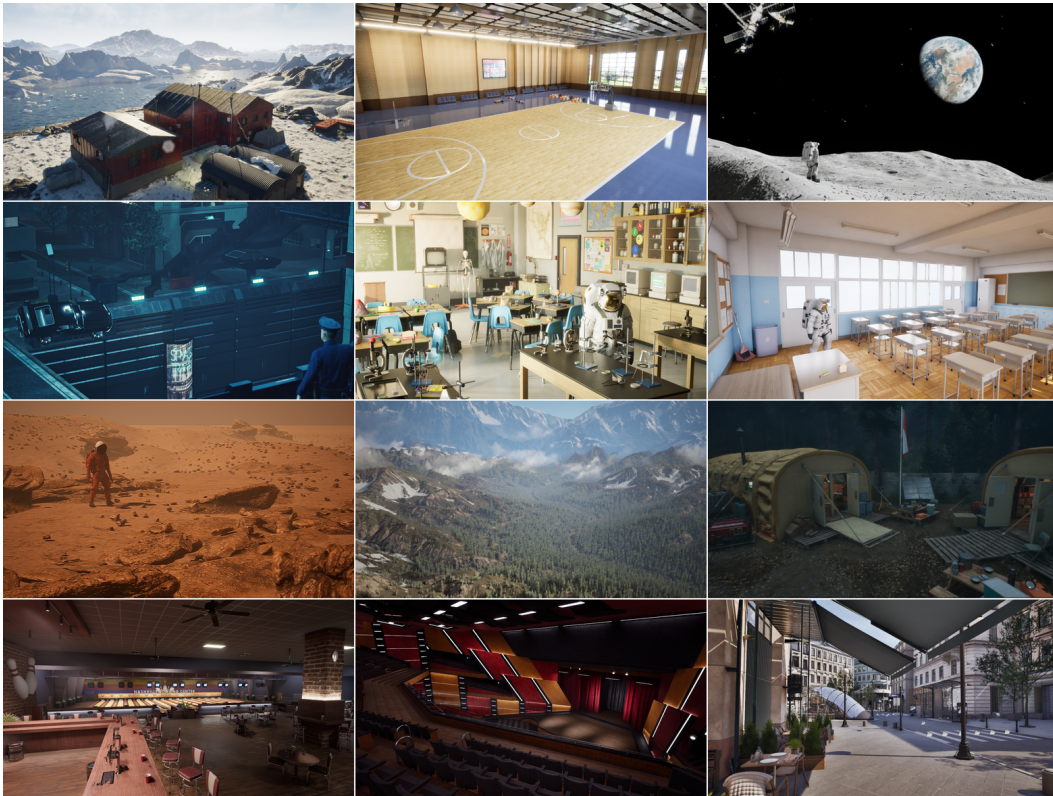


Figure 5: The 12 curated UE5 scene assets used in PhysEditWorld. The asset library covers indoor, outdoor, planetary, urban, sports, and object-interaction environments to support diverse gravity-sensitive motion patterns.



Figure 6: Character assets used in the current PhysEditWorld pipeline. Characters are configured with compatible controllers, animation retargeting, and low-gravity animation tuning to support reproducible replay under edited gravity settings.

603 PhysEditWorld is built from a manually curated UE5 asset library designed to cover diverse gravity-
604 sensitive interaction scenarios. The current asset collection contains 12 representative scene environ-
605 ments, including polar research bases, indoor sports courts, lunar and Martian surfaces, cyberpunk city
606 blocks, classrooms, laboratories, forest valleys, research camps, bowling alleys, theaters, and modern
607 urban streets. These assets span indoor, outdoor, planetary, urban, sports, and object-interaction
608 settings, allowing the dataset to expose different forms of gravity-dependent behavior such as free
609 fall, jump arcs, landing timing, object displacement, and contact response. Artists manually inspect
610 each scene for visual fidelity, collision quality, interaction suitability, and physical stability before it
611 is included in the replay pipeline.

612 In addition to scene assets, PhysEditWorld includes character assets with compatible controller,
613 animation, and retargeting setups. The character system uses a UE5 animation stack based on motion
614 matching, pose search, chooser-driven animation selection, blend stacks, orientation warping, and
615 animation warping. To support low-gravity scenarios, we implement a low-gravity animation tuner
616 that adjusts airborne-state detection and animation playback rate according to the gravity multiplier.
617 This prevents low-gravity clips from using visually implausible Earth-gravity animation timing.
618 External character meshes are connected to the shared animation library through IK Rig and IK
619 Retargeter assets, enabling different characters to share a consistent motion-control interface during
620 counterfactual replay.

621 A.2 Character-Centric Multi-Camera Capture

622 The capture system uses a character-centric multi-camera rig rather than placing independent camera
623 actors in the scene. Multiple UCameraComponents are embedded inside the character blueprint and
624 attached through skeletal sockets or bones, spring arms, and camera components. This topology
625 provides a shared character-centered reference frame for all camera views, so first-person, third-
626 person, side, front, back, and oblique views remain temporally aligned during walking, turning,
627 jumping, falling, and landing. In the current setup, the camera array includes FP, BK, FL, FW, FR, LF,
628 RT, and TP views. Rendering all views from the same replay instance avoids the synchronization drift
629 that would arise from repeated simulation or separately placed cameras.

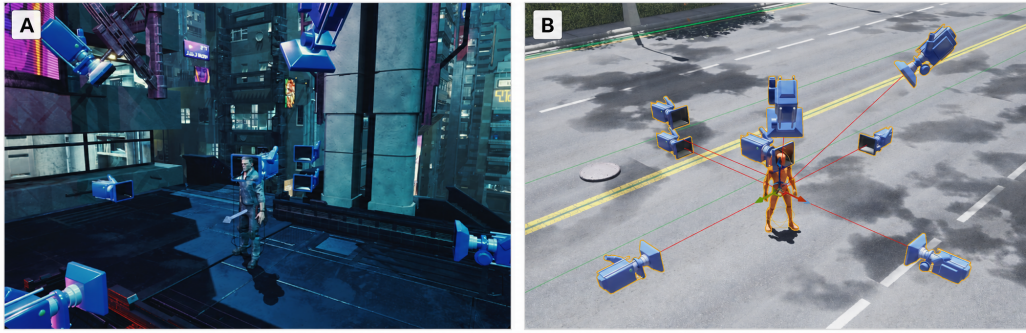


Figure 7: Character-centric multi-camera rig used for synchronized capture. (A) Cameras are embedded around the actor inside a UE5 scene. (B) The rig defines multiple named views around the same character reference frame, including first-person, third-person, side, front, and back views.

630 Socket-based attachment also supports stable close-range observation. Cameras can be mounted to
 631 head, spine, pelvis, or other skeletal sockets, while spring arms decouple the choice of skeletal anchor
 632 from the camera offset and near-body framing. During interactive preview, spring arms can provide
 633 collision handling and visual smoothing. During dataset replay and rendering, camera lag is disabled
 634 to prioritize deterministic frame alignment and reproducible multi-view sampling. This design is
 635 especially important for gravity-edited rollouts, where jump and fall timing are the measurement
 636 target rather than merely cinematic motion.

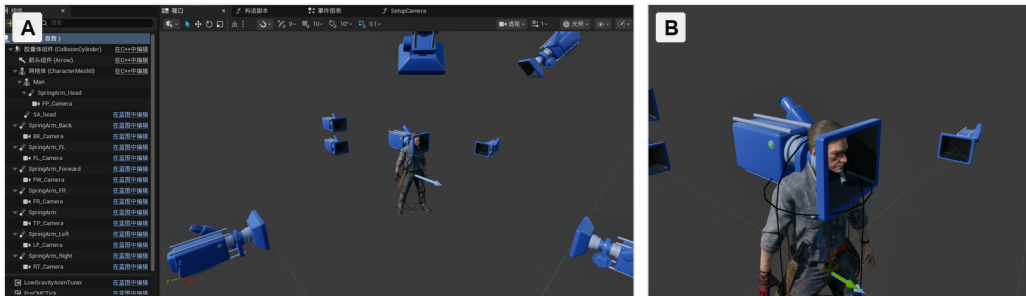


Figure 8: Socket and spring-arm based camera setup. (A) The character blueprint contains multiple spring-arm and camera components attached around the skeletal mesh. (B) Close-range cameras are mounted through skeletal sockets and spring arms, allowing first-person or near-body views to inherit character motion while preserving controllable camera offsets.

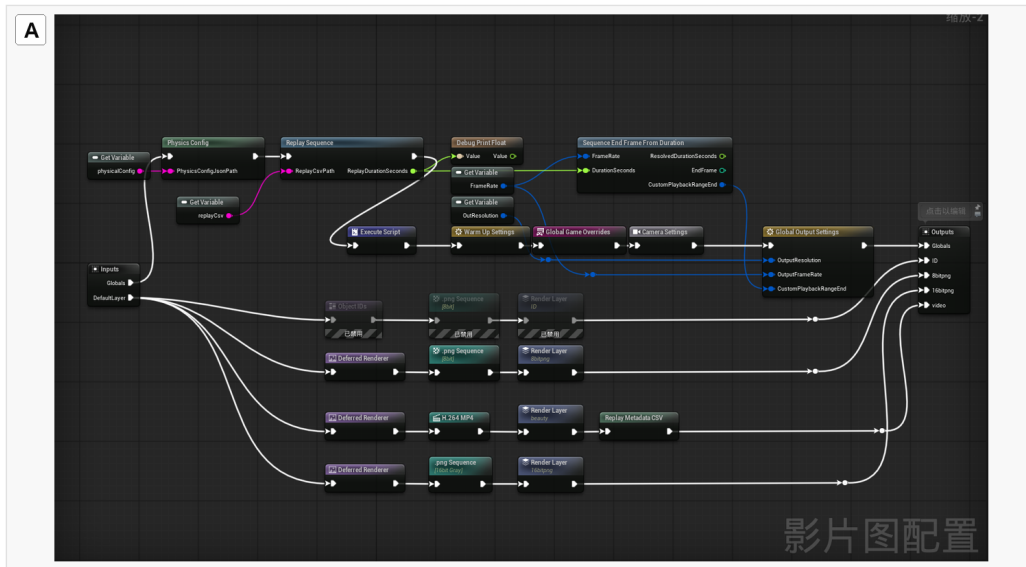


Figure 9: Movie Render Graph configuration for synchronized multimodal export. The graph combines replay sequence loading, warm-up, global game and camera settings, output settings, multi-branch visual exports, H.264 preview video, high-precision depth output, and replay metadata CSV export.

638 The UE5 rendering job is configured as a multi-branch Movie Render Graph. A shared global
 639 settings chain controls warm-up, game overrides, camera settings, and output settings. The current
 640 configuration uses a warm-up interval before formal capture so that animation state, temporal
 641 antialiasing, post-processing, and motion blur reach a stable state before frames are written. Captures
 642 are exported at 1280×720 resolution and 30 FPS.

643 The render graph produces several synchronized output branches from the same replay. An object-ID
 644 or segmentation branch exports per-frame ID supervision. An 8-bit image branch exports standard
 645 image sequences and auxiliary render passes such as motion vectors and world normals. A video
 646 branch exports H.264 MP4 previews with audio for manual inspection and semantic quality control.
 647 A 16-bit PNG branch exports high-precision depth-style supervision. In parallel, replay metadata is
 648 written as frame-level and time-level CSV files, together with a camera-name mapping JSON file.
 649 These metadata records include replay input events, player and camera transforms, physical state,
 650 camera identity, and capture termination signals.

651 B UE Editor Plugin and Data Generation Workflow

652 PhysEditWorld is generated through a UE5 Editor plugin that manages the data-generation workflow
 653 from scene preparation to large-scale synthetic data production. The plugin is designed to operate
 654 directly on artist-authored UE5 levels with minimal intrusion into existing game content. Instead
 655 of rebuilding scenes in a separate simulator, it augments prepared levels with the runtime contexts,
 656 replay components, camera bindings, and batch-generation interfaces required for reproducible
 657 replay-and-rendering.

658 A key design goal of the plugin is to make scene preparation lightweight. As illustrated in Figure ??,
 659 each artist-prepared level can be converted into a production-ready data-generation scene through
 660 four editor operations: selecting the required EnhancedInputContext, specifying the sequence and
 661 camera bindings, running the FactoryManager initialization, and saving the produce-ready level.

662 Since most updated character controller implementations already expose Enhanced Input support,
663 this design remains extensible across different controllers and input schemes.

664 After the level-specific configuration is completed, the FactoryManager initializes the scene with a
665 single command. It automatically registers the data-collection context, attaches replay components,
666 inserts the required capture and rendering modules, and validates the camera and input bindings. This
667 step injects only the minimal components needed for data generation, preserving the original authored
668 scene and gameplay logic as much as possible. Once initialization finishes, the level becomes ready
669 for reproducible replay and rendering.

670 The plugin then supports two ways of constructing interaction sequences. The first is Play-in-Editor
671 interaction capture, where users directly play the prepared level and record semantic input-action
672 traces, including movement axes, jump commands, camera deltas, and button states. The second is
673 procedural sequence generation, where PCG-based scripts generate scalable interaction sequences for
674 broader coverage. We also combined the Navmesh Agent in our pipeline. Both sources are converted
675 into replayable interaction sequences that can be executed under controlled physical configurations.

676 For large-scale production, the DataFactory pipeline exposes Python interfaces and YAML configura-
677 tion files for composing replay, rendering, cleaning, merging, and recovery stages. The batch system
678 supports checkpointing and resume functionality, allowing interrupted jobs to continue from the last
679 valid state rather than restarting from scratch. This makes the same editor-plugin workflow usable
680 both for interactive debugging inside UE5 and for large-scale offline dataset generation.

1. Assets Prepare

Artist-Prepared Ue5 Level

Controller and Character

Selected

2. Use Factory Manager To Init

- Add Replay Pawn
- Add Camera Comments
- Register Input Context
- Generate Physical Config
- ...

One-Click

3. Interact Sequence Construct

Play-in-Editor Interaction Caption

Navmesh Agent Sequence Generate

PCG-Based Sequence Generate

4. Scripts batch Produce

Config.yaml → Unreal Engine → PhysEdit World

Intermediate steps: Clean, Merge, Caption

1. Open the Factory Manager

2. Select the Input Content to record

3. Choose the Character Controller to Spawn

4. One-Click to bring this level to produce-ready

four step to bring a level to produce ready

helpful tools to generate configs, all in one

Choose which camera need to be recorded

Figure 10: UE Editor Plugin and Data Generation Workflow

681 **C Dataset Sample Format and Metadata Schema**

682 Each data sample in our dataset corresponds to one complete recording-rendering job under a fixed
 683 scene, action sequence, camera setup, and physical configuration. We use an 8-digit string as the
 684 sample identifier, e.g., 00000000. A sample is self-contained and consists of synchronized multi-view
 685 RGB videos, frame-level metadata, event-level metadata, and a physical-configuration descriptor.
 686 The dataset root contains a global index file, `Details.json`, which stores the list of samples and the
 687 relative paths to their associated files.

688 The cleaned dataset follows the directory layout below:

```
689 <dataset_root>/
690   Video/<sample_id>/<camera_name>.mp4
691   Meta/<sample_id>_meta_frame.csv
692   Meta/<sample_id>_meta_time.csv
693   PhysicalConfig/<physical_config_name>.json
694   Details.json
```

695 For auxiliary modalities, such as depth, surface normals, or object masks, the optional frame-wise
 696 outputs are stored as:

```
697 Aux/<aux_type>/<sample_id>/<camera_name>/<frame_id>.png
```

698 Table 4 summarizes the main fields used to describe each sample.

Table 4: Main fields of a dataset sample.

Field	Description
<code>sample_id</code>	Unique 8-digit identifier of the sample.
<code>scene_name</code>	Name of the simulated scene or environment.
<code>action_sequence_name</code>	Identifier of the replayed action or motion sequence.
<code>physical_config_name</code>	Name of the physical condition used for simulation.
<code>physical_config_file</code>	Relative path to the physical-configuration JSON file.
<code>frame_meta_file</code>	Relative path to frame-level metadata.
<code>time_meta_file</code>	Relative path to event-level temporal metadata.
<code>camera_names</code>	List of camera names available in the sample.
<code>cameras</code>	Per-camera video paths and attributes.
<code>aux_types</code>	Optional auxiliary modalities associated with the sample.

699 A representative sample entry is shown below. All paths are relative to the dataset root.

```
700 {
701   "sample_id": "00000000",
702   "scene_name": "Map_MarsRover",
703   "action_sequence_name": "MarsSequence_0001",
704   "physical_config_name": "G_0.05",
705   "physical_config_file": "PhysicalConfig/G_0.05.json",
706   "frame_meta_file": "Meta/00000000_meta_frame.csv",
707   "time_meta_file": "Meta/00000000_meta_time.csv",
708   "camera_names": ["BK_Camera", "FL_Camera", "FP_Camera", ... ],
709   "aux_types": ["Depth", "Normals"],
710   "cameras": [
711     {
712       "camera_name": "BK_Camera",
713       "rgb_file": "Video/00000000/BK_Camera.mp4",
714       "aux": {
715         "Depth": "Aux/Depth/00000000/BK_Camera",
716         "Normals": "Aux/Normals/00000000/BK_Camera",
717         "ObjectMask": "Aux/ObjectMask/00000000/BK_Camera"
718       }
719     }
720   ]
721 }
```

```

719     "resolution": {"width": 1280, "height": 720},
720     "fps": 30
721   },
722   ...
723 ]
724 }

```

725 The two metadata files provide complementary temporal annotations. The frame-level file records
726 information indexed by rendered frame, which is used to align visual observations with simulation
727 states. The event-level file records simulation and capture events over time, such as the start and end
728 of recording or replay. The physical-configuration file stores the controlled simulation parameters,
729 enabling downstream methods to condition learning or evaluation on explicit physical settings.

730 **D Example usage.**

731 The dataset generation pipeline is driven by a YAML configuration file and is launched through the
732 repository-level entry point `Scripts/cli.py`. The configuration specifies the simulated scenes,
733 replay trajectories, physical-configuration files, Movie Render Graph preset, output resolution, frame
734 rate, and the enabled pipeline stages.

735 Before running the full pipeline, we first validate the configuration and inspect the planned jobs:

```

736 uv run --with pyyaml python Scripts/cli.py \
737   --config Scripts/configs/base.yaml \
738   --dry-run

```

739 After verifying the generated plan, we run the complete rendering and cleaning pipeline:

```

740 uv run --with pyyaml python Scripts/cli.py \
741   --config Scripts/configs/base.yaml

```

742 A typical configuration contains the following fields:

```

743 base_path: <project_root>
744 unreal_exe: <path_to_UnrealEditor>
745 dataset_path: <raw_dataset_root>
746 output_path: <cleaned_dataset_root>
747 input_sequence: <replay_csv_root>
748 physical_config: <physical_config_root>
749
750 scenes:
751   - name: <scene_name>
752     level_path: <unreal_level_asset_path>
753     level_sequence_path: <unreal_level_sequence_asset_path>
754
755 job:
756   selection: auto
757   mrg_path: <movie_render_graph_asset_path>
758   params:
759     fps: 30
760     resolution: [1280, 720]
761     mp4: true
762
763 pipeline:
764   render:
765     enable: true
766     auto_start: true
767   clean:
768     enable: true
769     input_root: <stage1_render_output_root>

```

```
770     output_root: <cleaned_dataset_root>
771     physical_config_root: <physical_config_root>
772     append: false
773     workers: 8
774     overwrite: false
775     dataset_name: <dataset_name>
```

776 The command above expands the configuration into rendering jobs over the Cartesian product of
777 scenes, replay trajectories, and physical configurations. It then launches Unreal Engine to execute
778 the rendering jobs and, after rendering, converts the stage-1 outputs into the canonical dataset layout
779 described in Appendix C.

780 When the stage-1 rendering outputs have already been generated, the cleaning step can also be
781 executed independently:

```
782 uv run python Scripts/utools/clean_mrqa_graph_dataset.py \
783     <stage1_render_output_root> \
784     --output-root <cleaned_dataset_root> \
785     --physical-config-root <physical_config_root> \
786     --workers 8 \
787     --dataset-name <dataset_name>
```

788 This cleaning-only command scans the stage-1 output directories, identifies valid rendering jobs by
789 the presence of `camera_name_map.json` and camera videos, copies the RGB videos and metadata
790 into the canonical layout, and writes the global index file `Details.json`.

791 **E More Asset**

792 A small set of data is showed in the supplementary material with a html for the limitation of size.

793 **F Release Plan**

794 We plan to publicly release the PhysEditWorld dataset, including synchronized RGB videos, depth
795 maps, normal maps, gravity annotations, camera trajectories, action traces, and evaluation scripts used
796 in this work. We have provided a small subset in our supplementary material for the size limitation.

797 We also plan to release the UE5-based replay-and-rendering pipeline **as a plugin**, together with data
798 generation configurations and example replay assets, to facilitate reproducibility and future research
799 on physics-editable world modeling. Due to potential licensing restrictions associated with certain
800 third-party UE5 assets, some raw scene assets or commercial content may not be redistributed directly.
801 In such cases, we will provide replacement assets, asset lists, or reconstruction instructions whenever
802 possible.

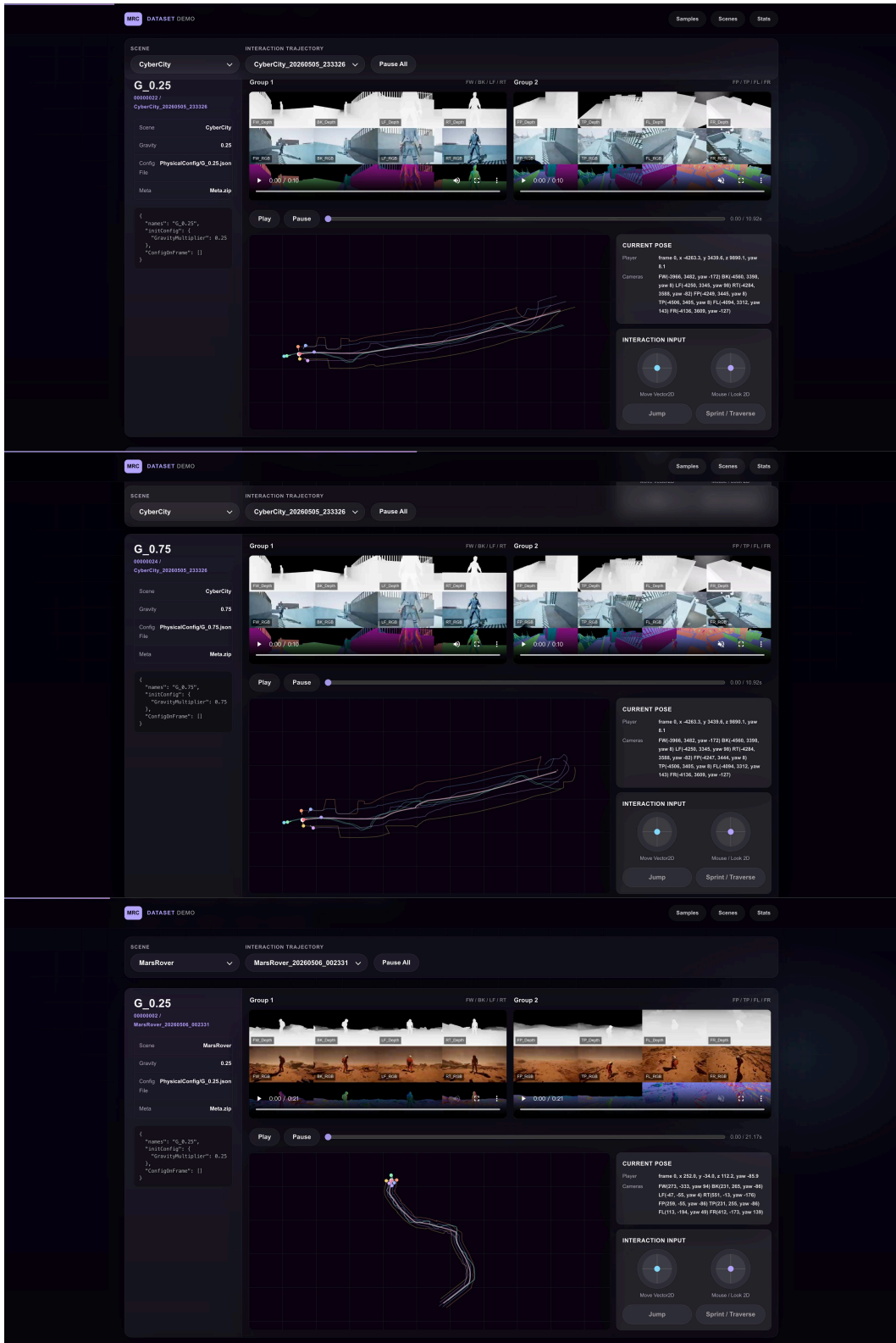


Figure 11: More PhysicsEdit-Data



Figure 12: More PhysicsEdit-Data with other physics config like friction, bounce coefficient



Figure 13: More PhysicsEdit-Data with other physics config like friction, bounce coefficient

803 NeurIPS Paper Checklist

804 1. Claims

805 Question: Do the main claims made in the abstract and introduction accurately reflect the
806 paper’s contributions and scope?

807 Answer: [Yes]

808 Justification: The abstract and introduction state that PhysEditWorld is a large-scale multi-
809 modal dataset for gravity-conditioned and physics-editable game world modeling, and the
810 claims are supported by the dataset construction, replay protocol, and utility studies. We
811 explicitly scope the current release to gravity as the editable physical attribute.

812 Guidelines:

- 813 • The answer [N/A] means that the abstract and introduction do not include the claims
814 made in the paper.
- 815 • The abstract and/or introduction should clearly state the claims made, including the
816 contributions made in the paper and important assumptions and limitations. A [No] or
817 [N/A] answer to this question will not be perceived well by the reviewers.
- 818 • The claims made should match theoretical and experimental results, and reflect how
819 much the results can be expected to generalize to other settings.
- 820 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
821 are not attained by the paper.

822 2. Limitations

823 Question: Does the paper discuss the limitations of the work performed by the authors?

824 Answer: [Yes]

825 Justification: We discuss that the current release focuses on gravity and does not yet cover
826 other physical attributes such as friction, drag, restitution, wind, or object-level physical
827 parameters. We also discuss limitations of video-based gravity evaluation, including scale
828 ambiguity and reliance on reliable camera-trajectory estimation.

829 Guidelines:

- 830 • The answer [N/A] means that the paper has no limitation while the answer [No] means
831 that the paper has limitations, but those are not discussed in the paper.
- 832 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 833 • The paper should point out any strong assumptions and how robust the results are to
834 violations of these assumptions (e.g., independence assumptions, noiseless settings,
835 model well-specification, asymptotic approximations only holding locally). The authors
836 should reflect on how these assumptions might be violated in practice and what the
837 implications would be.
- 838 • The authors should reflect on the scope of the claims made, e.g., if the approach was
839 only tested on a few datasets or with a few runs. In general, empirical results often
840 depend on implicit assumptions, which should be articulated.
- 841 • The authors should reflect on the factors that influence the performance of the approach.
842 For example, a facial recognition algorithm may perform poorly when image resolution
843 is low or images are taken in low lighting. Or a speech-to-text system might not be
844 used reliably to provide closed captions for online lectures because it fails to handle
845 technical jargon.
- 846 • The authors should discuss the computational efficiency of the proposed algorithms
847 and how they scale with dataset size.
- 848 • If applicable, the authors should discuss possible limitations of their approach to
849 address problems of privacy and fairness.
- 850 • While the authors might fear that complete honesty about limitations might be used by
851 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
852 limitations that aren’t acknowledged in the paper. The authors should use their best
853 judgment and recognize that individual actions in favor of transparency play an impor-
854 tant role in developing norms that preserve the integrity of the community. Reviewers
855 will be specifically instructed to not penalize honesty concerning limitations.

856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [N/A]

Justification:

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification:

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

910 **5. Open access to data and code**

911 Question: Does the paper provide open access to the data and code, with sufficient instruc-
912 tions to faithfully reproduce the main experimental results, as described in supplemental
913 material?

914 Answer: [Yes]

915 Justification: We offer a small subset of our data in the material for limitation of size.

916 Guidelines:

- 917 • The answer [N/A] means that paper does not include experiments requiring code.
- 918 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
919 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 920 • While we encourage the release of code and data, we understand that this might not
921 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
922 including code, unless this is central to the contribution (e.g., for a new open-source
923 benchmark).
- 924 • The instructions should contain the exact command and environment needed to run to
925 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 926 • The authors should provide instructions on data access and preparation, including how
927 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 928 • The authors should provide scripts to reproduce all experimental results for the new
929 proposed method and baselines. If only a subset of experiments are reproducible, they
930 should state which ones are omitted from the script and why.
- 931 • At submission time, to preserve anonymity, the authors should release anonymized
932 versions (if applicable).
- 933 • Providing as much information as possible in supplemental material (appended to the
934 paper) is recommended, but including URLs to data and code is permitted.
- 935

936 **6. Experimental setting/details**

937 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
938 rameters, how they were chosen, type of optimizer) necessary to understand the results?

939 Answer: [Yes]

940 Justification:

941 Guidelines:

- 942 • The answer [N/A] means that the paper does not include experiments.
- 943 • The experimental setting should be presented in the core of the paper to a level of detail
944 that is necessary to appreciate the results and make sense of them.
- 945 • The full details can be provided either with the code, in appendix, or as supplemental
946 material.

947 **7. Experiment statistical significance**

948 Question: Does the paper report error bars suitably and correctly defined or other appropriate
949 information about the statistical significance of the experiments?

950 Answer: [No]

951 Justification:

952 Guidelines:

- 953 • The answer [N/A] means that the paper does not include experiments.
- 954 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
955 intervals, or statistical significance tests, at least for the experiments that support the
956 main claims of the paper.
- 957 • The factors of variability that the error bars are capturing should be clearly stated (for
958 example, train/test split, initialization, random drawing of some parameter, or overall
959 run with given experimental conditions).

- 960 • The method for calculating the error bars should be explained (closed form formula,
961 call to a library function, bootstrap, etc.)
- 962 • The assumptions made should be given (e.g., Normally distributed errors).
- 963 • It should be clear whether the error bar is the standard deviation or the standard error
964 of the mean.
- 965 • It is OK to report 1-sigma error bars, but one should state it. The authors should
966 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
967 of Normality of errors is not verified.
- 968 • For asymmetric distributions, the authors should be careful not to show in tables or
969 figures symmetric error bars that would yield results that are out of range (e.g., negative
970 error rates).
- 971 • If error bars are reported in tables or plots, the authors should explain in the text how
972 they were calculated and reference the corresponding figures or tables in the text.

973 8. Experiments compute resources

974 Question: For each experiment, does the paper provide sufficient information on the com-
975 puter resources (type of compute workers, memory, time of execution) needed to reproduce
976 the experiments?

977 Answer: [No]

978 Justification:

979 Guidelines:

- 980 • The answer [N/A] means that the paper does not include experiments.
- 981 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
982 or cloud provider, including relevant memory and storage.
- 983 • The paper should provide the amount of compute required for each of the individual
984 experimental runs as well as estimate the total compute.
- 985 • The paper should disclose whether the full research project required more compute
986 than the experiments reported in the paper (e.g., preliminary or failed experiments that
987 didn't make it into the paper).

988 9. Code of ethics

989 Question: Does the research conducted in the paper conform, in every respect, with the
990 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

991 Answer: [Yes]

992 Justification:

993 Guidelines:

- 994 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
995 Ethics.
- 996 • If the authors answer [No], they should explain the special circumstances that require a
997 deviation from the Code of Ethics.
- 998 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
999 eration due to laws or regulations in their jurisdiction).

1000 10. Broader impacts

1001 Question: Does the paper discuss both potential positive societal impacts and negative
1002 societal impacts of the work performed?

1003 Answer: [Yes]

1004 Justification: We describe it that our data can help in physics editable generation.

1005 Guidelines:

- 1006 • The answer [N/A] means that there is no societal impact of the work performed.
- 1007 • If the authors answer [N/A] or [No], they should explain why their work has no societal
1008 impact or why the paper does not address societal impact.

- 1009 • Examples of negative societal impacts include potential malicious or unintended uses
1010 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1011 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1012 groups), privacy considerations, and security considerations.
- 1013 • The conference expects that many papers will be foundational research and not tied
1014 to particular applications, let alone deployments. However, if there is a direct path to
1015 any negative applications, the authors should point it out. For example, it is legitimate
1016 to point out that an improvement in the quality of generative models could be used to
1017 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
1018 that a generic algorithm for optimizing neural networks could enable people to train
1019 models that generate Deepfakes faster.
- 1020 • The authors should consider possible harms that could arise when the technology is
1021 being used as intended and functioning correctly, harms that could arise when the
1022 technology is being used as intended but gives incorrect results, and harms following
1023 from (intentional or unintentional) misuse of the technology.
- 1024 • If there are negative societal impacts, the authors could also discuss possible mitigation
1025 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1026 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1027 feedback over time, improving the efficiency and accessibility of ML).

1028 11. Safeguards

1029 Question: Does the paper describe safeguards that have been put in place for responsible
1030 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1031 image generators, or scraped datasets)?

1032 Answer: [N/A]

1033 Justification:

1034 Guidelines:

- 1035 • The answer [N/A] means that the paper poses no such risks.
- 1036 • Released models that have a high risk for misuse or dual-use should be released with
1037 necessary safeguards to allow for controlled use of the model, for example by requiring
1038 that users adhere to usage guidelines or restrictions to access the model or implementing
1039 safety filters.
- 1040 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1041 should describe how they avoided releasing unsafe images.
- 1042 • We recognize that providing effective safeguards is challenging, and many papers do
1043 not require this, but we encourage authors to take this into account and make a best
1044 faith effort.

1045 12. Licenses for existing assets

1046 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1047 the paper, properly credited and are the license and terms of use explicitly mentioned and
1048 properly respected?

1049 Answer: [Yes]

1050 Justification:

1051 Guidelines:

- 1052 • The answer [N/A] means that the paper does not use existing assets.
- 1053 • The authors should cite the original paper that produced the code package or dataset.
- 1054 • The authors should state which version of the asset is used and, if possible, include a
1055 URL.
- 1056 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1057 • For scraped data from a particular source (e.g., website), the copyright and terms of
1058 service of that source should be provided.
- 1059 • If assets are released, the license, copyright information, and terms of use in the
1060 package should be provided. For popular datasets, paperswithcode.com/datasets
1061 has curated licenses for some datasets. Their licensing guide can help determine the
1062 license of a dataset.

- 1063
- 1064
- 1065
- 1066
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

1067 **13. New assets**

1068 Question: Are new assets introduced in the paper well documented and is the documentation
1069 provided alongside the assets?

1070 Answer: [Yes]

1071 Justification:

1072 Guidelines:

- 1073
- 1074
- 1075
- 1076
- 1077
- 1078
- 1079
- 1080
- The answer [N/A] means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

1081 **14. Crowdsourcing and research with human subjects**

1082 Question: For crowdsourcing experiments and research with human subjects, does the paper
1083 include the full text of instructions given to participants and screenshots, if applicable, as
1084 well as details about compensation (if any)?

1085 Answer: [N/A]

1086 Justification:

1087 Guidelines:

- 1088
- 1089
- 1090
- 1091
- 1092
- 1093
- 1094
- 1095
- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

1096 **15. Institutional review board (IRB) approvals or equivalent for research with human
1097 subjects**

1098 Question: Does the paper describe potential risks incurred by study participants, whether
1099 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1100 approvals (or an equivalent approval/review based on the requirements of your country or
1101 institution) were obtained?

1102 Answer: [N/A]

1103 Justification:

1104 Guidelines:

- 1105
- 1106
- 1107
- 1108
- 1109
- 1110
- 1111
- 1112
- 1113
- 1114
- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
 - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
 - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [\[Yes\]](#)

Justification: We use llm to help the draft.

Guidelines:

- The answer [\[N/A\]](#) means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.