

GameCraft-2: Instruction-following Interactive Game World Model

Anonymous CVPR submission

Paper ID 2



Figure 1. *GameCraft-2* advances generative game world models from static scene synthesis to open-ended, instruction-following simulation. From a single image, we simulate action signals; left/right frames show key moments under different inputs. *GameCraft-2* generates interaction-aligned videos with temporal/3D consistency. Instructions (e.g., “draw a torch”, “hold a gun and fire”) and key/mouse actions control camera motion and content updates, producing coherent videos with consistent style. W/A/S/D indicate movement, \uparrow , \leftarrow , \downarrow , \rightarrow adjust view angles, and white denotes idle.

Abstract

001 Recent advances in generative world models have enabled
 002 remarkable progress in creating open-ended game environ-
 003 nments, evolving from static scene synthesis toward dynamic,
 004 interactive simulation. However, current approaches remain
 005 limited by rigid action schemas and high annotation costs,
 006 restricting their ability to model diverse in-game interactions
 007 and player-driven dynamics. To address these challenges,
 008 we introduce *GameCraft-2*, a new paradigm of instruction-
 009 driven interaction for generative game world modeling. In-
 010 stead of relying on fixed keyboard inputs, our model allows
 011 users to control game video contents through natural lan-
 012 guage prompts, keyboard, or mouse signals, enabling flexible
 013 and semantically rich interaction within generated worlds.
 014 We formally define the concept of Interactive Video Data
 015 and develop an automated pipeline that converts large-scale,
 016 unstructured text–video pairs into causally aligned interac-
 017 tive datasets. Built upon a 14B image-to-video Mixture-of-
 018 Experts (MoE) foundation model, our model incorporates a
 019 text-driven interaction injection mechanism for fine-grained
 020 control over camera motion, character behavior, and environ-
 021 ment dynamics. We introduce an interaction-focused bench-

mark, *InterBench* to evaluate interaction performance com- 022
 prehensively. Extensive experiments demonstrate that our 023
 model generates temporally coherent, and causally grounded 024
 interactive game videos that faithfully respond to diverse and 025
 free-form user instructions such as “open the door”, “draw 026
 a torch”, or “trigger an explosion”. 027

1. Introduction 028

The rapid progress of diffusion models [18, 27, 31, 41, 44, 029
 51] has accelerated dynamic game content creation [17, 28, 030
 64]. Beyond static images and short videos, recent systems 031
 such as RTFM [55] and the Genie series [1] demonstrate that 032
 world models can underpin immersive, controllable virtual 033
 experiences, advancing AI-driven “playable worlds” that 034
 simulate and respond to user intent. 035

Existing world models can be categorized into **3D-** 036
based and **video-based** approaches. 3D-based world mod- 037
 els [22, 30, 37, 47, 54, 55] emphasize geometric consistency 038
 and physical accuracy, enabling detailed world reconstruc- 039
 tion and memory persistence. However, they are often lim- 040
 ited to scripted or static interactions, lacking the creative 041
 flexibility and open-ended gameplay dynamics essential for 042
 interactive game environments. With recent improvements 043

| | | |
|-----|--|-----|
| 044 | in video foundation models [3, 27, 51, 53], the video-based | 097 |
| 045 | technical pathway [1, 11, 17, 28, 40, 56, 61, 64] has shown | 098 |
| 046 | remarkable potential. These works learn world dynamics | 099 |
| 047 | directly from large-scale video data [6, 11, 32, 66] through | 100 |
| 048 | implicit end-to-end representation learning. Notably, the | 101 |
| 049 | Genie series [1, 40] introduces latent action modeling to | 102 |
| 050 | simulate player-driven interactions, while Matrix-Game [64] | 103 |
| 051 | and Hunyuan-GameCraft [28] integrate discrete actions (e.g., | 104 |
| 052 | W/A/S/D, mouse movements) into a unified representation | 105 |
| 053 | space, achieving continuous, high-fidelity video generation | 106 |
| 054 | that responds to user inputs. | 107 |
| 055 | These frontier works mark a fundamental shift in focus | |
| 056 | from the world’s static appearance (“what the world looks | |
| 057 | like”) to its interactive dynamics (“how we interact with it”), | |
| 058 | necessitating a rigorous definition of “interaction” within | |
| 059 | world models, especially in game scenarios. | |
| 060 | We formally define interaction in world models as “ <i>actions</i> | |
| 061 | <i>executed by an explicit agent that trigger state transi-</i> | |
| 062 | <i>tions in the environment with clear causal relationships</i> | |
| 063 | <i>and physical or logical validity.”</i> This definition encom- | |
| 064 | passes diverse input modalities, from mouse and keyboard | |
| 065 | operations [11, 28, 56, 61, 64] to embodied motion sens- | |
| 066 | ing [38]. Grounded in this perspective, two key challenges | |
| 067 | hinder this progress: (1) the lack of a formal definition and a | |
| 068 | scalable construction pipeline for interactive video data, and | |
| 069 | (2) multi-turn interactions in long video generation while | |
| 070 | maintaining video quality and interaction accuracy. | |
| 071 | To address these challenges, we present GameCraft-2 , | |
| 072 | an interactive game world model for free-form instruction- | |
| 073 | following control. We begin by formally defining interaction | |
| 074 | within the context of generative world models, and develop | |
| 075 | two automated pipelines for interactive video data construc- | |
| 076 | tion and refinement. These pipelines, for the first time, en- | |
| 077 | able the efficient transformation of large-scale, unstructured | |
| 078 | text–video pairs into open-domain interactive datasets en- | |
| 079 | riched with implicit causal labels. | |
| 080 | For model training, our model integrates text-based in- | |
| 081 | structions and keyboard/mouse action signals into a unified | |
| 082 | controllable video generator, enabling flexible, semantically | |
| 083 | grounded, and causally consistent interaction within dy- | |
| 084 | namic game environments. To support efficient long-horizon | |
| 085 | video generation, we employ a comprehensive autoregres- | |
| 086 | sive distillation strategy that transfers the bidirectional video | |
| 087 | generator into a causal autoregressive model. Subsequently, | |
| 088 | a randomized image-to-long-video extension tuning scheme | |
| 089 | is introduced to alleviate error accumulation during extended | |
| 090 | rollouts, ensuring stable and coherent long-form genera- | |
| 091 | tion. For multi-turn interactive inference, we following | |
| 092 | LongLive [57] to employ a KV-recache mechanism to en- | |
| 093 | hance the accuracy and stability of multi-turn interactions in | |
| 094 | autoregressive long video generation. In addition, we incor- | |
| 095 | porate several engineering acceleration optimizations, boost- | |
| 096 | ing the model’s inference speed to 16 FPS on 8 NVIDIA | |
| | H2O GPUs, enabling real-time interactive video generation. | |
| | To comprehensively evaluate interactive performance | |
| | across models, we introduce InterBench, a benchmark that | |
| | measures key dimensions of interactive behavior: interaction | |
| | completeness, action effectiveness, causal coherence, and | |
| | physical plausibility. Extensive experiments on InterBench | |
| | demonstrate the effectiveness of our framework, achieving | |
| | SOTA performance in generating interactive videos that re- | |
| | spond faithfully to user instructions while maintaining high | |
| | visual fidelity and temporal coherence. | |
| | In general, our main contributions are as follows: | |
| | • We propose a Unified Interaction Framework and scal- | |
| | able data engine that unifies keyboard/mouse signals with | |
| | real-time text instructions . We then post-train Wan2.2 | |
| | 14B MoE on automatically generated/annotated/cleaned | |
| | data for strong controllability and emergent generaliza- | |
| | tion to unseen action combinations. | |
| | • We leverage autoregressive distillation and randomized | |
| | long-video tuning for efficient and stable long-horizon | |
| | generation, and introduce KV-recache for multi-turn in- | |
| | ference with real-time 16 FPS performance. | |
| | • We establish InterBench , a benchmark that combines | |
| | VLM-based evaluation with classical video metrics to | |
| | quantify both visual fidelity and instruction-following | |
| | accuracy in dynamic interactions. | |
| | 2. Related Works | |
| | 2.1. Long Video Extension | |
| | Long video generation primarily struggles with “train-short- | |
| | test-long” discrepancies, leading to semantic drift and ar- | |
| | tifacts. To mitigate this, Self-Forcing [23] and its fol- | |
| | lower [9, 35, 39, 57] align training with inference via self- | |
| | conditioning and rolling windows. Alternative paradigms | |
| | explore next-frame prediction [12, 13], hybrid diffusion- | |
| | autoregressive models like DiffusionForcing [7], or test-time | |
| | adaptation [10]. Recently, LongLive [57] introduced KV- | |
| | recache to improve interactive semantic control responsive- | |
| | ness. | |
| | 2.2. Interactive Video-based World Model | |
| | Unlike static video models, interactive world models respond | |
| | dynamically to inputs for explorable environments. Early | |
| | game-specific models like MineWorld [15], Matrix [11], | |
| | and GameFactory [61] focused on discrete actions (e.g., | |
| | keyboard/mouse). To enhance generalization and consis- | |
| | tency, Genie2 [40] proposed a 2D foundation model, while | |
| | WorldMem [56] utilized memory banks for long-term simu- | |
| | lation. Recent works such as GameGen-X [6], Genie3 [1], | |
| | and Hunyuan-GameCraft [28] unify control signals into | |
| | shared continuous spaces. However, these models mainly | |
| | use prompts for initialization rather than real-time interac- | |
| | tion, which is still limited by physical inputs. | |

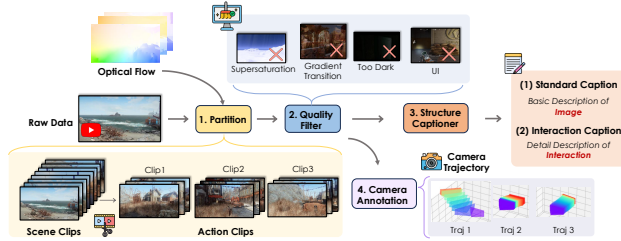


Figure 2. **Data Curation Pipeline.** (1) **Partition:** segment raw videos into scene/action clips via scene detection and optical flow; (2) **Quality Filtering:** discard low-quality frames using visual, luminance checks; (3) **Structured Captioning:** generate standard and interaction captions per clip; (4) **Camera Annotation:** recover 6-DoF camera trajectories to model viewpoint motion.

147 2.3. Text-guided Video Generation and Editing

148 Text-driven synthesis follows two paradigms: semantic enrichment and structured execution. The former enhances
149 prompts via LLM representations or adapters [36, 46, 63, 65]. The latter treats text as a script, using LLMs to decompose
150 prompts into frame-level descriptions or layouts [19, 20, 33, 34]. Related editing works leverage video-to-video frame-
151 works for zero-shot or end-to-end control [8, 25, 29, 42, 43]. Unlike these one-off, non-interactive methods, GameCraft-
152 2 enables continuous interaction where prompts iteratively evolve a dynamic world state through feedback loops.
153
154
155
156
157

158 3. Interactive Video Data Construction

159 3.1. Definition of Interactive Video Data

160 Unlike traditional video, **Interactive Video Data** is defined as a temporal sequence recording a *causally driven state-*
161 *transition process*. It captures how agents or environments move from a defined **initial state** to a distinct **final state**.
162 We categorize these interactions into: (1) **Environmental Interactions** (scene changes), (2) **Actor Actions** (embodied
163 agent-driven), and (3) **Entity Appearances** (subject introduction). Each includes varying complexity levels.
164
165
166
167

168 3.2. Controllable Synthetic Data Production

169 To address data scarcity, we develop a synthetic pipeline leveraging foundation models to generate high-quality inter-
170 active transitions. Given an initial frame F_t (which can be synthesized via T2I models to ensure diversity), we employ
171 two strategies based on the interaction type:
172
173

- 174 • **Start-End Frame Strategy:** For stationary scenes (e.g., weather changes), a VLM-guided editing model gener-
175 ates a target end-frame F'_t , ensuring strong controllability over the final state.
- 176 • **First-Frame-Driven Strategy:** For dynamic actions (e.g., opening a door), the model generates sequences
177 from F_t to ensure temporal continuity.
- 178
179
180

3.3. Game Scene Curation and Annotation

For raw videos, we implement a robust curation pipeline:

- **Action-aware Partition & Filtering:** We split videos into action-aware clips and perform a three-stage qual-
ity screen (fidelity, luminance, semantic consistency) to discard noisy or inconsistent clips.
- **Camera Annotation:** We reconstruct 6-DoF camera trajectories for each clip to provide precise translational
and rotational metadata for viewpoint modeling.
- **Structured Dual-Captioning:** To supervise both states and transitions, we generate a *Standard Caption* (C_t) for
static visual layout and an *Interaction Caption* ($I_{t \rightarrow t+1}$) capturing the semantic delta between clips. This enables
the model to jointly learn appearance perception and action-level reasoning.

4. Method

We present **GameCraft-2**, an interactive game video model for free-form instruction-based control (Fig. 3). It unifies
an action-injected causal architecture, image-conditioned autoregressive long-video generation, and multi-prompt in-
teraction within a single framework. We then describe the architecture, training, and inference procedures.

4.1. Model Architecture

Our model builds on a 14B image-to-video MoE diffusion foundation model [51], which we extend into an action-
controllable generator. As in Sec. 1, actions comprise both keyboard inputs and free-form text prompts.

For keyboard/mouse control (e.g., W/A/S/D, arrows), we map discrete inputs to continuous camera parameters fol-
lowing Hunyuan-GameCraft [28]. We encode annotated parameters as Plücker embeddings [16] and inject them via
token addition during training; at inference, user inputs are converted into camera trajectories to derive the parameters.

For prompt-based interaction injection, we find the base model underperforms on interactive verbs due to their higher
semantic/spatial complexity and tight grounding to specific regions or instances. We therefore use an MLLM [52] to ex-
tract, reason about, and inject interaction cues, strengthening fine-grained guidance and helping distinguish scene-level
instructions from interactive behaviors during training. Combined with camera-conditioned control, this yields a unified
mechanism for coherent navigation and interaction.

4.2. Training Procedure

To enable long-horizon, real-time interactive generation, we distill a bidirectional foundation model into a few-step causal
generator. Specifically, we scale Self-Forcing [23] to a 14B MoE image-to-video model [51], improving quality and
efficiency under rapid scene changes. To reduce error accumulation, we propose random extension tuning. Training

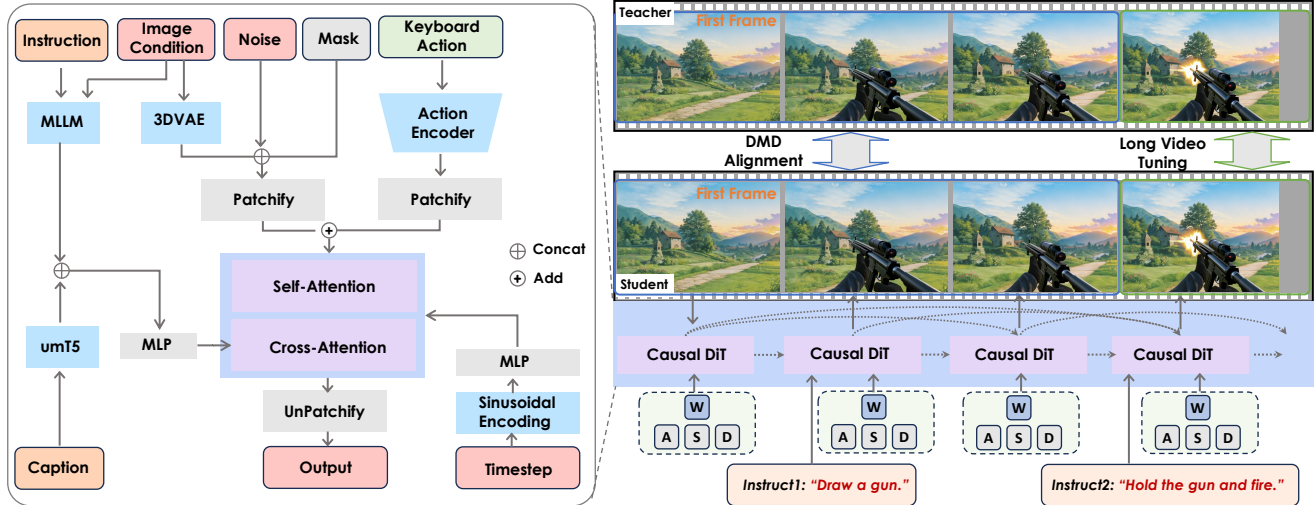


Figure 3. Model architecture of GameCraft-2. Given a reference image and the corresponding action, the keyboard/mouse signal, and prompt-based instruction, we inject these options to the main architecture. During training and inference, we leverage self-forcing post-training for long-video extension, and KV-cache/recache for multi-action switching. To maintain the long-term video quality, we design a randomized long video tuning scheme.

230 proceeds in four stages: (1) Action-Injected Training, (2)
 231 Instruction-Oriented SFT, (3) Autoregressive Distillation,
 232 and (4) Randomized Long-Video Extension Tuning.

233 4.2.1. Action-Injected Training

234 The primary objective of this stage is to establish a funda-
 235 mental understanding of 3D scene dynamics, lighting, and
 236 physics. We load the pre-trained weights and finetune the
 237 model with the flow-matching objective for architectural
 238 adaptation. In order to improve the long-term consistency,
 239 we adopted a curriculum learning strategy. Specifically, we
 240 organized the training into three phases, exposing the model
 241 to video data of 45, 81, and 149 frames in 480p in sequence.
 242 This stepped approach allows the model to first solidify its
 243 understanding of short-term motion dynamics before gradu-
 244 ally adapting its attention mechanisms to handle the complex
 245 dependencies required for longer-duration coherence. Be-
 246 sides, we randomly choose long and short captions during
 247 training, and concatenate interactive captions for interaction
 248 learning. This option will help the model to have an initial
 249 perception of the injection of interactive information.

250 4.2.2. Instruction-Oriented Supervised Fine-Tuning

251 To enhance interactivity, we built a 150K-sample dataset
 252 by augmenting real-world footage with procedurally gener-
 253 ated synthetic videos. These sequences provide high-fidelity
 254 supervision across diverse interactions, establishing a tight
 255 correspondence between actions and visual outcomes. Next,
 256 we freeze the camera encoder and fine-tune only the MoE
 257 experts to improve alignment with semantic control cues.

4.2.3. Autoregressive Generator Distillation

258 For interactive world models, extending fixed-length video
 259 generators to high-quality autoregressive long-video genera-
 260 tion is essential. Prior works have made preliminary attempts
 261 on long video generation [23, 45, 50, 57]. Building upon the
 262 high- and low-noise MoE architecture and camera parameter
 263 injection, we introduce targeted adaptations to the attention
 264 mechanism and the distillation protocol. These modifica-
 265 tions are specifically tailored to optimize performance within
 266 the autoregressive distillation process.

Sink Token and Block Sparse Attention: Sliding-window
 268 KV updates in prior causal attention [23, 60] can cause drift
 269 as later tokens lose access to the initial conditioning frame.
 270 Following [45, 50, 57], we keep the first frame as a *sink*
 271 *token* in the KV cache, improving long-horizon stability and
 272 preserving the coordinate origin so injected camera param-
 273 eters remain aligned, avoiding per-step recaching. We further
 274 use Block Sparse Attention [14] for local attention, where
 275 each target block attends to selected preceding blocks. Sink
 276 + block-sparse local attention together constitute the KV
 277 cache, boosting quality and speed for block-wise generation.

Distillation Schedule: Since the high-noise MoE expert is
 279 harder to train [51], especially under SFT/distillation, we use
 280 expert-specific learning rates. We also set timestep targets by
 281 the inter-expert noise boundary, ensuring teacher and student
 282 choose the same expert throughout distillation.

4.2.4. Randomized Extended Long-Video Tuning

284 Our approach to enabling long-form video generation is
 285 motivated by the observation that the foundation model, de-
 286 spite being pre-trained on short clips, implicitly captures the
 287 global visual data distribution. Previous methods [9, 57], roll
 288

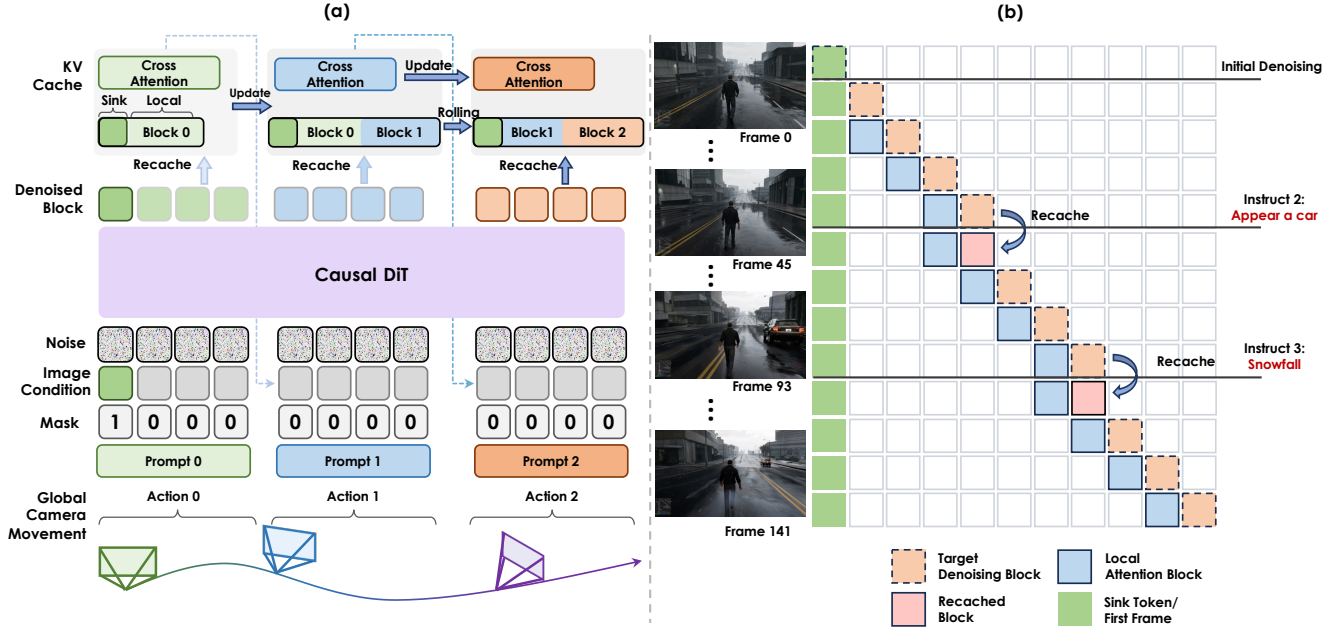


Figure 4. **Multi-turn Interactive Inference.** (a) shows block-wise autoregressive long-video inference with KV-cache updates: the initial frame is kept as sink tokens, while local attention uses recent blocks. (b) illustrates prompt re-caching, which improves interaction responsiveness and accuracy for new prompts.

289 out long video sequences from a causal generator and apply
290 distributional moment distance (DMD) [58, 59] alignment
291 on the extended frames. This strategy effectively mitigates
292 error accumulation during autoregressive generation.

293 Building upon this insight, we adopt a randomized exten-
294 sion tuning strategy using a dataset of long-form gameplay
295 videos exceeding 10 seconds. In this stage, the model au-
296 toregressively rolls out N frames, and contiguous T -frame
297 windows are uniformly sampled to align the predicted and
298 target distributions. Furthermore, we randomly extend the
299 predicted videos from the causal generator to varying lengths,
300 promoting robustness across different temporal horizons. In
301 practice, while rolling out at window $W = V[i : i + K - 1]$,
302 the student generator uses sink token and KV cache and
303 autoregressively extend long video, and the fake score teacher
304 model uses the last frame in the previous clean predicted
305 chunk $V[i - 1]$ as image condition; while the real score uses
306 the ground truth frame in the original video.

307 To mitigate the potential erosion of interactive capabil-
308 ities inherent in few-step distillation, we adopt a training
309 paradigm that interleaves self-forcing with teacher-forcing.
310 The rationale for this approach is to compel the model to
311 master state recovery and maintain temporal stability. Cruci-
312 ally, this is achieved by exposing it to diverse states at
313 arbitrary points along the generation trajectory, rather than
314 limiting such corrective training solely to the initial phase.

315 4.3. Multi-turn Interactive Inference

316 **Self-attention KV Cache.** We use a fixed-length KV cache
317 with rolling updates for efficient autoregressive generation.

Sink tokens are permanently retained at the cache start, while
318 a local window keeps the N frames preceding the target
319 denoising block across multi-turn interactions. The cache
320 comprises sink tokens and a block-sparse local attention com-
321 ponent, improving efficiency and preventing quality drift.
322

ReCache Mechanism. We employ a recache mechanism to
323 enhance the accuracy and stability of multi-turn interactions
324 in autoregressive long video generation. Upon receiving a
325 new interaction prompt, the model extracts the corresponding
326 interaction embeddings to recompute the last autoregressive
327 block and update both the self-attention and cross-attention
328 KV caches. This strategy provides precise historical context
329 for the subsequent target block with minimal computational
330 overhead, thereby ensuring accurate and responsive feedback
331 to facilitate a smoother user experience.
332

333 5. Experiments

334 5.1. Model and Dataset Configurations.

335 We compare our method against several SOTA image-to-
336 video generation foundation models, including **Hunyuan-**
337 **Video**, **Wan2.2 A14B**, and **LongCatVideo**. To compre-
338 hensively evaluate controllable generation, we generate 93-
339 frame videos at 832×448 resolution using a curated test set
340 spanning diverse scenes and specific action scenarios.

341 5.2. InterBench: A Comprehensive Benchmark for 342 Interactive Video Generation

343 To evaluate video generation, we adopt two metric families:
344 general video-quality metrics and our interaction-focused

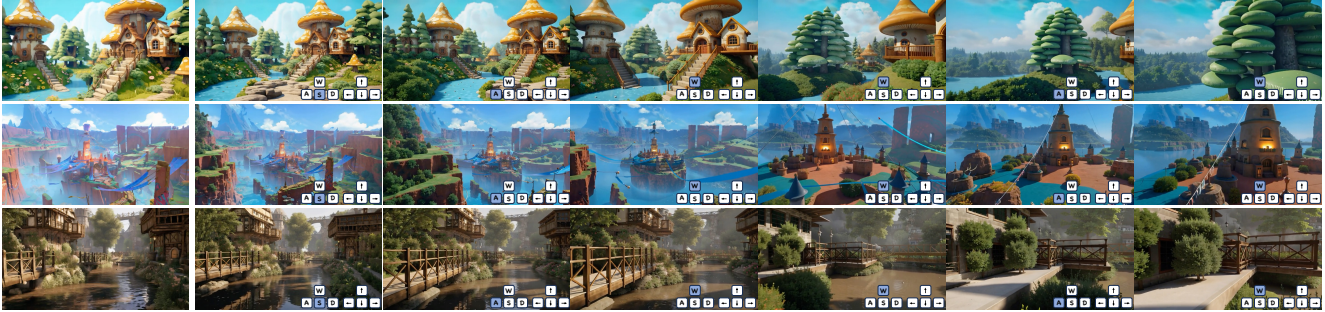


Figure 5. **Inference results by GameCraft-2 on multi-action control.** In our case, blue-lit keys indicate key presses. W, A, S, D represent transition movement and \uparrow , \leftarrow , \downarrow , \rightarrow denote changes in view angles.



Figure 6. **Inference results by GameCraft-2 on instruction-based multi-action control.** Except W, A, S, D, \uparrow , \leftarrow , \downarrow , and \rightarrow , the red notes denote action instructions.



Figure 7. **Inference results by GameCraft-2 on the third-perspective long-term game video generation.**

345 suite, **InterBench**. While general metrics capture fidelity,
 346 temporal consistency, and motion realism, they miss key
 347 interactive properties such as causality, action execution,
 348 and state transitions. InterBench addresses this with six
 349 interaction-centric dimensions (e.g., interaction complete-
 350 ness, action effectiveness, causal coherence), yielding a holistic
 351 evaluation of interactive video models.

General Quality Metrics To comprehensively assess our 352
 model, we adopt diverse metrics. For **video realism**, we 353
 use Fréchet Video Distance (FVD) [49] to capture spatial 354
 fidelity and temporal dynamics. **Visual quality** is measured 355
 by Image Quality and Aesthetic scores, reflecting percep- 356
 tual clarity and visual appeal. We evaluate temporal consis- 357
 tency for cross-frame coherence and artifacts (e.g., flicker- 358
 ing, structural instability). For **dynamic performance**, we 359
 adapt VBench’s Dynamic Degree [24] by reporting abso- 360
 lute optical-flow magnitudes (*Dynamic Average*) instead of 361
 binary motion classes, yielding a finer measure of motion 362
 intensity and naturalness. 363

Geometric Camera Control. For interactive camera control, 364
 we report Relative Pose Error (RPE) in translation and 365
 rotation after Sim3 Umeyama alignment between predicted 366
 and ground-truth trajectories. This removes global pose and 367
 scale ambiguity, so RPE reflects local frame-to-frame motion 368

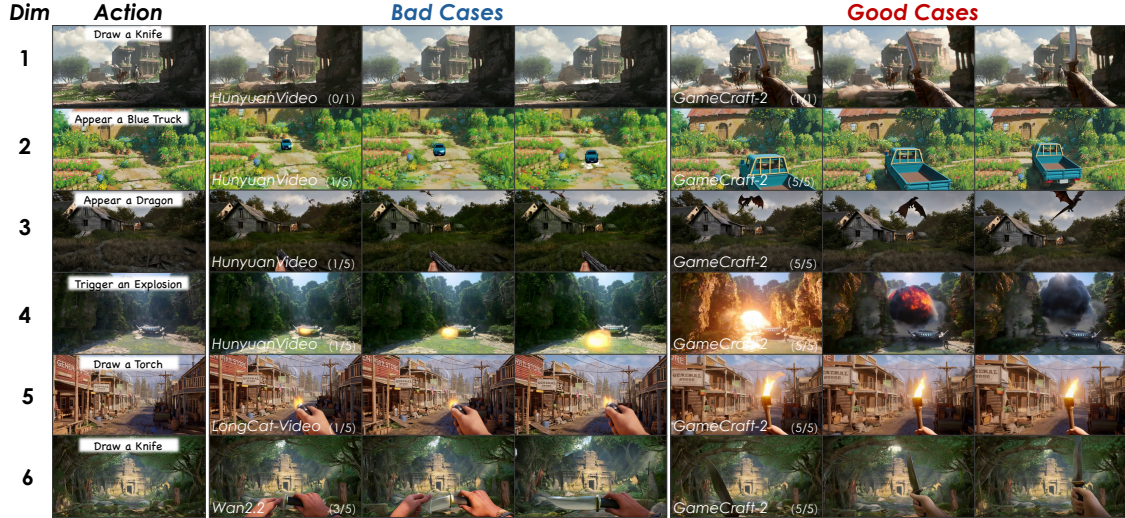


Figure 8. **Qualitative examples for the six InterBench dimensions.** Rows correspond to (1) **Trigger Rate**, (2) **Prompt-Video Alignment**, (3) **Fluency**, (4) **Scope Accuracy**, (5) **End-State Consistency**, and (6) **Object Physics**. For each, we compare a high-scoring result from GameCraft-2 (right) with a low-scoring baseline failure (left), illustrating the rating scale.

Table 1. **Quantitative comparison of recent controllable video generation models.**

| Model | Visual Quality | | | | Temporal Consistency \uparrow | RPE | |
|--------------------|------------------|--------------------------|--------------------|----------------------|---------------------------------|--------------------|------------------|
| | FVD \downarrow | Image Quality \uparrow | Dynamic \uparrow | Aesthetic \uparrow | | Trans \downarrow | Rot \downarrow |
| GameCraft | 1554.2 | 0.69 | 67.2 | 0.67 | 0.95 | 0.08 | 0.20 |
| GameCraft-PCM | 1883.3 | 0.67 | 43.8 | 0.65 | 0.93 | 0.08 | 0.20 |
| Matrix-Game | 2260.7 | 0.72 | 31.7 | 0.65 | 0.94 | 0.18 | 0.35 |
| Matrix-Game-2.0 | 1920.6 | 0.62 | 20.5 | 0.49 | 0.84 | 0.08 | 0.25 |
| GameCraft-2 | 1856.3 | 0.70 | 45.2 | 0.71 | 0.96 | 0.08 | 0.17 |

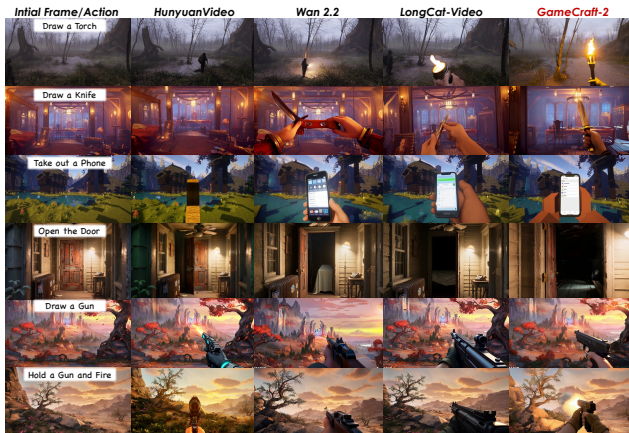


Figure 9. **Comparison of Actor-Action Interactions with Baseline Models.** Visual comparisons illustrating the quality of action-level interactions across representative prompts. Our method produces more coherent and physically consistent actions than all baselines.

Action-Level Interaction Semantics. To evaluate *action-level interaction*, we introduce **InterBench**, a six-dimensional protocol for interactive video generation. Using a Vision-Language Model (VLM) as an automatic evaluator, it measures interaction fidelity, smoothness, and physical plausibility. The six dimensions are:

1. **Interaction Trigger Rate.** A binary check of whether the requested interaction is initiated, distinguishing ignored prompts from attempted actions.
2. **Prompt-Video Alignment.** Measures semantic fidelity to the full prompt, including *static alignment* (scene/context consistency) and *dynamic alignment* (correct action execution).
3. **Interaction Fluency.** Assesses temporal naturalness and visual coherence, penalizing artifacts such as jumps, flicker, and object teleportation.
4. **Interaction Scope Accuracy.** Evaluates whether interaction effects have appropriate spatial extent (global events affect the whole scene; local actions remain localized).
5. **End-State Consistency.** Checks whether the interaction reaches a correct, stable end state.
6. **Object Physics Correctness.** Evaluates physical plausibility, including rigid-body integrity, realistic kinematics,

369 fidelity and control precision under interactive inputs.

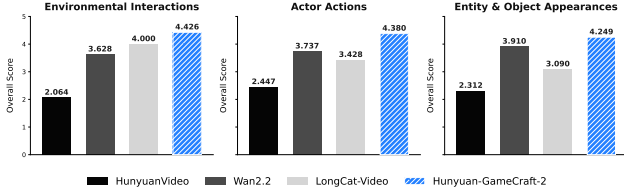


Figure 10. **Comparison of Environmental Interactions with Baseline Models.** Qualitative results showing the fidelity and consistency of environment-level effects. Our approach better preserves global influence and temporal stability.

and correct contact relationships (no penetration).

Hierarchical scoring combines a binary **Trigger** with five ordinal dimensions into a weighted *Overall* score for prioritization.

5.3. Interaction Evaluation

Quantitative Results on InterBench InterBench results show that GameCraft-2 best balances complex interaction logic with cinematic control. It achieves near-perfect semantic responsiveness with $\text{Trigger} > 0.96$ for both Environmental and Actor actions (Figure 10) outperforming Wan2.2 and HunyuanVideo, while also improving physical fidelity ($\text{Physics} +0.68$) and EndState stability. For detailed subclass scores, please see the appendix. Importantly, this does not sacrifice visual quality. As shown in Tab 1. GameCraft-2 attains the highest **Temporal Consistency** (0.96) and **Aesthetic** (0.71) among control-centric models. With mixed real and synthetic data, it further delivers precise camera control, yielding the lowest rotation error (**Rot** 0.17) and competitive translation accuracy. Overall, these metrics confirm that GameCraft-2 unifies high-level physical reasoning with high-fidelity generation.

Qualitative Analysis The results clearly highlight the superior performance of GameCraft-2 over baseline models. Baselines frequently exhibit noticeable deficiencies when handling complex interactions. For instance, environmental effects often lack dynamic evolution and realistic lighting interactions. Actor actions are commonly plagued by object deformation, motion incoherence, and inaccurate hand-object contact. Furthermore, newly generated entities tend to suffer from identity drift, unstable geometry, and poor integration with the scene. In contrast, GameCraft-2 demonstrates substantially higher fidelity and consistency across all interaction categories. For Actor Actions, GameCraft-2 produces more coherent action sequences, enabling characters to stably grasp and precisely manipulate objects while ensuring stable final states. These qualitative examples support the quantitative results and demonstrate GameCraft-2’s robust generation of semantically accurate, temporally coherent, physically plausible interaction videos.

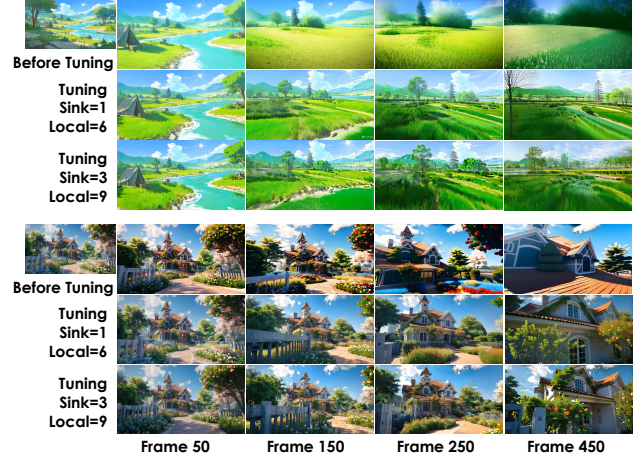


Figure 11. **Qualitative Analysis of Long-Video Tuning and Cache Settings.** Row 1: Baseline results without Long-Video Tuning (sink token size = 1, local attention size = 6). Row 2: Incorporates Long-Video Tuning upon the baseline. Row 3: Further modifies setting based on Row 2 by increasing the sink token size to 3 and local attention size to 9. Input prompts and camera parameters remain consistent across all samples.

Analysis of Long-Video Tuning and Cache Settings We qualitatively analyze the impact of long-video tuning and KV cache settings, specifically regarding the sink tokens and local attention. We compare generated frames at aligned time steps, the integration of randomized extended long-video tuning substantially improves video fidelity and motion consistency beyond the 450th frame. Moreover, expanding the sink tokens and local attention size can enrich the detail but increases the artifacts. These observations confirm both the efficacy of our tuning strategy and the importance of leveraging sink tokens and local attention to maintain robust context.

6. Conclusion

In this work, we present **GameCraft-2**, an interactive game world model that generates high-fidelity, controllable video from free-form text and keyboard/mouse actions. We formalize interactive video data and develop automated curation and synthesis pipelines to alleviate the data bottleneck. GameCraft-2 unifies multimodal controls under a robust training framework, combining randomized long-video tuning with efficient inference (e.g., KV-recache) for stable long-horizon and real-time generation. We further introduce **InterBench**, a benchmark for action-level interaction quality. Experiments show that GameCraft-2 consistently surpasses prior state-of-the-art methods in interaction fidelity, visual quality, and temporal coherence, advancing video synthesis from passive generation to user-driven, playable world creation.

460

References

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

- [1] Philip J. Ball, Jakob Bauer, Frank Belletti, Bethanie Brownfield, Ariel Ephrat, Shlomi Fruchter, Agrim Gupta, Kristian Holsheimer, Aleksander Holynski, Jiri Hron, Christos Kaplanis, Marjorie Limont, Matt McGill, Yanko Oliveira, Jack Parker-Holder, Frank Perbet, Guy Scully, Jeremy Shar, Stephen Spencer, et al. Genie 3: A new frontier for world models. 2025. 1, 2
- [2] Gary Bradski. The opencv library. *Dr. Dobbs's Journal: Software Tools for the Professional Programmer*, 25(11):120–123, 2000. 13
- [3] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. 2
- [4] Siyu Cao, Hangting Chen, Peng Chen, Yiji Cheng, Yutao Cui, Xincheng Deng, Ying Dong, Kipper Gong, Tianpeng Gu, Xiusen Gu, et al. Hunyuanimage 3.0 technical report. *arXiv preprint arXiv:2509.23951*, 2025. 12
- [5] Brandon Castellano. PySceneDetect. 12
- [6] Haoxuan Che, Xuanhua He, Quande Liu, Cheng Jin, and Hao Chen. Gamegen-x: Interactive open-world game video generation. In *International Conference on Learning Representations*, 2025. 2
- [7] Boyuan Chen, Diego Martí Monsó, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in Neural Information Processing Systems*, 37:24081–24125, 2024. 2
- [8] Jiaxin Cheng, Tianjun Xiao, and Tong He. Consistent video-to-video transfer using synthetic dataset, 2023. 3
- [9] Justin Cui, Jie Wu, Ming Li, Tao Yang, Xiaojie Li, Rui Wang, Andrew Bai, Yuanhao Ban, and Cho-Jui Hsieh. Self-forcing++: Towards minute-scale high-quality video generation. *arXiv preprint arXiv:2510.02283*, 2025. 2, 4
- [10] Karan Dalal, Daniel Kocaja, Gashon Hussein, Jiarui Xu, Yue Zhao, Youjin Song, Shihao Han, Ka Chun Cheung, Jan Kautz, Carlos Guestrin, et al. One-minute video generation with test-time training. *arXiv preprint arXiv:2504.05298*, 2025. 2
- [11] Ruili Feng, Han Zhang, Zhantao Yang, Jie Xiao, Zhilei Shu, Zhiheng Liu, Andy Zheng, Yukun Huang, Yu Liu, and Hongyang Zhang. The matrix: Infinite-horizon world generation with real-time moving control. *arXiv preprint arXiv:2412.03568*, 2024. 2
- [12] Kaifeng Gao, Jiabin Shi, Hanwang Zhang, Chunping Wang, and Jun Xiao. Vid-gpt: Introducing gpt-style autoregressive generation in video diffusion models. *arXiv preprint arXiv:2406.10981*, 2024. 2
- [13] Yuchao Gu, Weijia Mao, and Mike Zheng Shou. Long-context autoregressive video modeling with next-frame prediction. *arXiv preprint arXiv:2503.19325*, 2025. 2
- [14] Junxian Guo, Haotian Tang, Shang Yang, Zhekai Zhang, Zhi-jian Liu, and Song Han. Block Sparse Attention. <https://github.com/mit-han-lab/Block-Sparse-Attention>, 2024. 4
- [15] Junliang Guo, Yang Ye, Tianyu He, Haoyu Wu, Yushu Jiang, Tim Pearce, and Jiang Bian. Mineworld: a real-time and open-source interactive world model on minecraft, 2025. 2
- [16] Hao He, Yinghao Xu, Yuwei Guo, Gordon Wetzstein, Bo Dai, Hongsheng Li, and Ceyuan Yang. Cameractrl: Enabling camera control for text-to-video generation. *arXiv preprint arXiv:2404.02101*, 2024. 3
- [17] Xianglong He, Chunli Peng, Zexiang Liu, Boyang Wang, Yifan Zhang, Qi Cui, Fei Kang, Biao Jiang, Mengyin An, Yangyang Ren, et al. Matrix-game 2.0: An open-source real-time and streaming interactive world model. *arXiv preprint arXiv:2508.13009*, 2025. 1, 2
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models, 2020. arXiv:2006.11239 [cs]. 1
- [19] Susung Hong, Junyoung Seo, Heeseong Shin, Sunghwan Hong, and Seungryong Kim. Direct2v: Large language models are frame-level directors for zero-shot text-to-video generation, 2024. 3
- [20] Hanzhuo Huang, Yufan Feng, Cheng Shi, Lan Xu, Jingyi Yu, and Sibe Yang. Free-bloom: Zero-shot text-to-video generator with llm director and ldm animator, 2023. 3
- [21] Jiahui Huang, Qunjie Zhou, Hesam Rabeti, Aleksandr Korovko, Huan Ling, Xuanchi Ren, Tianchang Shen, Jun Gao, Dmitry Slepichev, Chen-Hsuan Lin, Jiawei Ren, Kevin Xie, Joydeep Biswas, Laura Leal-Taixe, and Sanja Fidler. Vipe: Video pose engine for 3d geometric perception. In *NVIDIA Research Whitepapers arXiv:2508.10934*, 2025. 13
- [22] Tianyu Huang, Wangguandong Zheng, Tengfei Wang, Yuhao Liu, Zhenwei Wang, Junta Wu, Jie Jiang, Hui Li, Rynson WH Lau, Wangmeng Zuo, and Chunchao Guo. Voyager: Long-range and world-consistent video diffusion for explorable 3d scene generation. *arXiv preprint arXiv:2506.04225*, 2025. 1
- [23] Xun Huang, Zhengqi Li, Guande He, Mingyuan Zhou, and Eli Shechtman. Self forcing: Bridging the train-test gap in autoregressive video diffusion, 2025. 2, 3, 4
- [24] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 6
- [25] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators, 2023. 3
- [26] KolorsTeam. Kolors: Effective training of diffusion model for photorealistic text-to-image synthesis. *arXiv preprint*, 2024. 13
- [27] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 1, 2
- [28] Jiaqi Li, Junshu Tang, Zhiyong Xu, Longhuang Wu, Yuan Zhou, Shuai Shao, Tianbao Yu, Zhiguo Cao, and Qinglin Lu. Hunyuan-gamecraft: High-dynamic interactive game video generation with hybrid history condition, 2025. 1, 2, 3

- 574 [29] Xirui Li, Chao Ma, Xiaokang Yang, and Ming-Hsuan Yang. Vidtome: Video token merging for zero-shot video editing, 2023. 3 630
- 575 Vidtome: Video token merging for zero-shot video editing, 2023. 3 631
- 576 Vidtome: Video token merging for zero-shot video editing, 2023. 3 632
- 577 [30] Xinyang Li, Tengfei Wang, Zixiao Gu, Shengchuan Zhang, Chunhao Guo, and Liujuan Cao. FlashWorld: High-quality 3D Scene Generation within Seconds, 2025. 1 633
- 578 FlashWorld: High-quality 3D Scene Generation within Seconds, 2025. 1 634
- 579 FlashWorld: High-quality 3D Scene Generation within Seconds, 2025. 1 635
- 580 [31] Zhimin Li, Jianwei Zhang, and and others Lin. Hunyuan-DiT: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding. 1 636
- 581 Hunyuan-DiT: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding. 1 637
- 582 Hunyuan-DiT: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding. 1 638
- 583 [32] Zhen Li, Chuanhao Li, Xiaofeng Mao, Shaoheng Lin, Ming Li, Shitian Zhao, Zhaopan Xu, Xinyue Li, Yukang Feng, Jianwen Sun, Zizhen Li, Fanrui Zhang, Jiaxin Ai, Zhixiang Wang, Yuwei Wu, Tong He, Jiangmiao Pang, Yu Qiao, Yunde Jia, and Kaipeng Zhang. Sekai: A video dataset towards world exploration, 2025. 2 639
- 584 Sekai: A video dataset towards world exploration, 2025. 2 640
- 585 Sekai: A video dataset towards world exploration, 2025. 2 641
- 586 Sekai: A video dataset towards world exploration, 2025. 2 642
- 587 Sekai: A video dataset towards world exploration, 2025. 2 643
- 588 Sekai: A video dataset towards world exploration, 2025. 2 644
- 589 [33] Long Lian, Boyi Li, Adam Yala, and Trevor Darrell. Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models, 2024. 3 645
- 590 Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models, 2024. 3 646
- 591 Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models, 2024. 3 647
- 592 Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models, 2024. 3 648
- 593 [34] Han Lin, Abhay Zala, Jaemin Cho, and Mohit Bansal. Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning, 2024. 3 649
- 594 Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning, 2024. 3 650
- 595 Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning, 2024. 3 651
- 596 [35] Kunhao Liu, Wenbo Hu, Jiale Xu, Ying Shan, and Shijian Lu. Rolling forcing: Autoregressive long video diffusion in real time, 2025. 2 652
- 597 Rolling forcing: Autoregressive long video diffusion in real time, 2025. 2 653
- 598 Rolling forcing: Autoregressive long video diffusion in real time, 2025. 2 654
- 599 [36] Mushui Liu, Yuhang Ma, Yang Zhen, Jun Dan, Yunlong Yu, Zeng Zhao, Zhipeng Hu, Bai Liu, and Changjie Fan. Llm4gen: Leveraging semantic representation of llms for text-to-image generation, 2024. 3 655
- 600 Llm4gen: Leveraging semantic representation of llms for text-to-image generation, 2024. 3 656
- 601 Llm4gen: Leveraging semantic representation of llms for text-to-image generation, 2024. 3 657
- 602 Llm4gen: Leveraging semantic representation of llms for text-to-image generation, 2024. 3 658
- 603 [37] Yifan Liu, Zhiyuan Min, Zhenwei Wang, Junta Wu, Tengfei Wang, Yixuan Yuan, Yawei Luo, and Chunhao Guo. World-mirror: Universal 3d world reconstruction with any-prior prompting. *arXiv preprint arXiv:2510.10726*, 2025. 1 659
- 604 World-mirror: Universal 3d world reconstruction with any-prior prompting. *arXiv preprint arXiv:2510.10726*, 2025. 1 660
- 605 World-mirror: Universal 3d world reconstruction with any-prior prompting. *arXiv preprint arXiv:2510.10726*, 2025. 1 661
- 606 World-mirror: Universal 3d world reconstruction with any-prior prompting. *arXiv preprint arXiv:2510.10726*, 2025. 1 662
- 607 [38] Xiaoxiao Long, Qingrui Zhao, Kaiwen Zhang, Zihao Zhang, Dingrui Wang, Yumeng Liu, Zhengjie Shu, Yi Lu, Shouzheng Wang, Xinzhe Wei, Wei Li, Wei Yin, Yao Yao, Jia Pan, Qiu Shen, Ruigang Yang, Xun Cao, and Qionghai Dai. A survey: Learning embodied intelligence from physical simulators and world models, 2025. 2 663
- 608 A survey: Learning embodied intelligence from physical simulators and world models, 2025. 2 664
- 609 A survey: Learning embodied intelligence from physical simulators and world models, 2025. 2 665
- 610 A survey: Learning embodied intelligence from physical simulators and world models, 2025. 2 666
- 611 A survey: Learning embodied intelligence from physical simulators and world models, 2025. 2 667
- 612 A survey: Learning embodied intelligence from physical simulators and world models, 2025. 2 668
- 613 [39] Yunhong Lu, Yanhong Zeng, Haobo Li, Hao Ouyang, Qiuyu Wang, Ka Leong Cheng, Jiapeng Zhu, Hengyuan Cao, Zhipeng Zhang, Xing Zhu, et al. Reward forcing: Efficient streaming video generation with rewarded distribution matching distillation. *arXiv preprint arXiv:2512.04678*, 2025. 2 669
- 614 Reward forcing: Efficient streaming video generation with rewarded distribution matching distillation. *arXiv preprint arXiv:2512.04678*, 2025. 2 670
- 615 Reward forcing: Efficient streaming video generation with rewarded distribution matching distillation. *arXiv preprint arXiv:2512.04678*, 2025. 2 671
- 616 Reward forcing: Efficient streaming video generation with rewarded distribution matching distillation. *arXiv preprint arXiv:2512.04678*, 2025. 2 672
- 617 Reward forcing: Efficient streaming video generation with rewarded distribution matching distillation. *arXiv preprint arXiv:2512.04678*, 2025. 2 673
- 618 [40] Jack Parker-Holder, Philip Ball, Jake Bruce, Vibhavari Dasagi, Kristian Holsheimer, Christos Kaplanis, Alexandre Moufarek, Guy Scully, Jeremy Shar, Jimmy Shi, Stephen Spencer, Jessica Yung, Michael Dennis, Sultan Kenjeyev, Shangbang Long, Vlad Mnih, Harris Chan, Maxime Gazeau, Bonnie Li, Fabio Pardo, Luyu Wang, Lei Zhang, Frederic Besse, Tim Harley, Anna Mitenkova, Jane Wang, Jeff Clune, Demis Hassabis, Raia Hadsell, Adrian Bolton, Satinder Singh, and Tim Rocktäschel. Genie 2: A large-scale foundation world model. 2024. 2 674
- 619 Genie 2: A large-scale foundation world model. 2024. 2 675
- 620 Genie 2: A large-scale foundation world model. 2024. 2 676
- 621 Genie 2: A large-scale foundation world model. 2024. 2 677
- 622 Genie 2: A large-scale foundation world model. 2024. 2 678
- 623 Genie 2: A large-scale foundation world model. 2024. 2 679
- 624 Genie 2: A large-scale foundation world model. 2024. 2 680
- 625 Genie 2: A large-scale foundation world model. 2024. 2 681
- 626 Genie 2: A large-scale foundation world model. 2024. 2 682
- 627 Genie 2: A large-scale foundation world model. 2024. 2 683
- 628 [41] William Peebles and Saining Xie. Scalable Diffusion Models with Transformers, 2023. *arXiv:2212.09748* [cs]. 1 684
- 629 Scalable Diffusion Models with Transformers, 2023. *arXiv:2212.09748* [cs]. 1 685
- [42] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing, 2023. 3 686
- [43] Bosheng Qin, Juncheng Li, Siliang Tang, Tat-Seng Chua, and Yueting Zhuang. Instructvid2vid: Controllable video editing with natural language instructions, 2024. 3 687
- [44] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-Resolution Image Synthesis with Latent Diffusion Models, 2022. *arXiv:2112.10752* [cs]. 1 688
- [45] Joonghyuk Shin, Zhengqi Li, Richard Zhang, Jun-Yan Zhu, Jaesik Park, Eli Shechtman, and Xun Huang. Motionstream: Real-time video generation with interactive motion controls. *arXiv preprint arXiv:2511.01266*, 2025. 4 689
- [46] Shuai Tan, Biao Gong, Yutong Feng, Kecheng Zheng, Dandan Zheng, Shuwei Shi, Yujun Shen, Jingdong Chen, and Ming Yang. Mimir: Improving video diffusion models for precise text understanding, 2024. 3 690
- [47] HunyuanWorld Team, Zhenwei Wang, Yuhao Liu, Junta Wu, Zixiao Gu, Haoyuan Wang, Xuhui Zuo, Tianyu Huang, Wenhuan Li, Sheng Zhang, Yihang Lian, Yulin Tsai, and Wang-gang others. HunyuanWorld 1.0: Generating Immersive, Explorable, and Interactive 3D Worlds from Words or Pixels, 2025. 1 691
- [48] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020. 13 692
- [49] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Fvd: A new metric for video generation. 2019. 6 693
- [50] Florentina Voboril, Vaidyanathan Peruvemba Ramaswamy, and Stefan Szeider. Streamllm: Enhancing constraint programming with large language model-generated streamliners. In *2025 IEEE/ACM 1st International Workshop on Neuro-Symbolic Software Engineering (NSE)*, pages 17–22. IEEE Computer Society, 2025. 4 694
- [51] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Fei Wu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wente Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. 1, 2, 3, 4 695
- [52] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 3, 13 696

- 688 [53] Thaddäus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane
689 Gu, Nick Matarese, Kevin Swersky, Been Kim, Priyank Jaini,
690 and Robert Geirhos. Video models are zero-shot learners and
691 reasoners. *arXiv preprint arXiv:2509.20328*, 2025. 2
- 692 [54] WorldLabs. Generating worlds. <https://www.worldlabs.ai/blog>, 2024. 1
- 693
694 [55] WorldLabs. Rtfm: A real-time frame model. <https://www.worldlabs.ai/blog/rtfm>, 2025. 1
- 695
696 [56] Zeqi Xiao, Yushi Lan, Yifan Zhou, Wenqi Ouyang, Shuai
697 Yang, Yanhong Zeng, and Xingang Pan. WORLDMEM:
698 Long-term Consistent World Simulation with Memory, 2025.
699 *arXiv:2504.12369 [cs]*. 2
- 700 [57] Shuai Yang, Wei Huang, Ruihang Chu, Yicheng Xiao, Yuyang
701 Zhao, Xianbang Wang, Muyang Li, Enze Xie, Yingcong
702 Chen, Yao Lu, Song Han, and Yukang Chen. Longlive: Real-
703 time interactive long video generation, 2025. 2, 4
- 704 [58] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang,
705 Eli Shechtman, Fredo Durand, and Bill Freeman. Improved
706 distribution matching distillation for fast image synthesis.
707 *Advances in neural information processing systems*, 37:47455–
708 47487, 2024. 5
- 709 [59] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shecht-
710 man, Fredo Durand, William T Freeman, and Taesung Park.
711 One-step diffusion with distribution matching distillation. In
712 *Proceedings of the IEEE/CVF conference on computer vision*
713 *and pattern recognition*, pages 6613–6623, 2024. 5
- 714 [60] Tianwei Yin, Qiang Zhang, Richard Zhang, William T. Free-
715 man, Fredo Durand, Eli Shechtman, and Xun Huang. From
716 slow bidirectional to fast autoregressive video diffusion mod-
717 els, 2025. 4
- 718 [61] Jiwen Yu, Yiran Qin, Xintao Wang, Pengfei Wan, Di Zhang,
719 and Xihui Liu. Gamefactory: Creating new games with gener-
720 ative interactive videos. *arXiv preprint arXiv:2501.08325*,
721 2025. 2
- 722 [62] Jintao Zhang, Jia Wei, Haofeng Huang, Pengle Zhang, Jun
723 Zhu, and Jianfei Chen. Sageattention: Accurate 8-bit attention
724 for plug-and-play inference acceleration, 2025. 14
- 725 [63] Xiangjun Zhang, Litong Gong, Yinglin Zheng, Yansong Liu,
726 Wentao Jiang, Mingyi Xu, Biao Wang, Tiezheng Ge, and
727 Ming Zeng. Rise-t2v: Rephrasing and injecting semantics
728 with llm for expansive text-to-video generation, 2025. 3
- 729 [64] Yifan Zhang, Chunli Peng, Boyang Wang, Puyi Wang,
730 Qingcheng Zhu, Zedong Gao, Eric Li, Yang Liu, and Yahui
731 Zhou. Matrix-game: Interactive world foundation model.
732 *arXiv*, 2025. 1, 2
- 733 [65] Shihao Zhao, Shaozhe Hao, Bojia Zi, Huaizhe Xu, and Kwan-
734 Yee K. Wong. Bridging different language models and gener-
735 ative vision models for text-to-image generation, 2024. 3
- 736 [66] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe,
737 and Noah Snavely. Stereo magnification: Learning view
738 synthesis using multiplane images, 2018. 2

739 A. Interactive Video Data Construction

740 Current video data suitable for training interactive world
741 models remain scarce. Real-world captured videos offer
742 high realism but are costly, time-consuming to collect, and
743 difficult to scale. Simulation-based generation using engines
744 such as *Unreal Engine* provides strong controllability over
745 viewpoint and interactions, yet the heavy modeling and rendering
746 costs restrict scene diversity. Internet videos from
747 platforms like *YouTube* offer massive volume and variety,
748 but their highly inconsistent quality and abundant noise demand
749 complex cleaning pipelines. Public academic datasets,
750 while well-annotated and reliable, are limited in scale and
751 domain coverage. As a result, none of these sources simultaneously
752 satisfy the requirements of interactivity, large scale,
753 and broad diversity, leaving high-quality interactive video
754 data fundamentally insufficient. This scarcity highlights the
755 need for a clearer understanding of what truly qualifies as *interactive*
756 data, which we formalize in the following analysis.

757 A.1. Definition of Interactive Video Data

Interactive Video Data refers to a temporal sequence that explicitly records a **causally driven state-transition process**, in which agents or the environment transition from a clearly defined **initial state** to a significantly different **final state**. The importance of such data lies in its ability to faithfully capture *how an event evolves over time*, rather than in visual complexity.

758
759 A video segment is considered *interactive* if it satisfies
760 any of the following properties:

- 761 • **Significant State Transition.** The video must contain a
762 recognizable and non-trivial macroscopic change of state.
763 It should present clearly distinguishable *pre-condition* and
764 *post-condition* states, with the temporal content between
765 them forming the *transition process*.
- 766 • **Subject Emergence or Interaction.** The main content
767 involves explicit subjects, including:
 - 768 1. *Emergence*: a new subject appears in a previously
769 empty context.
 - 770 2. *Action-driven*: a subject performs an action that
771 changes its own state or affects the environment.
- 772 • **Scene Shift or Evolution.** The video records a fundamental
773 shift or evolution of the scene or background, rather
774 than minor or random perturbations.

775 Interactive videos thus possess **explicit causal structure**,
776 **clear state transitions**, and **perceivable action agents**, enabling
777 world models to learn interpretable action–outcome mappings.
778 Following this definition, we systematically organize interactive
779 data into three principal categories to structure

780 our analysis: (1) **Environmental Interactions**, which
781 encompass global or local scene changes; (2) **Actor Actions**,
782 which are driven by an embodied agent; and (3) **Entity and
783 Object Appearances**, which involve the introduction of new
784 subjects. To facilitate a nuanced evaluation, each category is
785 further divided into *simple* and *complex* settings, reflecting
786 varying degrees of difficulty. Specific examples for each
787 category are provided in Appendix D.2.

788 A.2. Synthetic Data Construction

789 To address the scarcity and high annotation cost of interactive
790 video data, we propose a controllable Synthetic Interaction
791 Video Pipeline for large-scale automated production. While
792 generating synthetic data for training video models has been
793 underexplored, we argue it is now feasible by leveraging
794 the advanced world knowledge and visual representation
795 capabilities of recent foundation models. The effectiveness
796 of our pipeline in producing diverse, high-quality data is
797 showcased in Appendix E (Figs. 18-20).

798 We generate interactive videos starting from an initial
799 frame F_t . To handle diverse visual contexts, we first employ
800 a Vision-Language Model (VLM) to analyze F_t and, guided
801 by a high-level instruction (e.g., “*taking out a torch*”), generate
802 a customized, scene-specific prompt. Based on the
803 interaction type, we then apply one of two distinct strategies:

- 804 1. **Start-End Frame Strategy:** For stationary scenes requiring
805 explicit state transitions (e.g., environmental changes
806 like “*making it snow*”), a VLM guides an image editing
807 model to generate a target end-frame F_t' . This provides
808 strong controllability over the final state.
- 809 2. **First-Frame-Driven Strategy:** For dynamic actions
810 involving significant camera motion (e.g., “*opening a
811 door*”), the model generates freely from only the initial
812 frame. This approach avoids distortions and yields
813 smoother camera movement and temporal continuity.

814 Sourcing specific initial frames for certain interactions,
815 such as “*opening a door*”, is a significant bottleneck, as
816 manual curation is both costly and inefficient. To address
817 this, we leverage an advanced text-to-image model (e.g.,
818 HunyuanImage-3.0 [4]), to synthesize these requisite frames
819 on demand, providing a scalable source of high-quality inputs
820 for our video generation pipeline.

821 A.3. Game Scene Data Curation

822 We build our dataset from over 150 AAA games (e.g., *Assassin's
823 Creed, Cyberpunk 2077*), which provides extensive diversity in
824 environments, lighting, artistic styles, and camera viewpoints
825 is showcased in Appendix E Figs. 17.

826 **Scene and Action-aware Data Partition.** We employ a
827 two-stage partitioning strategy to process the raw videos.
828 First, PySceneDetect [5] segments long videos into visually
829 coherent 6-second clips. Subsequently, we use RAFT-based

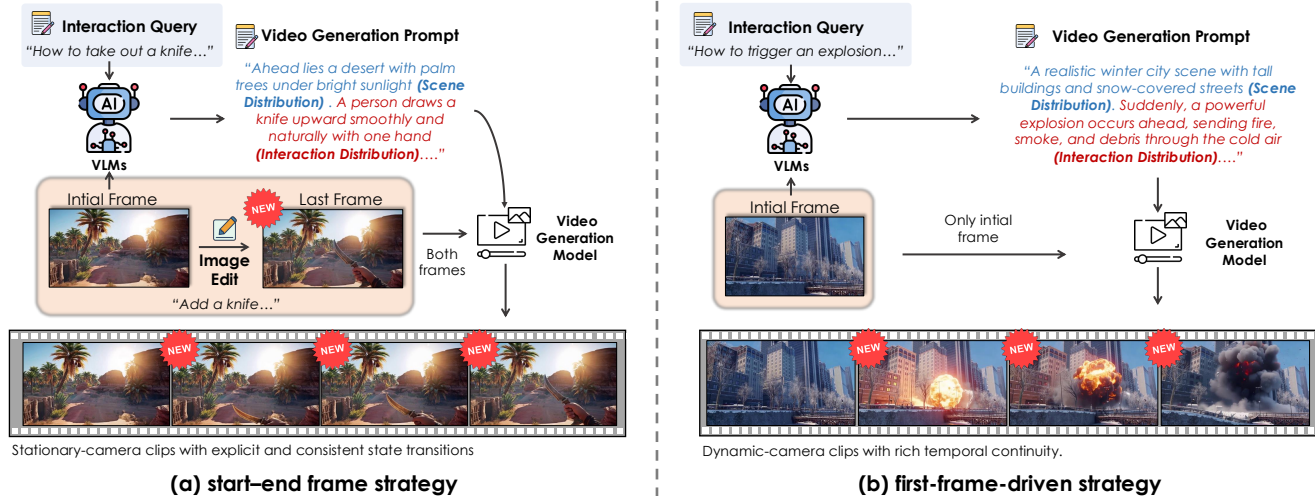


Figure 12. **Showcase of our Synthetic Interaction Video Pipeline.** (a) *The start–end frame strategy* uses a VLM and an image-editing model to construct both initial and edited target frames, enabling controlled state transitions for stationary-camera scenarios. (b) *The first-frame-driven strategy* relies solely on the initial frame and VLM-generated prompts, allowing the video generator to create dynamic, motion-rich interactions with flexible camera movement.

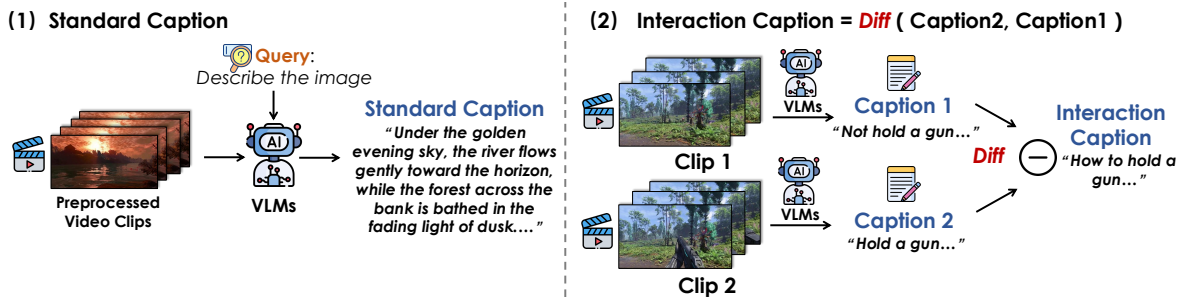


Figure 13. **Pipeline of the Caption Generation System.** The system produces two forms of captions: a *standard caption* that describes the visual content of each clip, and an *interaction caption* derived by computing the semantic difference between consecutive clips. This enables both scene-level descriptions and explicit interaction-oriented annotations for supervision.

830 optical flow [48] to localize fine-grained action boundaries,
831 ensuring each clip preserves temporal integrity for training.

832 **Data Filtering.** To ensure data quality, we perform a three-
833 stage filtering process. A learning-based model first removes
834 low-fidelity or artifact-heavy frames [26]. Next, luminance
835 filtering eliminates scenes that are poorly lit [2]. Finally,
836 a VLM-based semantic check verifies content consistency
837 across frames, retaining only clips with clean visual structure
838 and accurate motion alignment [52].

839 **Camera Annotation.** We reconstruct 6-DoF camera trajec-
840 tories for each clip using VIPE [21]. This process yields
841 frame-by-frame translational and rotational motion esti-
842 mates, providing precise metadata for training camera-aware
843 models and enforcing spatio-temporal consistency.

844 **Structured Captioning.** To provide interaction-aware su-
845 pervision, we devise a structured captioning scheme with
846 two components. First, a **Standard Caption** (C_t), gener-
847 ated by a VLM for each clip, describes the static visual content.
848 Second, an **Interaction Caption** ($I_{t \rightarrow t+1}$) captures the state
849 transition between adjacent clips. This interaction is com-

puted as the semantic difference between their respective
850 standard captions:
851

$$I_{t \rightarrow t+1} = \Delta(\Phi(C_{t+1}), \Phi(C_t)), \quad 852$$

where Φ is a semantic encoder and Δ is a difference operator.
853 This dual-component approach enables the model to jointly
854 learn **appearance-level perception** (from C_t) and **action-**
855 **level reasoning** (from $I_{t \rightarrow t+1}$).
856

B. Detailed Experimental Setup 857

B.1. Model Configurations 858

• **HunyuanVideo.** We use the official inference configu-
859 ration to ensure reproducibility, with FLOW_SHIFT=7.0
860 and EMBEDDED_CFG_SCALE=6.0 controlling motion
861 dynamics and conditional guidance, respectively. All ex-
862 periments are conducted with 50 denoising steps. In ad-
863 dition, both flow_reverse and i2v_stability are enabled to
864 improve temporal robustness across frames.
865

- 866 • **Wan2.2 A14B.** We employ the UniPC sampler with sam- 914
- 867 ple_shift=5.0, sample_steps=40, and boundary=0.900. A 915
- 868 dual-stage classifier-free guidance (CFG) strategy is used, 916
- 869 with guidance scales of (3.5, 3.5) applied to the two noise 917
- 870 regimes to maintain consistent conditioning. 918
- 871 • **LongCatVideo.** We adopt the default high-quality infer- 919
- 872 ence setup, using a guidance_scale of 4 and 50 denoising 920
- 873 steps. Compilation optimizations are enabled to improve 921
- 874 inference efficiency while preserving generation quality. 922

875 B.2. Prompt Design. 914

876 To ensure fair and controlled evaluations, we designed a 915

877 standardized, two-part prompt strategy. This approach con- 916

878 structs two complementary components for each test image: 917

879 an *interaction prompt* to specify the dynamic target action 918

880 or event, and a *base prompt* to describe static scene attributes 919

881 and anchor the generation process to the input image’s ap- 920

882 pearance. During inference, these two prompts are concate- 921

883 nated into a single conditioning sentence and fed directly to 922

884 each model. This decoupled design not only ensures that 923

885 all models receive identical instructions for fair comparison 924

886 but, critically, it also enables controlled evaluations by allow- 925

887 ing us to systematically vary interaction instructions while 926

888 keeping the visual context constant. 927

889 B.3. Test Datasets. 914

890 To comprehensively evaluate interaction-aware video gener- 915

891 ation, we construct two complementary test datasets, each 916

892 targeting different aspects of interaction capability. 917

893 **Viewpoint-Oriented Interaction Test Set.** The first test 914

894 set is designed to evaluate interaction robustness under di- 915

895 verse viewpoints and camera motions. It consists of **75** 916

896 **base images**, where each image is associated with multi- 917

897 ple viewpoint and camera control instructions, including 918

898 translations (front, back, left, right, up, down), diagonal 919

899 viewpoints (upper-left, upper-right, lower-left, lower-right), 920

900 and rotational motions (turn left and turn right). For each 921

901 image–viewpoint combination, a corresponding interaction 922

902 prompt is applied, resulting in a total of **900 generated** 923

903 **videos**. This dataset is used for the quantitative evaluation 924

904 reported in Table 1, enabling a systematic assessment of 925

905 model performance across varied camera perspectives. 926

906 **Fine-Grained Interaction Test Set.** In addition to the 914

907 viewpoint-oriented dataset, we construct a second test set 915

908 focusing on detailed and fine-grained interaction behaviors. 916

909 This dataset is organized around three core interaction dimen- 917

910 sions: (1) **Environmental Interactions**, (2) **Actor Actions**, 918

911 and (3) **Entity and Object Appearances**. Each dimension 919

912 is further divided into fine-grained categories and interaction 920

913 types to support controlled and diagnostic evaluation. 921

For **Environmental Interactions**, the dataset includes 914

weather-related effects (e.g., snow, rain, and lightning) as 915

well as large-scale physical events (e.g., explosions), assess- 916

ing whether models can produce globally consistent and 917

temporally stable scene-level changes. The **Actor Actions** 918

dimension covers both primitive actions, such as drawing 919

a weapon or taking out a tool, and composite, multi-step 920

actions, such as drawing and firing a gun or taking out and 921

operating a phone. To further probe action-specific gener- 922

alization, we additionally include a specialized subset of 923

20 closed-door images for evaluating the *open door* action, 924

which requires consistent state transitions across frames. The 925

Entity and Object Appearances dimension spans a diverse 926

set of entities, including animals, vehicles, and human sub- 927

jects, as well as extended unseen entities such as dragons, 928

enabling evaluation of object emergence, appearance consis- 929

tency, and generalization to novel concepts. 930

Overall, this fine-grained interaction test set consists of 931

100 images, covering a wide range of scenes (indoor/outdoor, 932

natural/urban), lighting conditions, and visual styles (real- 933

istic, game-like, and cartoon). For all evaluations, models 934

are required to generate videos at a unified resolution of 935

832×448 with a fixed length of **93 frames**, ensuring consis- 936

tent and fair comparison across methods. 937

C. Real-time Acceleration Details 938

To further accelerate inference and minimize latency, we 939

incorporate several system-level optimizations: 940

- **FP8 Quantization** reduces memory bandwidth and lever- 941
- ages GPU acceleration while preserving visual quality; 942
- **Parallelized VAE decoding** enables simultaneous latent- 943
- frame reconstruction, mitigating bottlenecks in long- 944
- sequence decoding; 945
- **SageAttention** [62] replaces FlashAttention with an op- 946
- timized quantized attention kernel for faster transformer 947
- computation; and 948
- **Sequence parallelism** distributes video tokens across mul- 949
- tiple GPUs, supporting efficient long-context generation. 950

Together, these techniques boost inference speed to 16 FPS 951

on 8 NVIDIA H20 GPUs, achieving real-time interactive 952

video generation with stable quality and low latency. 953

D. Illustrative Examples of Interactive Video Data 954

To provide a comprehensive understanding of our definition, 955

we present representative examples that clarify the boundary 956

between interactive and non-interactive video data. 957

D.1. Positive Examples: Interactive Video Data 958

The following examples satisfy one or more properties of 959

interactive video data, exhibiting clear causal structures and 960

961

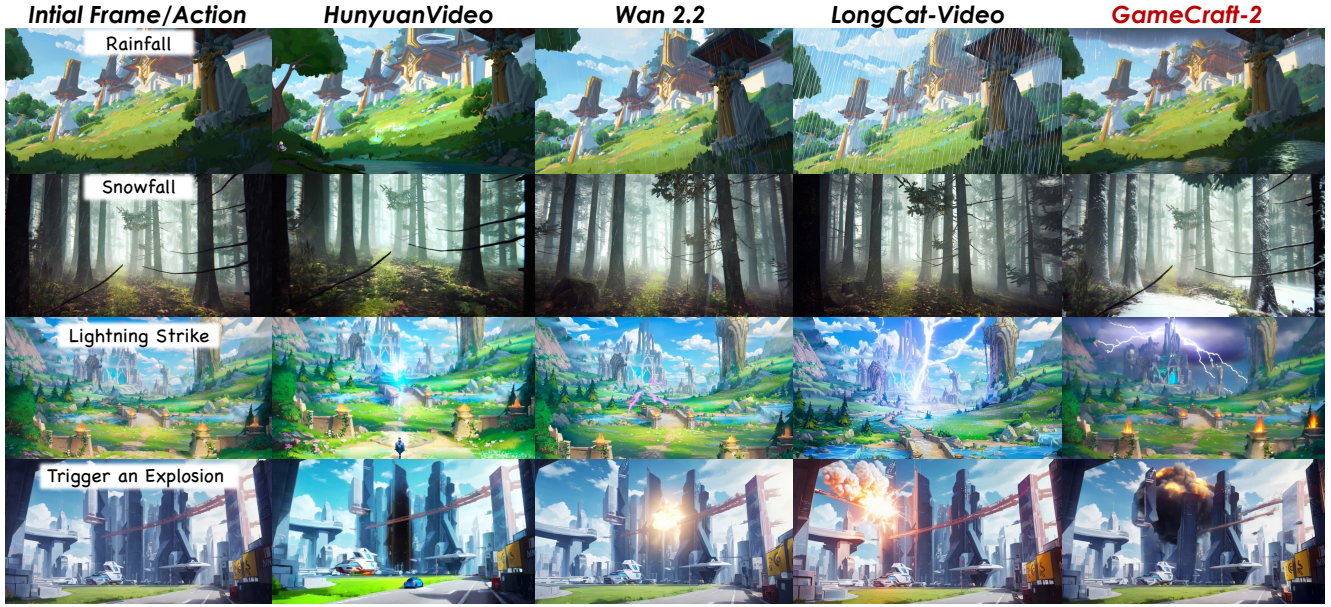


Figure 14. **Comparison of Environmental Interactions with Baseline Models.** Qualitative results showing the fidelity and consistency of environment-level effects. Our approach better preserves global influence and temporal stability.

962 perceivable state transitions.

963 Subject Emergence.

- 964 • *Example 1: Vehicle Appearance.* An empty street (ini-
965 tial state) transitions as a car enters from off-screen and
966 parks at the roadside (transition process), culminating in a
967 scene depicting “a car parked on the street” (final state).
968 The automobile constitutes the emergent core subject,
969 transforming the scene from vacant to occupied.
- 970 • *Example 2: Object Retrieval.* From a first-person per-
971 spective, the frame initially contains only a pair of hands
972 (initial state). The hands retrieve a key from a pocket
973 and hold it prominently (transition process), resulting in
974 a final state of “hands holding a key” (final state). The
975 key represents the emergent core subject.

976 Action-Driven Interaction.

- 977 • *Example 3: Door Opening.* The scene begins with a
978 subject standing before a closed door (initial state). The
979 subject pushes the door open (transition process), lead-
980 ing to a fully open door (final state). This exemplifies
981 direct interaction where the subject acts upon an object,
982 inducing a clear state change.
- 983 • *Example 4: Weapon Discharge.* A character aims a
984 firearm at a target (initial state), pulls the trigger (tran-
985 sition process), resulting in projectile impact and tar-
986 get destruction (final state). This demonstrates action-
987 consequence causality with observable physical effects.

Environmental State Evolution.

- 988 • *Example 5: Weather Transition.* A clear sky (initial state) 989
undergoes gradual cloud accumulation followed by snow- 990
fall of increasing intensity (transition process), ultimately 991
blanketing the entire scene in heavy snow (final state). 992
This represents a fundamental transformation of the envi- 993
ronmental weather attribute. 994
- 995 • *Example 6: Spatial Transition.* Upon opening a door, the 996
camera view shifts from an interior room (initial scene) 997
to an exterior courtyard (final scene). This exemplifies a 998
discrete scene transition driven by subject action, funda- 999
mentally altering the observational context.

These examples, though visually dynamic, lack the defin- 1000
ing characteristics of interactive video data. 1001

Continuous Static Process.

- 1002 • *Example 1: Sustained Blizzard.* A 10-second video seg- 1003
ment depicting continuous heavy snowfall. Although 1004
visually dynamic, the macroscopic state remains constant 1005
as “actively snowing throughout,” lacking a transition 1006
from “no snow” to “snow present.” The absence of state 1007
evolution disqualifies this as interactive data. 1008

Stochastic Background Activity.

- 1009 • *Example 2: Busy Intersection.* A scene featuring continu- 1010
ous pedestrian and vehicular traffic at a crowded intersec- 1011
tion. While abundant motion exists, there is no singular 1012

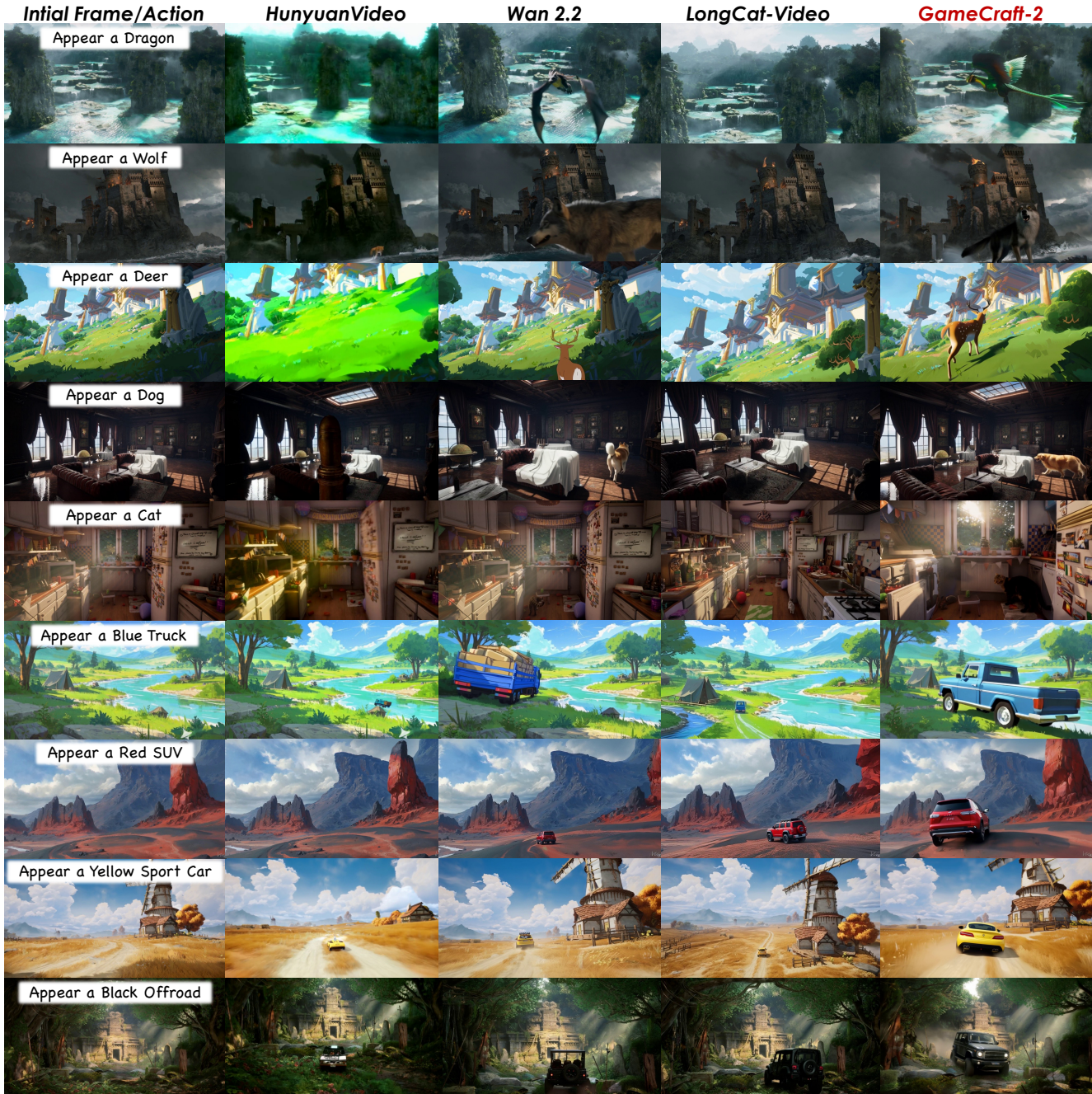


Figure 15. **Comparison of Entity and Object Appearance Interactions with Baseline Models.** Visual comparisons of object emergence and interaction correctness. Our method delivers more accurate, stable, and physically plausible object behaviors.

1013 event-driven macroscopic state change with definitive be-
 1014 ginning and end points. The scene’s overarching state
 1015 persistently remains “*busy intersection*,” lacking a coher-
 1016 ent causal narrative.

1017 **Generalized Motion without Core Subject.**

1018 • *Example 3: Ambient Environmental Fluctuations.* Rip-

ples propagating across a water surface or leaves swaying
 in wind. These phenomena typically constitute random
 environmental perturbations rather than state transitions
 driven by specific subjects or events with explicit causal
 chains. They lack the purposeful, agent-driven transfor-
 mation characteristic of interactive data.

1019
 1020
 1021
 1022
 1023
 1024

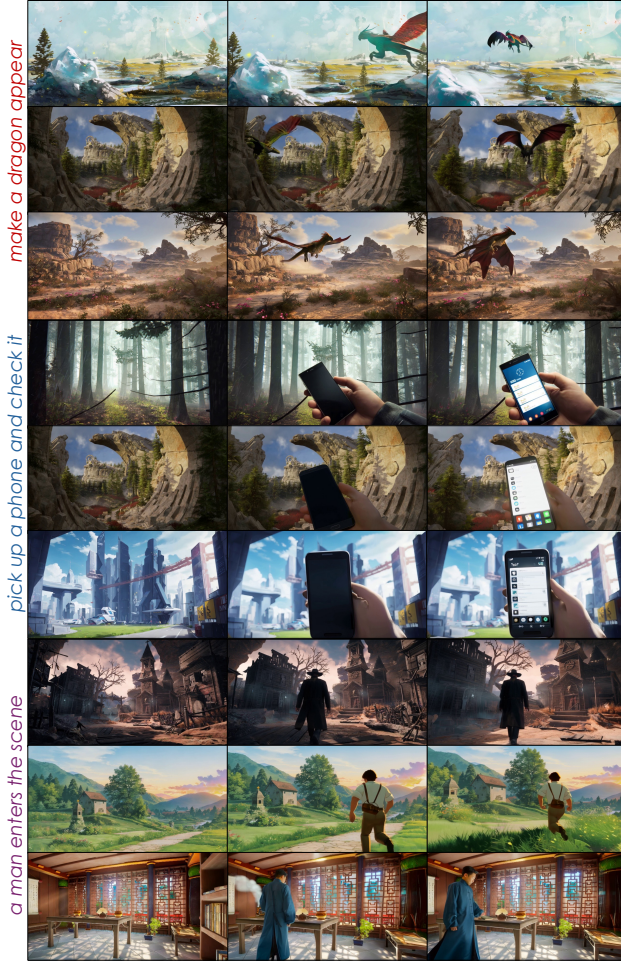


Figure 16. **Generalization Beyond Training Distribution.** GameCraft-2 generates coherent, physically plausible interactions for unseen entities/actions (e.g., a man suddenly appearing, a dragon emerging, taking out a phone), despite such cases not occurring in training, indicating it learns transferable interaction dynamics rather than memorizing visuals.

1025

D.2. Interaction Categories

1026

1027

1028

1029

1030

1031

Following the definition of interactive data in the main text, we provide here a detailed breakdown of the three principal interaction categories used to structure our dataset and analysis. Each category includes both simple and complex settings to reflect different levels of difficulty and to facilitate a fine-grained evaluation of model capabilities.

1032

1033

1034

1035

1036

1037

(1) Environmental Interactions. These interactions reflect global or local scene changes. *Simple cases* include atmospheric effects such as *snowfall* and *rainfall*. *Complex cases* involve more substantial causal transformations, such as *lightning strikes* or *triggering an explosion*, which require coherent illumination changes, particle dynamics, and

physically plausible propagation.

1038

(2) Actor Actions. These interactions are driven by an embodied or first-person actor. *Simple cases* include basic manipulation actions such as *drawing a gun* or *drawing a knife*. *Complex cases* require multi-step or environment-affecting interactions, such as *drawing a torch to illuminate the surroundings*, *firing a gun*, *taking out a phone and operating it*, or *opening a door*. These demand consistent body-object coordination and temporal stability.

1039

1040

1041

1042

1043

1044

1045

1046

(3) Entity and Object Appearances. These interactions introduce new entities into the scene. *Simple cases* include the appearance of a single human or common object. *Complex cases* involve entities with more distinct geometry or motion priors, such as animals (cat, dog, deer, wolf, dragon) or vehicles (red SUV, yellow sports car, blue truck, black off-road vehicle), which require accurate spatial placement, scale consistency, and stable identity preservation.

1047

1048

1049

1050

1051

1052

1053

1054

E. Dataset Showcase

1055

This appendix provides visual examples from our constructed dataset, which is composed of two primary sources: curated real-world gameplay footage and synthetically generated interactive videos. The following sections showcase the diversity and quality of each data type.

1056

1057

1058

1059

1060

E.1. Curated Gameplay Data

1061

The following figures illustrate the rich diversity of our curated gameplay data, collected from over 150 AAA games. As shown, the dataset covers a wide array of interaction contexts, including both first-person and third-person viewpoints (Fig. 17), as well as a comprehensive range of environments spanning natural and urban scenes under various lighting, weather. This diversity is crucial for training robust and generalizable world models.

1062

1063

1064

1065

1066

1067

1068

1069

E.2. Synthetic Interaction Data

1070

Generated by our synthetic data pipeline, the following examples demonstrate the pipeline’s capability to create controlled and high-quality interactive videos. These examples cover the three main interaction categories defined in our work: **Environmental Interactions** such as weather changes and explosions (Fig. 18), **Actor Actions** involving complex body-object coordination (Fig. 19), and **Entity/Object Appearances** that introduce new subjects into the scene with high fidelity (Fig. 20).

1071

1072

1073

1074

1075

1076

1077

1078

1079

F. Detailed Comparison Across Interaction Dimensions

1080

1081

In this section, we present a concise comparison of our approach against baseline models across three key dimensions

1082

1083



Figure 17. **Examples of First-person and Third-person Interactive Gameplay Videos.** Samples showing diverse actor actions under different viewpoints, illustrating rich interactive semantics captured from our gameplay collection.

1084 of interactive video generation: environmental interactions,
 1085 actor–action dynamics, and entity or object appearance be-
 1086 haviors. As shown in Tab 2, our method achieves higher
 1087 temporal stability, more coherent action execution, and more
 1088 accurate object emergence, consistently outperforming exist-
 1089 ing models across diverse interaction scenarios.

1090 G. InterBench: A Detailed Protocol for Bench- 1091 marking Action-Level Interaction

1092 **Motivation and Design Philosophy.** Existing video genera-
 1093 tion benchmarks, such as Fréchet Video Distance or CLIP
 1094 Score, primarily assess perceptual quality, temporal consis-
 1095 tency, and static text-video alignment. While valuable,
 1096 they are ill-suited for evaluating *interactive* video genera-
 1097 tion, where the primary task is to render a causal change in re-
 1098 sponse to a specific action command. These metrics cannot
 1099 distinguish between a correctly executed action and a visu-
 1100 ally plausible but semantically incorrect video. To fill this
 1101 critical gap, we designed **InterBench**, an evaluation protocol
 1102 specifically tailored to measure the fidelity of action-level

interactions. Its philosophy is to deconstruct the complex
 1103 concept of a “good interaction” into a set of distinct, measur-
 1104 able, and interpretable dimensions, enabling a fine-grained
 1105 analysis of model capabilities and failure modes. 1106

Interaction Trigger Rate. This dimension serves as the
 1107 most fundamental, gateway assessment. It asks the question:
 1108 *Did the requested interaction happen at all?* This metric
 1109 is designed to isolate the model’s basic ability to acknowl-
 1110 edge and act upon an instruction, separating cases where the
 1111 model successfully initiated the action from those where it
 1112 completely failed to respond. This is a binary metric: 1113

- **1 (Success):** The requested interaction is initiated in the
 1114 video. For instance, for the prompt *draw a gun*, this score
 1115 is given if a gun becomes visible. If this score is given, the
 1116 subsequent dimensions are evaluated on their respective
 1117 scales. 1118
- **0 (Failure):** The requested interaction does not occur at
 1119 all. The model ignores or completely misunderstands the
 1120 interaction prompt. If this score is given, all subsequent
 1121

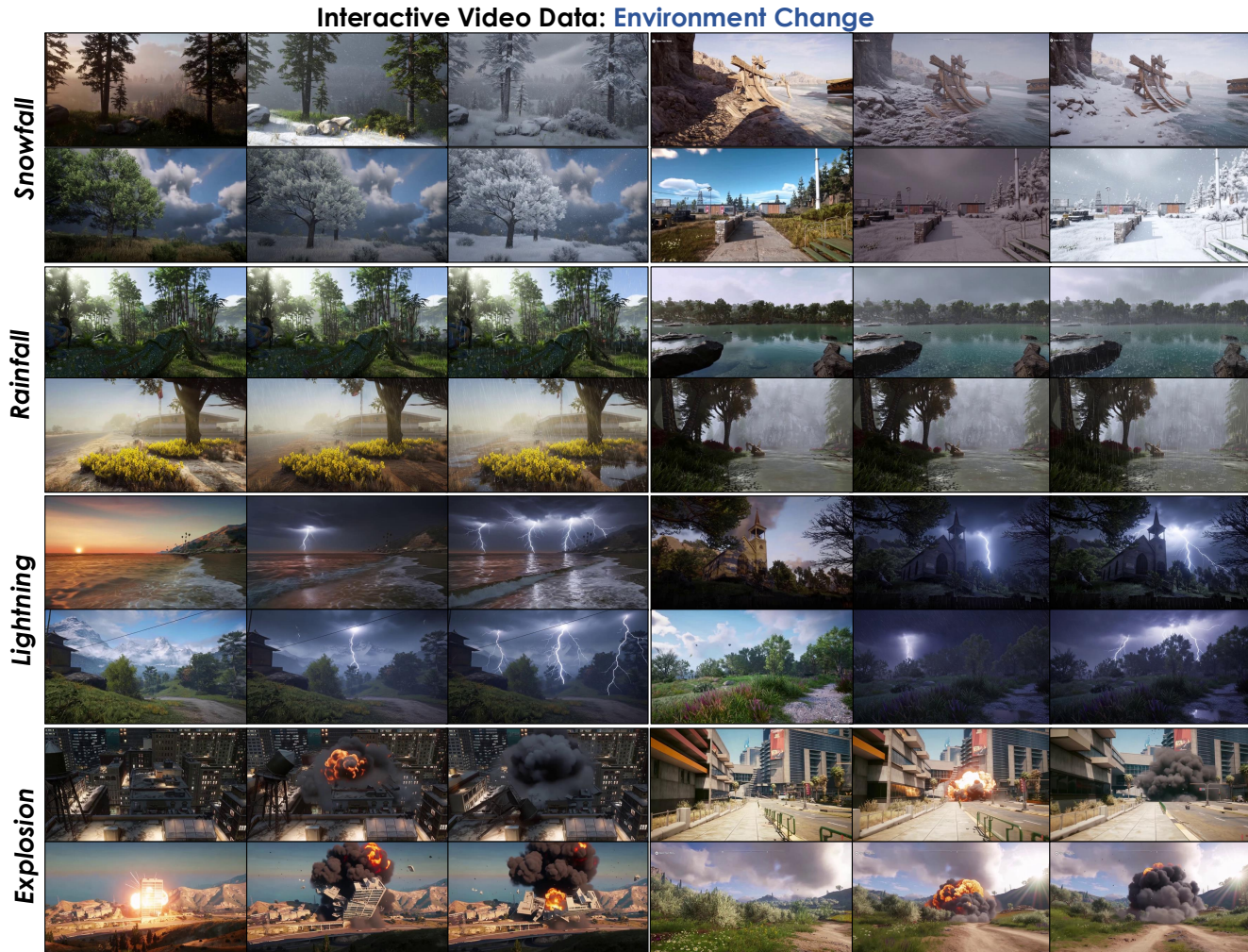


Figure 18. **Synthetic Examples of Environmental Interactions.** Examples of synthetic scene-change interactions generated by our pipeline, covering *snowfall*, *rainfall*, *lightning*, and *explosions*.

Table 2. **Quantitative performance evaluation**, against state-of-the-art competitors on the InterBench protocol. Scores are presented across six key interaction dimensions, with our model’s superior results highlighted.

| Category | Method | Trigger | Align | Fluency | Scope | EndState | Physics |
|-----------------------------|---------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Environmental Interactions | Wan2.2 A14B | 0.799 | 3.511 | 3.579 | 3.722 | 3.951 | 3.008 |
| | LongCat-Video | 0.897 | 3.963 | 3.777 | 4.188 | 4.377 | 3.210 |
| | HunyuanVideo | 0.490 | 1.950 | 1.940 | 2.065 | 2.308 | 1.670 |
| | GameCraft-2 | 0.962 | 4.342 | 4.247 | 4.578 | 4.688 | 3.893 |
| Actor Actions | Wan2.2 A14B | 0.836 | 3.490 | 3.488 | 4.036 | 4.054 | 3.175 |
| | LongCat-Video | 0.806 | 3.089 | 3.005 | 3.832 | 3.771 | 2.839 |
| | HunyuanVideo | 0.587 | 2.147 | 2.202 | 2.717 | 2.748 | 1.931 |
| | GameCraft-2 | 0.983 | 4.087 | 4.191 | 4.576 | 4.686 | 3.828 |
| Entity & Object Appearances | Wan2.2 A14B | 0.874 | 3.943 | 3.545 | 4.281 | 4.265 | 3.054 |
| | LongCat-Video | 0.712 | 3.050 | 2.758 | 3.340 | 3.482 | 2.352 |
| | HunyuanVideo | 0.607 | 2.037 | 1.870 | 2.736 | 2.734 | 1.462 |
| | GameCraft-2 | 0.944 | 4.292 | 3.978 | 4.410 | 4.514 | 3.578 |

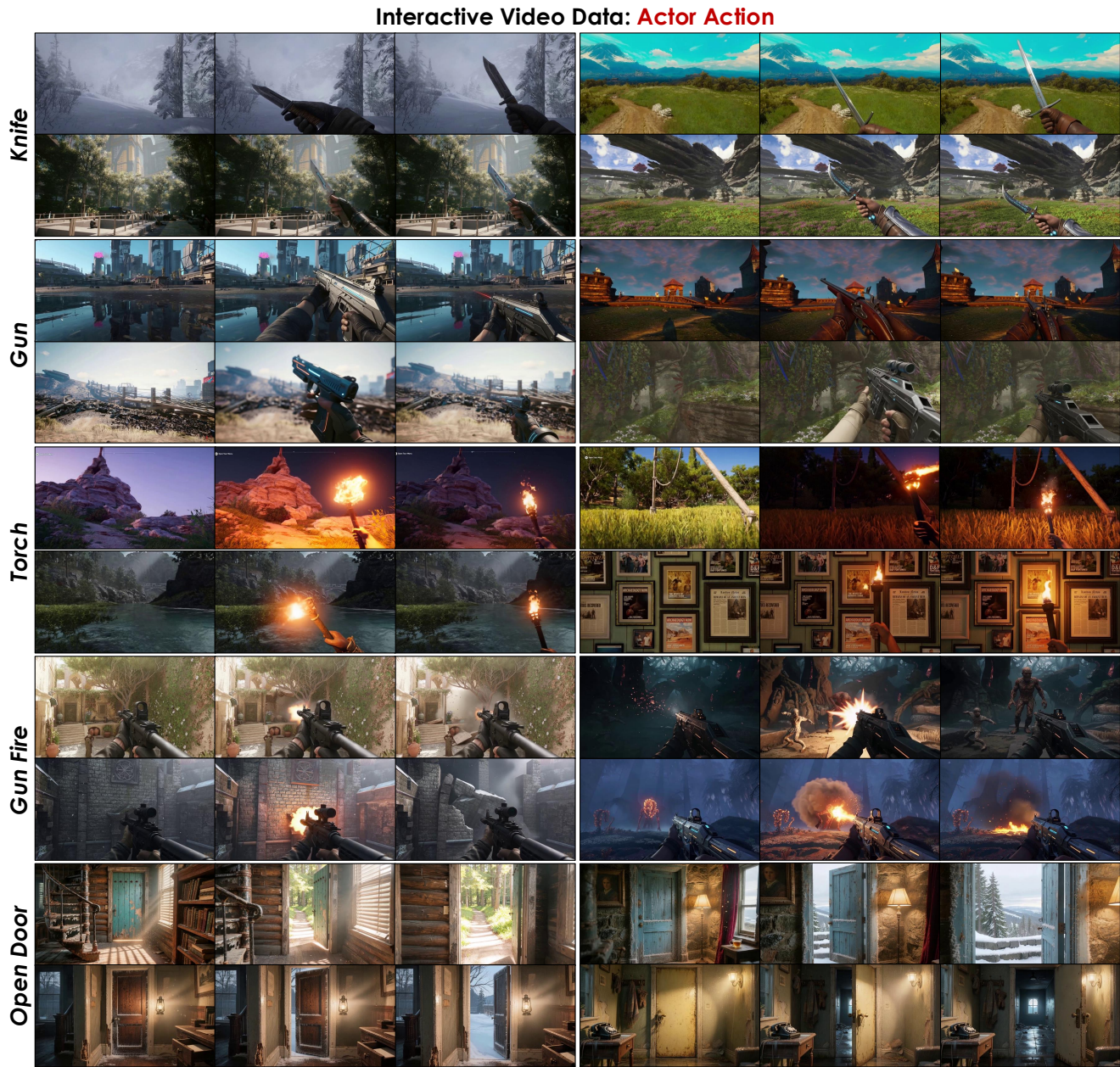


Figure 19. **Synthetic Examples of Actor Actions.** Examples illustrating the range of interactive behaviors synthesized by our pipeline. *Actor-driven interactions* include *drawing a knife*, *drawing a gun*, *drawing a torch*, *firing a weapon*, and *opening a door*, which require consistent body–object coordination and temporal coherence.

1122 dimensions are automatically scored 0.

1123 **Prompt–Video Alignment.** Beyond simply triggering an
1124 action, this dimension evaluates the semantic fidelity of the
1125 generated video with respect to the *entire* prompt (both the
1126 base scene description and the interaction command). It
1127 ensures the interaction happens in the *right way* and the *right*
1128 *context*, encompassing both static and dynamic alignment.

This metric is scored on a 0-1-3-5 ordinal scale, contingent
on the interaction being triggered:

- **5 (Excellent):** Both the static context (scene, style) and the dynamic action perfectly match the prompt’s description.
- **3 (Moderate):** The primary action is correct, but there are minor semantic deviations in the scene’s context or the specifics of the action’s execution.
- **1 (Poor):** A recognizable interaction occurs, but it involves

Interactive Video Data: Entity and Object Appearances

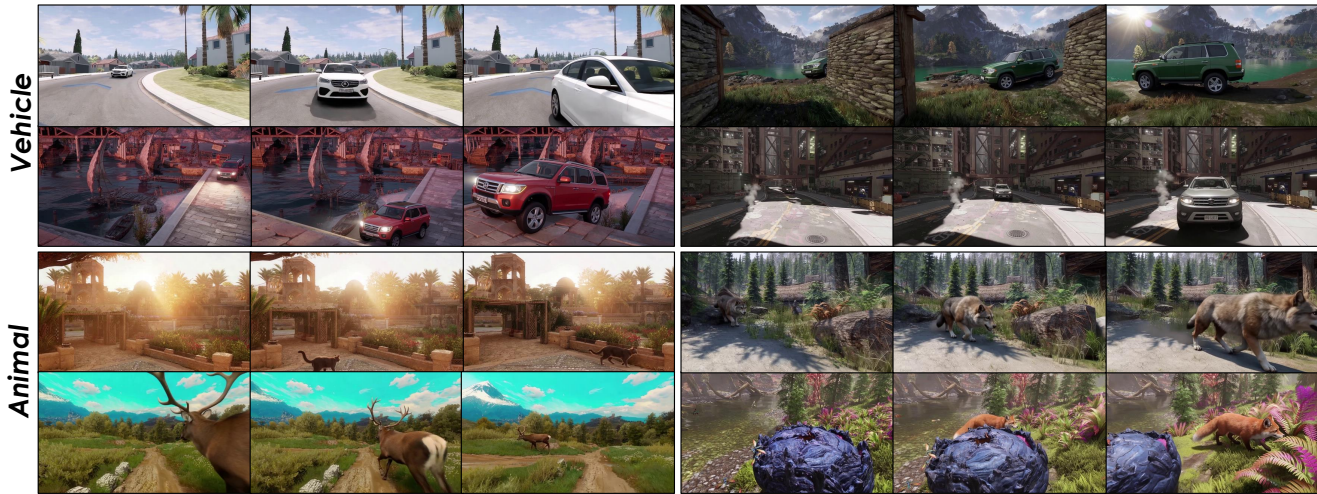


Figure 20. **Synthetic Examples of Entity/Object Appearances.** Examples illustrating the range of interactive behaviors synthesized by our pipeline. *Entity- and object-level interactions* include *animal intrusion* and *vehicle entry*, showing the pipeline’s capability to introduce new entities with realistic geometry, scale consistency, and stable identity across frames.

- 1137 a major semantic error, such as performing the wrong action (e.g., closing instead of opening a door) or generating
 1138 a scene that bears no resemblance to the base prompt.
 1139
- 1140 • **0 (Failure):** The triggered video content shows no meaningful semantic alignment with either the prompt’s context
 1141 or its specified action.
 1142
- 1143 **Interaction Fluency.** This dimension measures the temporal naturalness and continuity of the interaction process. It
 1144 specifically penalizes temporal discontinuities such as abrupt teleportation of objects, noticeable frame jumps, unrealistic
 1145 motion jitter, and structural tearing of geometry, particularly around the interacting regions. This metric is scored on a
 1146 0-1-3-5 ordinal scale:
 1147
- 1148 • **5 (Excellent):** The motion is perfectly smooth, continuous, and natural, with no temporal artifacts present.
 1149
 - 1150 • **3 (Moderate):** The motion is generally continuous but contains minor, non-disruptive artifacts like slight jitter or
 1151 a single inconspicuous jump-cut.
 1152
 - 1153 • **1 (Poor):** The interaction is plagued by severe temporal artifacts (e.g., constant flickering, object teleportation) that
 1154 significantly disrupt the viewing experience.
 1155
 1156
 1157
- 1158 **Interaction Scope Accuracy.** This metric assesses a model’s spatial reasoning by examining whether the spatial
 1159 extent and environmental influence of an interaction are plausible and consistent with its expected scope (global or
 1160 local). This metric is scored on a 0-1-3-5 ordinal scale:
 1161
- 1162 • **5 (Excellent):** The spatial influence of the interaction is physically and semantically correct (e.g., global effects are
 1163 global, local effects are local and propagate realistically).
 1164
- 1165 • **3 (Moderate):** The scope is generally correct but with minor inaccuracies, such as a global effect not covering
 1166 the entire scene or a local effect having a slightly incorrect area of influence.
 1167
 - 1168 • **1 (Poor):** The scope is fundamentally wrong. For example, a global event is rendered as a tiny local patch, or a local
 1169 effect implausibly affects the entire scene.
 1170
 1171
 1172
- 1173 **End-State Consistency.** A successful interaction must not only be initiated correctly but also *converge* to a stable and
 1174 correct outcome. This dimension evaluates the final state of the video to ensure the result of the action persists as
 1175 expected. This metric is scored on a 0-1-3-5 ordinal scale:
 1176
- 1177 • **5 (Excellent):** The interaction converges to the correct final state, which remains stable until the end of the video.
 1178
 - 1179 • **3 (Moderate):** The final state is mostly correct but exhibits minor instability, such as slight flickering, object drift, or
 1180 subtle geometric inconsistencies.
 1181
 - 1182 • **1 (Poor):** The interaction fails to converge correctly. The final state is incorrect, highly unstable (e.g., oscillating),
 1183 or the effects of the action vanish prematurely.
 1184
 1185
- 1186 **Object Physics Correctness.** This dimension focuses on the physical plausibility and structural integrity of the objects
 1187 and agents involved in the interaction, evaluating whether their behavior adheres to basic physical principles like object
 1188 permanence, rigidity, and kinematics. This metric is scored on a 0-1-3-5 ordinal scale:
 1189
 1190
 1191

- 1192 • **5 (Excellent):** All objects and agents maintain structural
1193 integrity and interact in a physically plausible manner.
1194 There is no unnatural deformation, interpenetration, or
1195 kinematic errors.
- 1196 • **3 (Moderate):** Minor physical inaccuracies are present,
1197 such as slight object warping during movement or brief,
1198 non-critical interpenetration between an agent and an ob-
1199 ject.
- 1200 • **1 (Poor):** Severe physical violations occur. Objects unnat-
1201 urally deform, agents pass through solid objects, or motion
1202 is kinematically impossible.

1203 H. Limitation and Future Work

1204 Despite its advancements, our framework has several limita-
1205 tions that highlight avenues for future research. First, while
1206 our randomized long-video tuning strategy alleviates error
1207 accumulation in autoregressive generation, it does not en-
1208 tirely eliminate it, and semantic drift may still manifest in
1209 long sequences (More than 500 frames). This is partly at-
1210 tributable to our model’s lack of an explicit long-term mem-
1211 ory mechanism, a crucial component for advanced world
1212 models, as it relies instead on the finite capacity of its KV
1213 cache. Furthermore, the scope of supported interactions is
1214 currently centered on single-step, immediate-effect actions.
1215 Enabling multi-stage tasks that require logical reasoning or
1216 planning remains a significant future challenge. Finally, al-
1217 though we achieve real-time performance at 16 FPS, further
1218 optimization is required to reduce latency for highly reac-
1219 tive gameplay and to enable deployment on more accessible
1220 hardware.