

From Virtual Games to Real-World Play

Wenqiang Sun^{*1,2}, Fangyun Wei^{*2}, Jinjing Zhao³, Xi Chen², Zilong Chen⁴,
Hongyang Zhang⁵, Jun Zhang^{†1}, Yan Lu^{†2}

¹HKUST ²Microsoft Research ³University of Sydney

⁴Tsinghua University ⁵University of Waterloo

wsunap@connect.ust.hk, {fawei,xichen6,yanlu}@microsoft.com,
jzha0100@uni.sydney.edu.au, chenz122@mails.tsinghua.edu.cn,
hongyang.zhang@uwaterloo.ca, eejzhang@ust.hk

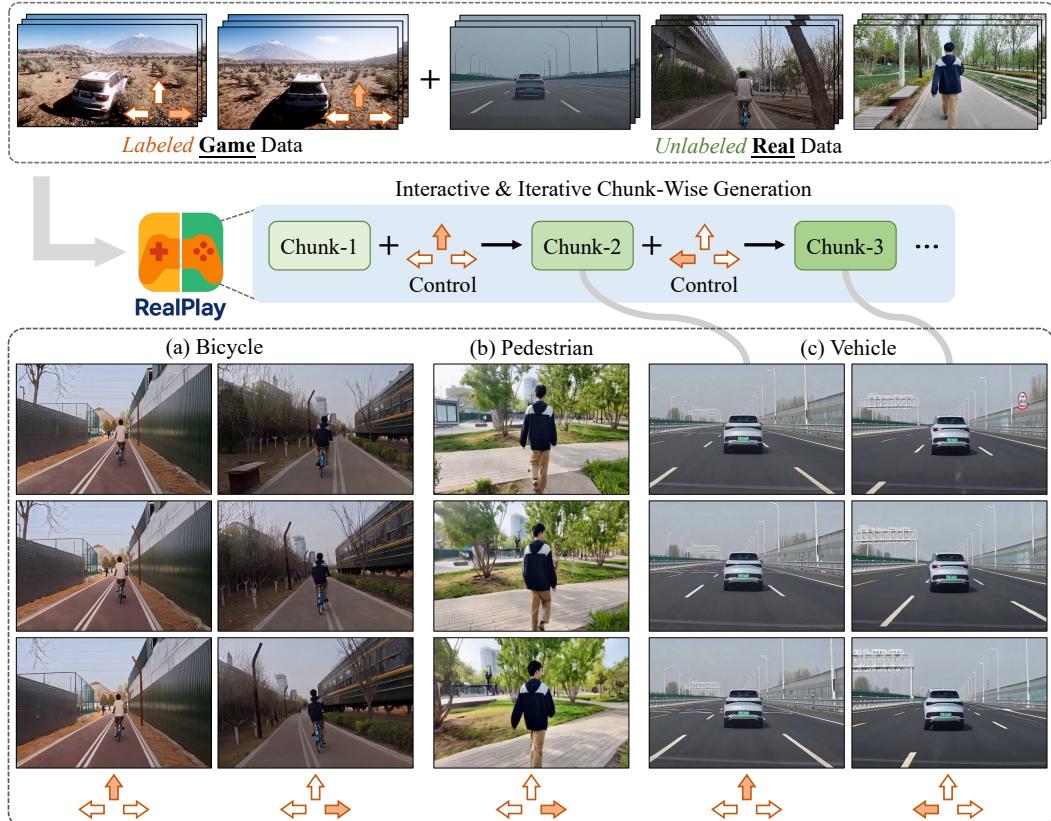


Figure 1: *RealPlay* is a neural-network-driven real-world game engine with three key characteristics: (1) It supports iterative interaction—at each iteration, users observe the current visual scene, provide control signals, and receive control-accurate, temporally consistent, and realistic video chunks in response. (2) It eliminates the need for annotated real-world data while exhibiting strong control transfer capabilities, effectively mapping control signals (e.g., “move forward”, “turn left” and “turn right”) from the game environment to the real world. (3) It demonstrates entity transfer capabilities: although the labeled game data are sourced exclusively from the car racing game *Forza Horizon 5*, *RealPlay* successfully generalizes these control signals to other real-world entities such as (a) bicycles and (b) pedestrians, beyond (c) vehicles. **Additional visualizations are provided in the appendix.**

^{*}Equal Contribution. [†]Corresponding author.

Abstract

We introduce RealPlay, a neural network-based real-world game engine that enables interactive video generation from user control signals. Unlike prior works focused on game-style visuals, RealPlay aims to produce photorealistic, temporally consistent video sequences that resemble real-world footage. It operates in an interactive loop: users observe a generated scene, issue a control command, and receive a short video chunk in response. To enable such realistic and responsive generation, we address key challenges including iterative chunk-wise prediction for low-latency feedback, temporal consistency across iterations, and accurate control response. RealPlay is trained on a combination of labeled game data and unlabeled real-world videos, without requiring real-world action annotations. Notably, we observe two forms of generalization: (1) control transfer—RealPlay effectively maps control signals from virtual to real-world scenarios; and (2) entity transfer—although training labels originate solely from a car racing game, RealPlay generalizes to control diverse real-world entities, including bicycles and pedestrians, beyond vehicles. Project page can be found: <https://wenqsun.github.io/RealPlay/>

1 Introduction

Recent works such as GameFactory [51] and GameNGen [39] demonstrate that neural networks can effectively model games and even function as game engines. These approaches typically adopt a two-stage pipeline: (1) collecting large-scale game data where each frame or chunk is annotated with control signals (e.g., “move forward” or “turn left” in a car racing game); and (2) training an interactive visual content generator—such as a video generation model—on the labeled data to predict future visual frames conditioned on the current visual context and the corresponding control signal.

Although neural networks can effectively model and replicate game environments, the upper bound of the visual quality they generate is ultimately constrained by the underlying game engine. While modern game engines—such as Unreal Engine 5—can produce highly realistic graphics, humans can still easily distinguish between game-rendered visuals and real-world footage. In other words, current game engines struggle to generate visuals that are indistinguishably close to reality or faithfully capture the complex physical laws of the real world. This motivates us to explore the feasibility of using neural networks to develop games whose visual outputs appear more realistic than those rendered by conventional game engines.

In this work, we present *RealPlay*, an interactive video generation model that functions as a real-world game engine, built upon recent advances in diffusion-based video generation models [18, 22, 23, 29, 41, 48, 53]. Notably, RealPlay does not rely on labeled real-world data; instead, it leverages labeled game data and unlabeled real-world data, demonstrating strong control transfer capabilities from virtual environments to real-world scenarios. As illustrated in Figure 1, RealPlay enables users to engage in an interactive loop: they observe the current visual scene, provide control signals, and receive temporally consistent, photorealistic frames in response. The newly generated frames then serve as the observation for the next interaction step. RealPlay tackles four key challenges in interactive video generation:

Chunk-Wise Generation. RealPlay is powered by a pre-trained image-to-video generator [48], originally designed to produce long-horizon videos in a single forward pass. However, generating long clips at each iteration may result in noticeable delays, reducing the responsiveness of the interaction. To address this, we adapt the pre-trained video generator to produce shorter clips—consisting of only a few frames—at each iteration, enabling a more responsive user experience.

Iterative Generation. Existing video generation models are typically designed for one-shot image-to-video generation, which does not align with our use case where video chunks must be generated iteratively, with RealPlay receiving a new control signal at each iteration. To support this iterative generation process, we adapt a pre-trained video generator—originally conditioned on a single image to produce an entire video—into a chunk-wise generation framework, where each new segment is conditioned on the previously generated video chunk rather than just a static image.

Consistency Across Iterations. During inference, RealPlay generates visual outputs conditioned on the current observation—which is itself generated in the previous iteration—and the user-provided control signals. In other words, each new video chunk is generated based on *predicted observations*,

rather than the *ground-truth observations* used during training. This discrepancy introduces a distribution gap between training and inference, often resulting in accumulated artifacts, visual drift, or temporal inconsistencies across iterations. To mitigate this issue, we adopt the Diffusion Forcing [4] strategy, which has been shown effective in narrowing this gap. Specifically, during training, we introduce noise to the conditional chunk, encouraging the model to better handle imperfect conditions during inference.

Precise Control. Annotating game data is relatively easy and highly scalable, as control signals can be automatically recorded during gameplay. In contrast, labeling real-world data is time-consuming, ambiguous, and often requires manual efforts. RealPlay explores a mixed training paradigm that combines labeled game data with unlabeled real-world data. Our findings show that this approach enables RealPlay to learn effective control strategies in the game domain and successfully transfer them to real-world scenarios—demonstrating promising control transfer capability without requiring explicit real-world annotations. Interestingly, although our labeled game data come solely from the car racing game *Forza Horizon 5*, the control signals (i.e., “move forward”, “turn left”, and “turn right”) can be effectively transferred to control real-world entities beyond vehicles—such as bicycles and pedestrians—as illustrated in Figure 1.

Our experimental results demonstrate that RealPlay is capable of generating visually compelling, control-accurate, and realistic video sequences, while supporting iterative user interaction. RealPlay represents an initial step toward validating the feasibility of using neural networks as real-world game engines—paving the way for a new paradigm where interactive, high-fidelity simulations are driven purely by data and learned dynamics, rather than traditional graphics engines.

2 Related Work

Video Diffusion Models. The field of video generation has advanced rapidly with the emergence of diffusion models [6, 7, 17], with modern architectures demonstrating exceptional temporal coherence and physical plausibility [22, 24, 37, 41, 48]. These models support multi-modal conditioning—ranging from text-to-video synthesis [18, 34, 48] to image animation [13, 33, 46]—while maintaining stable dynamics over extended sequences. Beyond creative applications, their ability to capture realistic dynamics positions them as promising candidates for world simulators [2, 8, 9, 27, 49], with early demonstrations in tasks such as game engines and embodied planning.

Within this landscape, bidirectional diffusion models have gained prominence by employing global spatio-temporal attention to capture inter-frame dependencies, enabling the generation of high-fidelity videos with strong temporal coherence [22, 23, 25, 29, 41, 53]. However, their quadratic computational complexity—stemming from full-sequence attention—limits them to generating only short-horizon videos. To overcome this limitation, autoregressive diffusion frameworks have been introduced, generating videos iteratively by conditioning on previously generated frames [1, 11, 14, 20, 21, 39, 42]. While early variants successfully extend sequence lengths, they often suffer from performance degradation over time, leading to diminished visual quality and temporal consistency. To address this, diffusion forcing techniques are proposed, introducing flexible temporal conditioning mechanisms that significantly mitigate these issues [4, 36]. Building upon this foundation, recent autoregressive methods have further improved long-term coherence and visual fidelity, enabling efficient and stable generation of significantly longer videos [5, 12, 32, 43, 45, 50].

World Models. Building on the generative strengths of video diffusion models, recent research has explored their potential as interactive world models capable of simulating complex, controllable environments [1–3, 8, 39, 44, 47]. These efforts move beyond passive video generation toward agent-environment interaction, using diffusion-based architectures to synthesize responsive and dynamic virtual worlds. Foundational works [9, 27, 51, 52] highlight the feasibility of generating playable environments from diverse inputs—such as images, text, or user actions—thus enabling real-time interaction and embodied planning. However, while many of these models produce visually rich simulations, they often suffer from limited generalization due to overfitting on specific virtual datasets. In contrast, RealPlay targets real-world scenarios, enabling photorealistic and controllable video generation from user input, while demonstrating strong generalization across domains and entities.

3 Method

Problem Formulation. The objective is to train *RealPlay* using a labeled car racing dataset, in which user actions (i.e., “move forward”, “turn left”, and “turn right”) are recorded at a

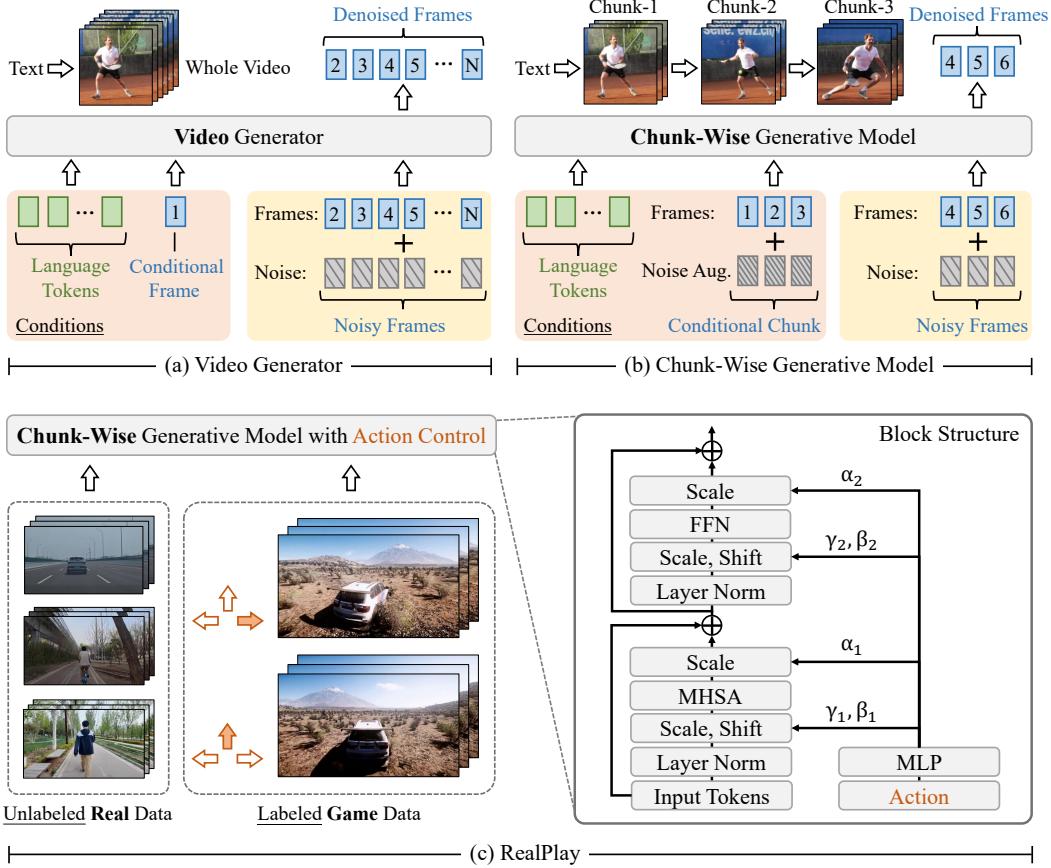


Figure 2: RealPlay involves a two-stage training process. *Stage-1*: We adapt a pre-trained image-to-video generator (Figure (a))—which generates an entire video in a single pass conditioned on a single frame—into a chunk-wise generation model (Figure (b)), which generates video chunks iteratively, conditioned on the previously generated chunk. This adaptation includes several key modifications detailed in Section 3.1. *Stage-2*: RealPlay (Figure (c)) is trained on a combination of a labeled game dataset and an unlabeled real-world dataset, enabling action transfer from controlling a car in the game environment to manipulating various entities in the real world. This is achieved by modifying the chunk-wise generation model to incorporate action control through an adaptive LayerNorm mechanism. In all figures, “frames” refer to frame latents encoded by the video VAE encoder from CogVideoX [48]. For clarity, we omit the details of injecting noise timestep embeddings.

fixed frequency, alongside an unlabeled dataset consisting of real-world videos that capture featuring vehicles in motion, bicycles being ridden, and pedestrians walking, as shown in Figure 1. Once trained, RealPlay enables control transfer from car-based actions in the game environment to real-world entities including vehicles, bicycles, and pedestrians. It supports interactive and iterative video generation: at each iteration, users observe the scene generated in the previous step and issue new control commands. RealPlay then generates a new, temporally consistent video chunk that reflects the user’s input.

Training Data. Formally, let $\mathcal{G} = \{G_i\}_{i=1}^{M_1}$ denote the *labeled* game dataset comprising M_1 training samples. Each sample $G_i = \{(C_k, a_k)\}_{k=1}^K$ consists of K pairs of video chunks and corresponding action commands. An inherent temporal association exists between adjacent pairs: the k -th video chunk C_k combined with the action a_k leads to the $(k+1)$ -th video chunk C_{k+1} . Each action a_k is represented as a 3-dimensional one-hot vector, corresponding to the discrete set {move forward, turn left, turn right}. Let $\mathcal{R} = \{R_i\}_{i=1}^{M_2}$ represent the *unlabeled* real-world dataset including M_2 samples. Each sample $R_i = \{C_k\}_{k=1}^K$ contains a sequence of K video chunks, where adjacent chunks exhibit temporal consistency. Each sequence captures real-world motion patterns involving vehicles, bicycles, or pedestrians, each appearing exclusively within a given sequence.

Overview. Figure 2 presents an overview of RealPlay, a two-stage training pipeline:

1. In stage 1, we adapt a pre-trained image-to-video generator—originally designed to generate long videos conditioned on clean, ground-truth images—into a chunk-wise generative model capable of generating video segments conditioned on noisy, previously generated chunks. Details are provided in Section 3.1.
2. In stage 2, we fine-tune this chunk-wise model using both the *labeled* game dataset \mathcal{G} and the *unlabeled* real-world dataset \mathcal{R} , ultimately training RealPlay to perform chunk-wise generation while enabling control transfer from the game to the real world, as detailed in Section 3.2.

3.1 Chunk-Wise Generation

Image-to-Video Generator. We begin by adapting a pre-trained image-to-video generator into a chunk-wise generative model. Specifically, we leverage CogVideoX, a large-scale model capable of generating videos with a duration of 49 frames. CogVideoX [18] employs T5 [30] as its text encoder, which encodes a language instruction into a fixed-length token sequence denoted as \mathcal{T} . It also includes a video VAE encoder with a downsampling rate of 4, which encodes each video into a sequence of latent representations $\mathcal{F} = \{\mathbf{F}_i\}_{i=1}^N$, where N denotes the number of video latents (with $N = 13$ in the case of CogVideoX). Due to the causal design of the video VAE encoder, the latent \mathbf{F}_1 captures information solely from the first frame. The remaining latents $\{\mathbf{F}_i\}_{i=2}^N$ correspond to the future frames to be generated. As illustrated in Figure 2(a), CogVideoX’s backbone is a DiT [28], which takes as input the concatenation of the language tokens \mathcal{T} , the first frame latent \mathbf{F}_1 , and a sequence of noisy future frame latents. All frame latents are flattened in raster scan order. Each token, regardless of its modality, attends to every other token. The model is trained using a DDPM-style [16] diffusion loss applied over $\{\mathbf{F}_i\}_{i=2}^N$. Further details are available in the original CogVideoX paper.

Adaptation to Chunk-Wise Generation. As illustrated in Figure 2(b), our adaptation of the original model to enable chunk-wise generation involves four key modifications:

1. **Chunk-Level Conditioning.** We replace the original first-frame condition with a chunk-level condition, enabling the model to condition on a previously generated video chunk rather than a single frame.
2. **Masking Strategy.** The attention mask is tailored for chunk-wise generation: latents in the conditioning chunk can attend only to the language tokens and themselves, while latents of future frames can attend to all tokens (language tokens, conditioning chunk latents, and themselves).
3. **Temporal Resolution Adjustment.** We reduce the number of future frame latents to be generated from 13 (covering 49 frames in CogVideoX) to 4 (covering 16 frames), enabling finer-grained iterative generation with lower latency.
4. **Noise Augmentation.** During training, the model is conditioned on ground-truth video chunks, whereas at inference time, it relies on its own generated outputs, which may deviate from the true distribution. To address this mismatch, we inject noise into the conditioning inputs during training—a strategy shown to effectively improve robustness [4]. At inference, the model directly uses previously generated latents as conditioning inputs for the current iteration.

All other configurations remain consistent with those of CogVideoX [48], including the use of its backbone, language encoder, video VAE encoder, DDIM-style [35] inference, and the strategy for applying multi-modal positional embeddings. The training in this stage is performed on a general-domain video dataset. Specifically, we select 100K high-quality videos from the OpenViD [26] dataset to serve as the training corpus.

3.2 From Game Environment to Real World

RealPlay is trained on a combination of a labeled game dataset \mathcal{G} and an unlabeled real-world dataset \mathcal{R} , enabling control transfer from operating a car in the game environment to controlling vehicles, bicycles, and pedestrians in real-world scenarios. We fine-tune the chunk-wise generation model introduced in Section 3.1 on both \mathcal{G} and \mathcal{R} , introducing a control module as the only architectural modification to enable action-conditioned chunk-wise generation.

Specifically, given a sample $G_i = \{(C_k, \mathbf{a}_k)\}_{k=1}^K$ from the labeled dataset \mathcal{G} —where each C_k denotes the k -th video chunk and \mathbf{a}_k is a 3-dimensional one-hot vector indicating the action (from the discrete set {move forward, turn left, turn right}) that transitions C_k to C_{k+1} —the control module learns to condition the generation process on the specified action.

As illustrated in Figure 2(c), RealPlay extends the chunk-wise generation architecture described in Section 3.1 by introducing an additional input: the action a_k . This action is injected via an adaptive LayerNorm mechanism. Specifically, a_k is first projected into a 512-D feature vector using an MLP. This vector is then added to the noise timestep embedding, and the resulting feature is passed through a linear layer to produce two sets of modulation parameters: $\{\alpha_1, \gamma_1, \beta_1\}$ and $\{\alpha_2, \gamma_2, \beta_2\}$, where α and γ are scale factors and β is a shift bias. These parameters modulate the network activations before and after the attention and feed-forward layers, respectively, enabling action-aware generation.

For a sample $R_i = \{C_k\}_{k=1}^K$ from the real-world dataset \mathcal{R} , action annotations are unavailable, so we only observe the visual transition from C_k to C_{k+1} . During training, we represent the absence of explicit action supervision by using an all-zero vector $\mathbf{0}$ as the action input.

In summary, during training, we adopt the formulation $C_k + a_k \rightarrow C_{k+1}$ for labeled game samples, while we use $C_k + \mathbf{0} \rightarrow C_{k+1}$ for unlabeled real-world samples. During inference, given a previously generated video chunk C_k depicting a real-world scene (e.g., a moving car, a riding bicycle, or a walking pedestrian), we apply an action command originating from the game environment to generate the next video chunk C_{k+1} , which faithfully follows the specified action.

We now provide an intuitive explanation for the emergence of *game-to-real-world transfer capabilities*—from controlling a car in a game environment to controlling various entities such as vehicles, bicycles, and pedestrians in the real world—through the lenses of classifier-free guidance (CFG) [15] and network generalization. CFG is a technique used in conditional generation models that in the training stage, CFG randomly drops conditioning inputs of a training sample with a pre-defined probability, while during inference it steers outputs toward a desired condition by interpolating between conditional and unconditional predictions.

We begin by introducing RealPlay trained solely on the labeled game dataset \mathcal{G} , with the incorporation of CFG. Under this setting, each sample from \mathcal{G} has a predefined probability of being used with or without its action condition during training. This mechanism is effectively equivalent to partitioning \mathcal{G} into two disjoint subsets: \mathcal{G}_L , which retains action labels, and \mathcal{G}_U , which omits them. RealPlay is jointly trained on both subsets, resulting in controllable generation performance that closely matches the variant trained on \mathcal{G} with conventional CFG.

In our actual scenario, we use the labeled game dataset \mathcal{G} as the counterpart to \mathcal{G}_L , and the unlabeled real-world dataset \mathcal{R} as the counterpart to \mathcal{G}_U . Although \mathcal{G} and \mathcal{R} originate from different distributions, they share important commonalities: (1) both datasets capture entity motion—while \mathcal{G} primarily depicts cars in motion, \mathcal{R} includes additional entities such as bicycles and pedestrians; and (2) the game data used in our setup are derived from Forza Horizon 5, a photorealistic AAA game with high visual realism, which significantly narrows the domain gap between synthetic and real-world data. Combined with the strong generalization capabilities of the large-scale pre-trained video generator, RealPlay enables control transfer from operating a car in the game environment to manipulating diverse entities in the real world.

4 Experiment

Implementation Details. We use CogVideoX-5B [48] as our pre-trained model. First, we fine-tune CogVideoX-5B for chunk-wise generation, as described in Section 3.1, using 100K samples from the OpenViD [26] dataset over 10K iterations on $8 \times$ A100 GPUs. We then further fine-tune the model for 15K iterations, also on $8 \times$ A100 GPUs, following the mixed training strategy outlined in Section 3.2. This stage uses a combined dataset comprising a labeled game dataset from The Matrix [9]—which contains 80K samples from the car racing game Forza Horizon 5—and our own collected unlabeled real-world dataset, consisting of 18K samples each for vehicles, bicycles, and pedestrians. Each game sample is a 32-frame video clip annotated with an action label that controls the transition from the first 16 frames to the next 16. In contrast, each real-world sample is an unannotated 32-frame clip, capturing unconstrained motion of entities. Additional training details are provided in the appendix.

Evaluation. Our evaluation covers three key aspects:

1. **Visual Quality.** We adopt four metrics proposed by VBench [19]—motion consistency, aesthetic appeal, imaging quality, and scene dynamics—to quantitatively evaluate the visual quality of the generated video chunks.
2. **Control Effectiveness.** We report the success rate based on human evaluation, where evaluators judge whether the generated video accurately reflects the intended control command.

Table 1: Comparison of *RealPlay variants* with: (1) *Single-forward-pass models*, which generate an entire video sequence in a single forward pass; (2) *Chunk-wise generation models*, where CogVideoX-5B is first adapted to a chunk-wise generator using the approach described in Section 3.1, and then fine-tuned on either manually labeled real-world data or pseudo-labeled data generated by LAPA [49]. RealPlay-AdaLN: uses Adaptive LayerNorm to fuse action signals. RealPlay-Text: actions are expressed in text prompt. ‘L’ and “U” denote labeled data and unlabeled data, respectively.

Method	Real Data	Game Data	Motion	Aesthetic	Visual Imaging	Quality Dynamic	Control Rate	Elo
<i>Single-Forward-Pass Models:</i>								
CogVideoX-5B [48]			98.3	47.4	71.3	49.4	33.9	1025
Hunyuan-720P [22]	-	-	99.3	49.1	91.9	52.8	31.7	1102
Wan-2.1 [41]	-	-	97.9	47.5	72.9	61.1	32.2	929
OpenSora-2.0 [29]			99.2	46.9	69.5	100.0	26.7	795
<i>Chunk-Wise Generation Models:</i>								
CogVideoX-5B (Human)	L	-	97.9	46.7	69.7	98.5	58.9	1045
CogVideoX-5B (LAPA [49])	U	-	98.7	44.5	64.3	62.3	36.2	779
<i>RealPlay Variants:</i>								
RealPlay-Text	U	L	97.6	46.3	69.0	99.4	69.9	1143
RealPlay-AdaLN (Default)	U	L	97.5	47.8	68.9	100.0	90.0	1184

3. **Comprehensive Human Evaluation.** To holistically evaluate model performance, we compute an Elo score for each model using 500 pairwise comparisons. In each comparison, two video outputs—generated from the same initial observation and control condition by different models—are displayed side by side. Human evaluators are asked to select the video that appears more realistic and better aligned with the intended control. Elo scores are updated using a standard rating adjustment scheme, where the winner gains points from the loser. All models are initialized with an Elo score of 1000.

Unless otherwise specified, we iteratively generate 3 chunks using 3 control commands, each randomly selected from the set {move forward, turn left, turn right}, for comparisons with other models and for ablation studies. We evaluate the performance on vehicles, bicycles, and pedestrians that were observed during training.

4.1 Main Results

In Table 1, we compare our RealPlay model with several baselines, which can be categorized into four classes. All models receive the same initial visual input, but differ in their approach to action control and the nature of their training data:

1. **Single-Forward-Pass Models.** We directly evaluate four state-of-the-art pre-trained (text,image)-to-video models—CogVideoX-5B [48], Hunyuan-720P [22], Wan-2.1 [41], and OpenSora-2.0 [29]—without any fine-tuning. The text prompt is “Control the [Entity] to first [Action-1], then [Action-2], and finally [Action-3]”, where Entity $\in \{\text{vehicle, bicycle, pedestrian}\}$, and Action-1,2,3 $\in \{\text{move forward, turn left, turn right}\}$.
2. **Chunk-Wise Generation Models Trained on Labeled Real-World Data.** We adapt the original CogVideoX-5B [48] into a chunk-wise generation model using the method described in Section 3.1, and fine-tune it on our collected real-world dataset using a selective annotation strategy. Specifically, we manually label only clear samples and exclude ambiguous ones to ensure annotation quality, ultimately selecting approximately 20% of the total data. Each training sample is a short video clip consisting of two temporally continuous chunks, denoted as C_1 and C_2 . A sample is considered clear if, given C_1 , human annotators can confidently determine whether a single action a leads to C_2 (e.g., moving forward throughout C_2). In contrast, ambiguous samples contain multiple actions within C_2 , making it difficult to assign a single action label. Using this filtering strategy, we construct a training set in which we have reliable transitions of the form $C_1 + a \rightarrow C_2$.
3. **Chunk-Wise Generation Models Trained on Pseudo-Labeled Real-World Data.** This model is identical to the one described in point 2, except for the annotation strategy applied to the real-world data. Instead of relying on manual labels, we use LAPA [49] to generate pseudo-action labels a' . Details of the LAPA method are provided in the appendix. Using these pseudo labels, we construct a training set with transitions of the form $C_1 + a' \rightarrow C_2$.

Table 2: Per-entity evaluation shows that real-world entities with larger motion amplitudes achieve higher control success rates, as their more distinctive motion patterns make it easier for the network to transfer control from the game environment to the real world.

Entity	Visual Quality				Control Rate	Motion Amplitude
	Motion	Aesthetic	Imaging	Dynamic		
In-Game Car	97.8	52.5	70.8	100.0	100.0	3.3
Vehicle	97.1	52.6	71.0	100.0	83.3	1.9
Bicycle	97.1	51.2	66.8	100.0	91.7	2.8
Pedestrian	98.4	39.7	68.8	100.0	95.0	5.7

Table 3: Cross-entity training significantly improves control success for individual entities. The ablation study is conducted on the *bicycle* entity.

Game Data	Real-World Data			Visual Quality				Control Rate
	Bicycle	Vehicle	Pedestrian	Motion	Aesthetic	Imaging	Dynamic	
✓				96.1	50.5	59.1	100.0	0.0
✓	✓			97.1	50.3	67.0	100.0	72.5
✓	✓	✓		96.7	52.7	66.7	100.0	78.4
✓	✓	✓	✓	97.1	51.2	66.8	100.0	91.7

4. **RealPlay.** We evaluate two variants of RealPlay: (1) *RealPlay-AdaLN*: the default version, which uses Adaptive LayerNorm to fuse action signals. It is trained on both labeled game data and unlabeled real-world data. (2) *RealPlay-Text*: trained with the same data as RealPlay-AdaLN, but action signals are reflected in the form of text instructions. Specifically, we use the template: “Control the car to [Action]” where Action $\in \{\text{move forward}, \text{turn left}, \text{turn right}\}$.

From the results presented in Table 1, we draw several key conclusions: (1) All single-forward-pass models exhibit some level of control capability through text prompts. However, their control success rates remain low, ranging from 26.7% to 33.9%. We choose CogVideoX-5B [18] as our pre-trained model due to its highest zero-shot control performance and acceptable visual quality. Additionally, we observe that models such as CogVideoX-5B [48], Hunyuan-720P [22], and Wan-2.1 [41] sometimes generate static videos, which contributes to lower dynamic scores. (2) The two chunk-wise generation models, fine-tuned on a human-labeled dataset and the LAPA [49] pseudo-labeled dataset respectively, both show improvements in control success rate. However, the overall success rates remain low. This is primarily due to two factors: the limited size of the human-labeled dataset, which is insufficient to effectively train a large network, and the noise present in the pseudo labels, which hampers learning in the LAPA-fine-tuned model. (3) Both RealPlay variants significantly improve the control success rate. Our default implementation, which uses adaptive LayerNorm for action fusion, achieves a success rate of 90%. Additionally, neither variant generates static videos; instead, they produce more visually appealing results than the original CogVideoX-5B, while introducing two key enhancements: chunk-wise generation and control over real-world entities in the video.

4.2 Analysis

Larger Motion Amplitudes Yield Higher Control Rates. RealPlay is trained using a labeled game dataset (car racing) and an unlabeled real-world dataset containing vehicles, bicycles, and pedestrians. In Table 2, we report the visual quality and control rate for each entity. A key observation is that control rates vary across real-world entities, ranging from 83.3% to 95.0%. We attribute this variation to differences in motion amplitude across real-world entities. Specifically, real-world pedestrians typically exhibit larger motion amplitudes than bicycles and vehicles—especially during dynamic movements such as turning—whereas vehicles tend to move more smoothly and slowly. Our study reveals that a larger motion amplitude leads to a higher control rate, as it results in more distinctive motion patterns—making it easier for the network to transfer control from the game environment to the real world. Details on computing motion amplitudes are provided in the appendix.

Cross-Entity Training Improves Control Success for Individual Entities. In Table 3, we evaluate the effect of cross-entity training by training our model on game data combined with real-world data from various entities—bicycles, vehicles, and pedestrians—and testing specifically on bicycles. We make three key observations: (1) RealPlay trained solely on game data fails to control the real-world entity (bicycle); in fact, during generation, the bicycle gradually transforms into a game-style car due to the model’s lack of exposure to real-world entity distributions. (2) While labeled game data provides essential supervision, the inclusion of unlabeled real-world data plays a critical role by

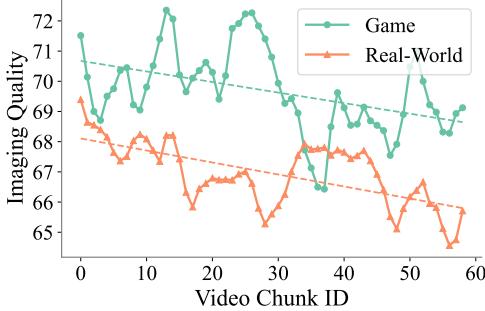


Figure 3: Visual quality degrades in both game and real-world settings, but the image quality when controlling a game entity consistently remains higher than that of a real-world entity (e.g., the bicycle in this study), highlighting the greater challenge of modeling real-world entities.

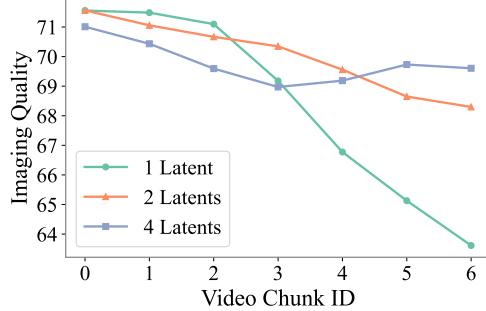


Figure 4: Reducing the number of video latents per chunk leads to visual quality degradation, as the pre-trained video generator—originally optimized for long-horizon generation—loses temporal coherence and consistency when adapted to extremely short-horizon outputs (e.g., 1 latent).

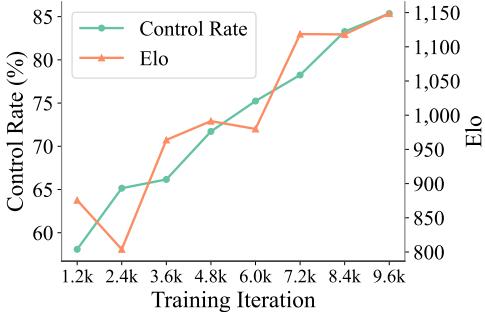


Figure 5: Both the control success rate and Elo scores steadily improve as training progresses. The evaluation is performed on the bicycle entity.

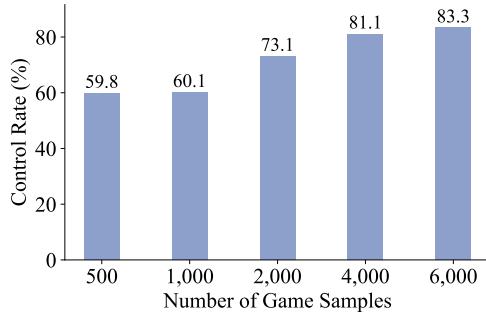


Figure 6: Increasing the number of labeled game samples consistently improves control success on real-world entities.

supplying realistic motion patterns as well as the appearance of entities. As a result, training RealPlay on a combination of game data and bicycle-specific real-world data significantly boosts the control rate to 72.5%. (3) The additional involvement of other real-world cross-entity data (i.e., vehicles and pedestrians) during training enhances the control of bicycles. This improvement arises because cross-entity data introduces similar motion patterns and shares underlying motion dynamics.

Real-World Environments Are Significantly Harder to Fit Than Virtual Games. As a chunk-wise generation model, RealPlay generates each video chunk conditioned on the previously generated chunk and an action signal. This design introduces an error accumulation issue, as the model conditions on its own generated outputs rather than ground-truth data. We analyze how this affects long-term generation in both the game and real-world settings. Figure 3 illustrates the degradation of image quality as the number of generated video chunks increases.

The Number of Video Latents per Chunk Affects Visual Quality. At each iteration, RealPlay generates a video chunk composed of 4 video latents (equivalent to 16 frames). To enhance user experience with low control latency, it is desirable to minimize the number of frames generated per iteration. However, we observe that generating too few video latents—such as 1 latent (4 frames) per iteration—leads to a noticeable drop in visual quality, as illustrated in Figure 4. This degradation arises because RealPlay leverages a large-scale pre-trained video generation model, CogVideoX, which is originally trained to generate 13 video latents in a single forward pass. Adapting such a model to extremely low-latency, short-horizon generation (e.g. 1 video latent) disrupts its learned temporal coherence and internal consistency, resulting in lower-quality outputs.

Co-Improvement of Visual Quality and Control Rate During Training. Figure 5 illustrates how both the control success rate and Elo scores—which jointly reflect control performance and visual quality based on human evaluation—consistently improve as training progresses.

More Labeled Game Samples Improve Control Transfer. We investigate how the amount of labeled game data affects stable control transfer. Specifically, we use a fixed subset of the real-world training set containing 4K samples, and vary the amount of labeled game data. The control success rate for each setting is reported in Figure 6.

5 Conclusion

We present RealPlay, a neural network-based real-world game engine that enables interactive, photorealistic video generation conditioned on user control signals. Unlike traditional game engines limited by handcrafted graphics and physics rules, RealPlay leverages a chunk-wise generation framework to support iterative interaction with strong temporal consistency. Through a mixed training paradigm—combining labeled game data with unlabeled real-world videos—RealPlay achieves effective control transfer from virtual to real-world settings. Remarkably, it generalizes to diverse real-world entities, despite being trained only with car-based game supervision. Our results demonstrate RealPlay’s potential to bridge simulation and reality, marking a first step toward neural game engines that learn to simulate the real world from data.

A More Implementation Details

More Training Details. Both training stages share the same optimization settings: we use the AdamW optimizer with a learning rate of 1e-5 and a weight decay of 1e-3. A DDPM-style diffusion loss is employed with the default noise scheduler, and training is conducted on 8 NVIDIA A100 (40G) GPUs. The chunk-wise generation model (Stage 1) is trained for 10K iterations, while the final model (Stage 2) is fine-tuned for 15K iterations. During inference, we adopt DDIM with 50 denoising steps.

LAPA. In Table 1 of the main paper, we present a baseline named CogVideoX-5B (LAPA[49]). LAPA is an unsupervised model that learns a pseudo-action a' to enable the transition $C_1 + a' \rightarrow C_2$. The definitions of C_1 and C_2 are provided in Section 4.1. LAPA adopts a VQ-VAE [40] architecture, where the objective is to reconstruct C_2 given C_1 and a discrete action dictionary that facilitates the reconstruction. In our implementation, we set the dictionary size to 3. We first fine-tune LAPA on our real-world dataset, and then use it to generate a pseudo-action label a' for each training sample. These pseudo-labeled samples are then used to finetune the chunk-wise generation model described in Section 3.1.

Motion Amplitude Computation. To compute motion amplitude, we first use SAM-2 [31] to segment the entities in each training video, and then apply RAFT [38] to estimate the optical flow. The motion amplitude is defined as the average optical flow magnitude across all training videos for a given entity.

B More Experiments

Table 4 compares three action signal injection strategies: (1) the default adaptive LayerNorm described in Section 3.2; (2) self-attention, where the action signal is treated as a token and appended to the input sequence after the language tokens; and (3) cross-attention, where a cross-attention layer is added after the self-attention layer in each attention block, using the output of the self-attention layer as the query and the action token as the key.

Table 4: Comparison of various action signal injection strategies.

Strategy	Visual Quality				Control Success Rate (%)
	Motion	Aesthetic	Imaging	Dynamic	
Self-Attention	97.7	47.3	69.2	97.5	78.3
Cross-Attention	97.6	47.4	68.8	100.0	77.3
Adaptive LayerNorm (Default)	97.5	47.8	68.9	100.0	90.0

C More Visualizations

Visualizations of Videos Generated by the Chunk-Wise Generation Model. Section 3.1 details how a pre-trained image-to-video generator is adapted into a chunk-wise generation model. Trained on a general-domain dataset, this model is capable of generating diverse videos in a chunk-wise manner. Figure 7 showcases several examples produced by our chunk-wise generation model.

Visualizations of Videos Generated by RealPlay. RealPlay is trained on a combination of a labeled car-racing game dataset and an unlabeled real-world dataset containing three types of entities: bicycles, pedestrians, and vehicles. This setup enables control transfer from the game environment to the real world. Figure 8 demonstrates RealPlay’s ability to control cars across diverse game scenes, while Figure 9 illustrates its effectiveness in controlling various real-world entities. Additionally, Figure 10 presents six long-horizon videos generated by RealPlay, demonstrating its capability in sustaining control over real-world entities.

What Happens When Control Signals Are Applied to Videos Without a Clear Entity Focus? RealPlay enables control over real-world entities to move forward, turn left, and turn right. In typical videos, the target entity is centrally positioned within the frame. Notably, when applying these control signals, the camera is not static—instead, it dynamically follows the trajectory of the entity. In Figure 11, we visualize several examples where there is no clear entity focus. Interestingly,

we observe that when control signals are applied to these examples, the camera itself moves in accordance with the control direction. This suggests that RealPlay learns to transfer both entity and camera control: when a focused entity is present in the frame, it controls both the entity and the camera; when no focused entity is present, it controls only the camera. In both cases, the results are visually coherent and pleasant.

Qualitative Comparison with Baseline Methods. We provide a qualitative comparison with baseline methods listed in Table 1 of our main paper. As shown in Figure 12 and Figure 13, in comparison to baselines, our approach achieves more precise control while maintaining superior visual quality.

D Limitations and Broader Impacts

This work represents an initial step toward leveraging video generation techniques to build real-world games, where visuals are photorealistic and the game engine is powered by a data-driven neural network. RealPlay depends on a pre-trained general-purpose video generator—specifically, CogVideoX-5B in this study. Due to the large model size and the relatively high resolution of the generated videos, RealPlay cannot yet operate in real-time. Several techniques could potentially accelerate inference, such as model distillation and few-step denoising methods like the Shortcut [10] model. We consider these as promising directions for future optimization and deployment.

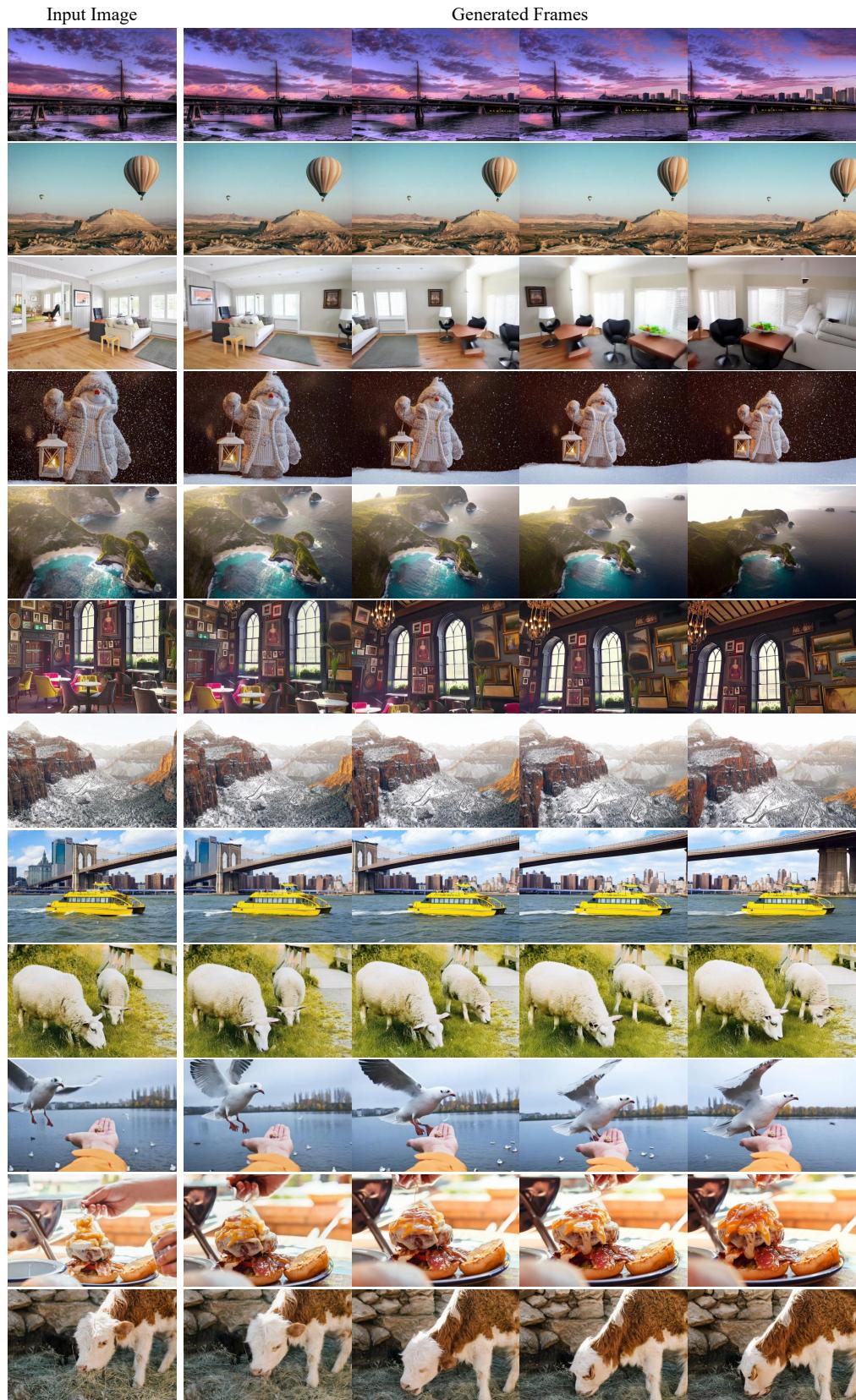


Figure 7: Visualizations of videos generated by the chunk-wise generation model.

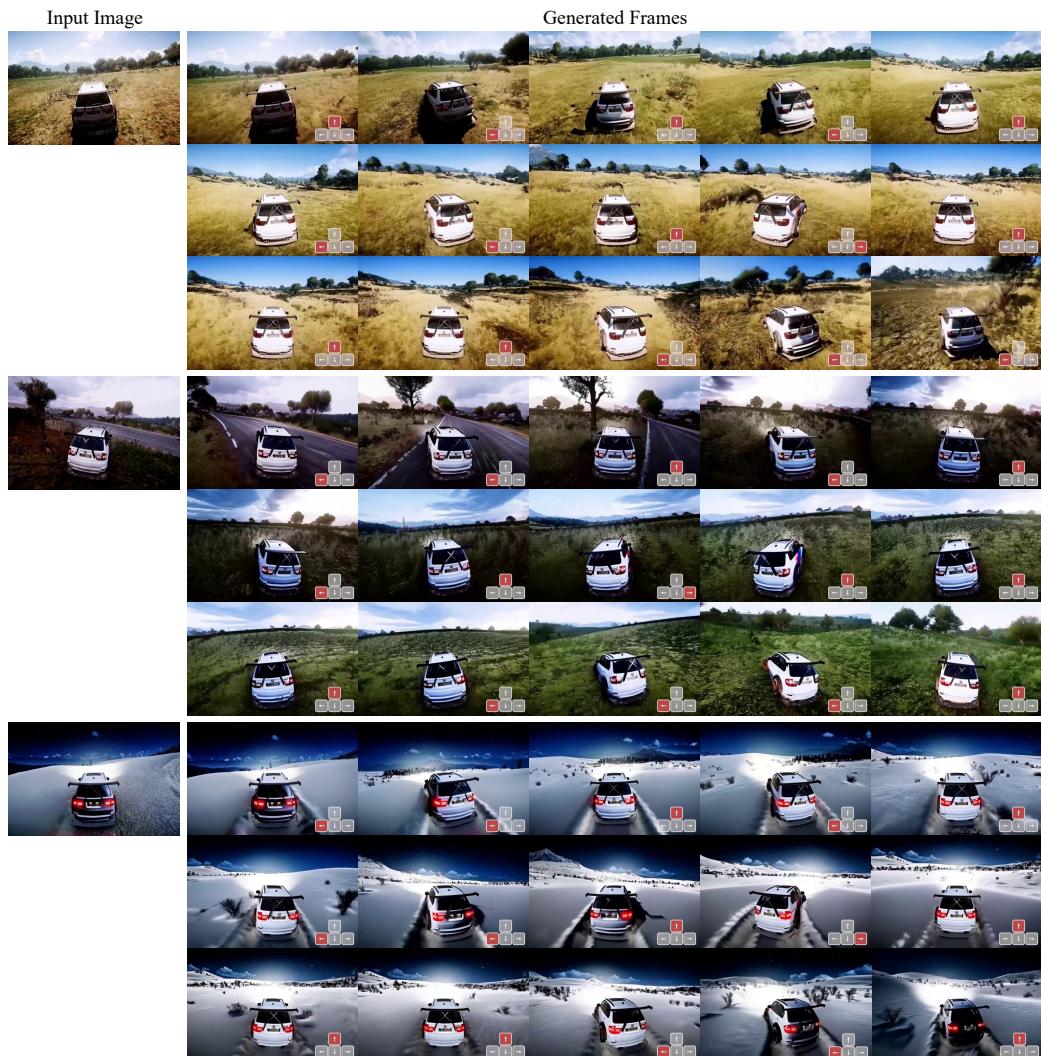


Figure 8: Visualization of RealPlay’s capability to control cars in diverse game environments.

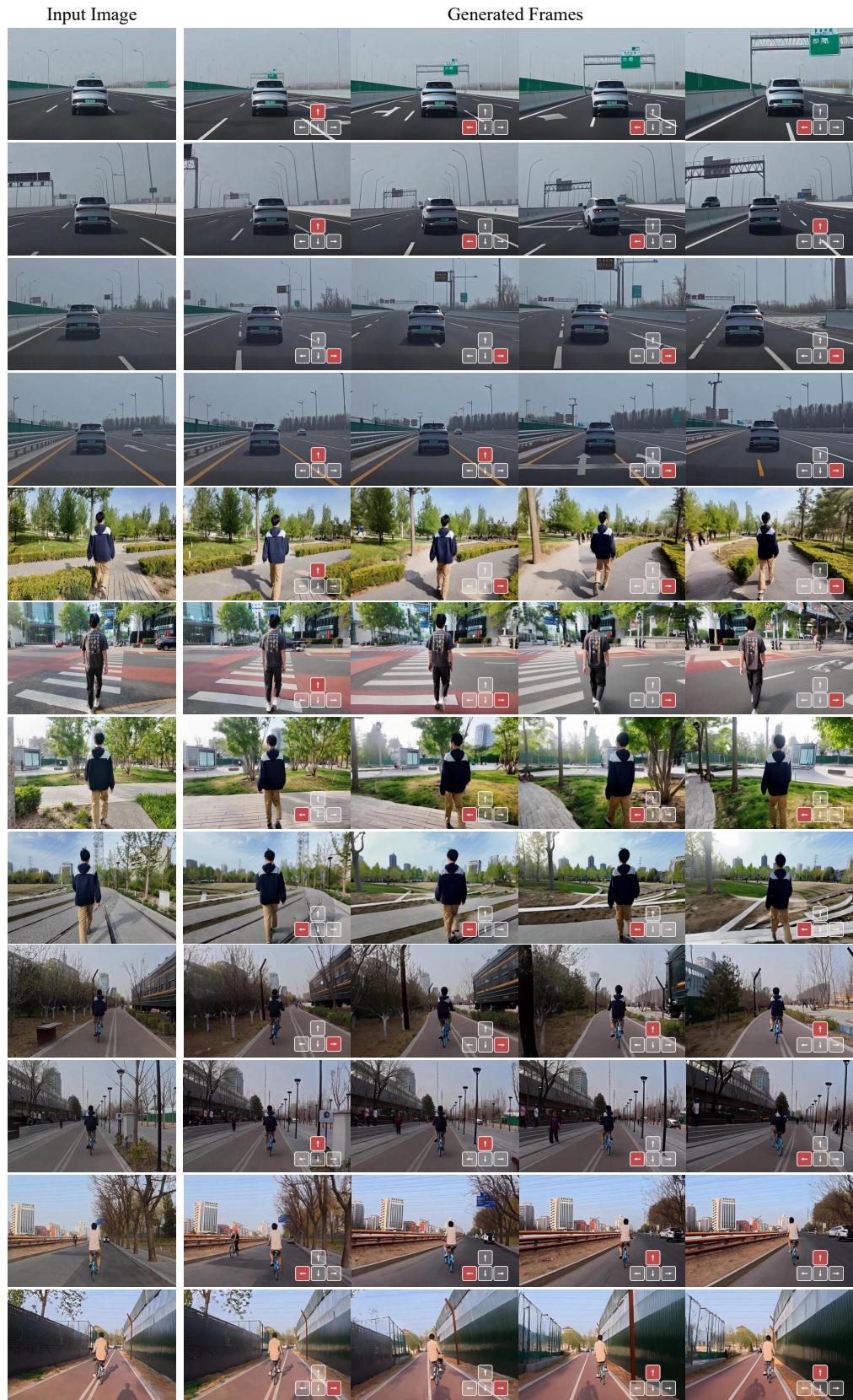


Figure 9: RealPlay’s effectiveness in controlling diverse real-world entities.

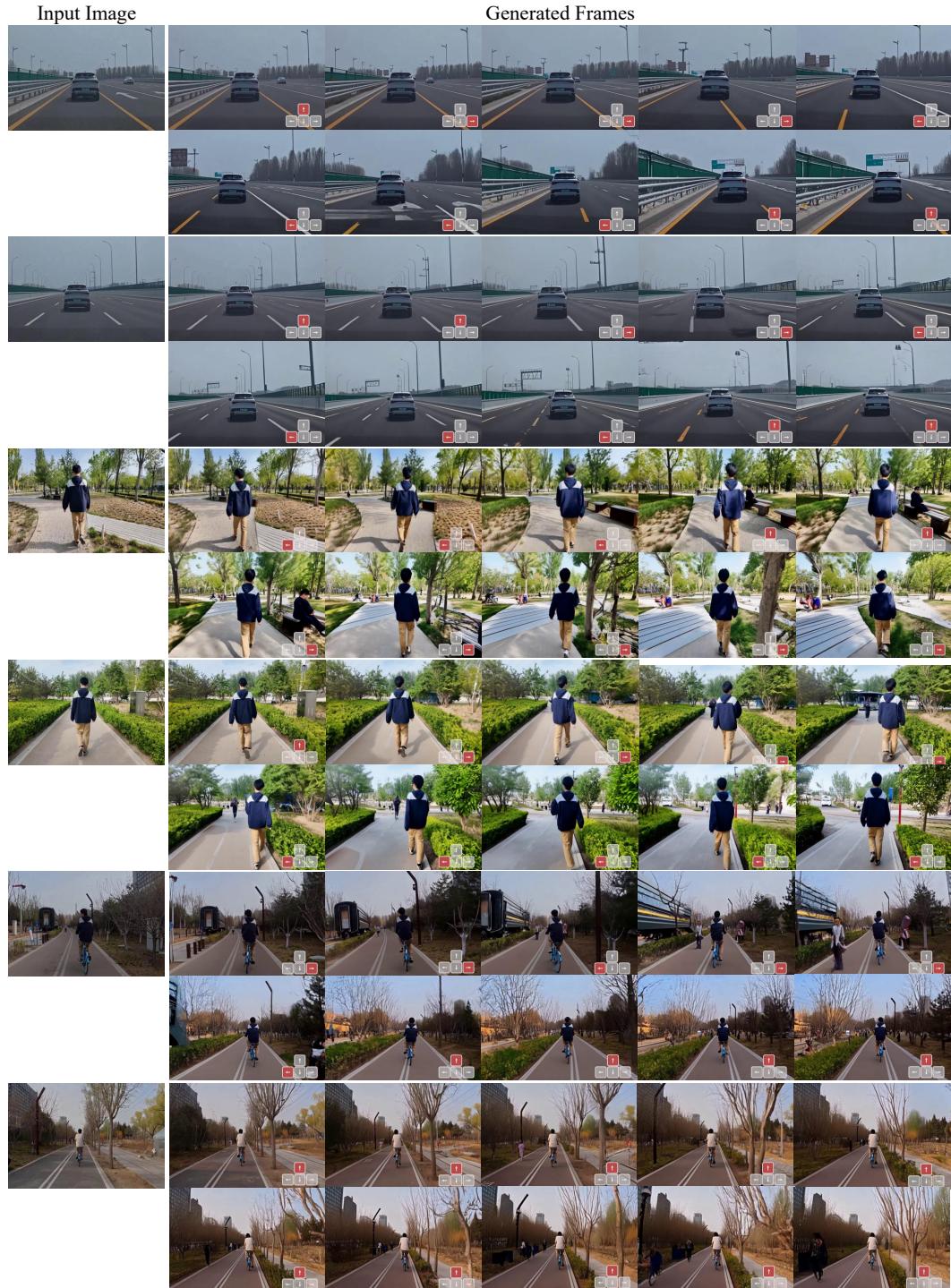


Figure 10: Long-horizon videos generated by RealPlay demonstrating sustained control over real-world entities.

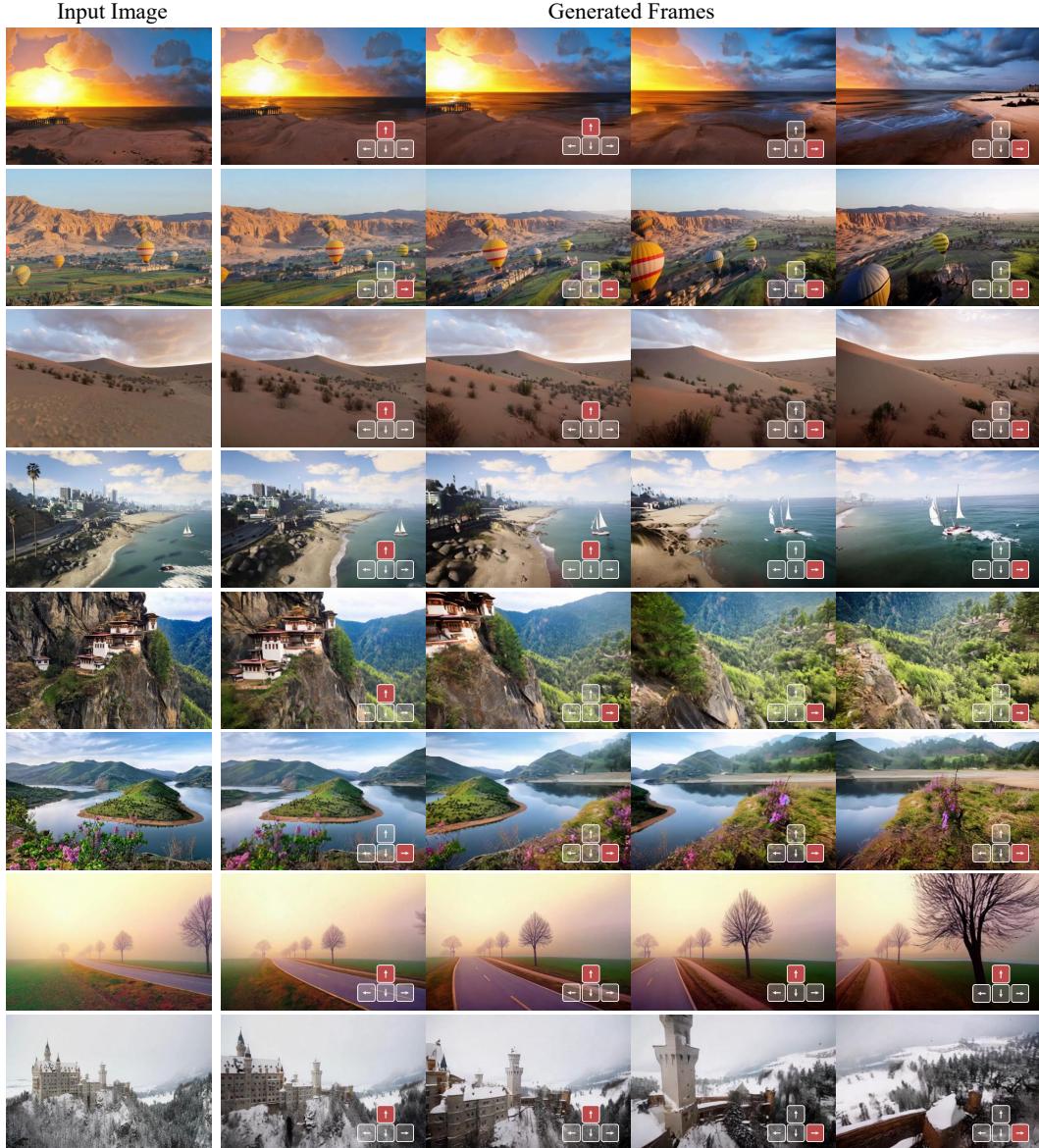


Figure 11: RealPlay learns both entity and camera control; when no focused entity is present, it controls the camera alone.

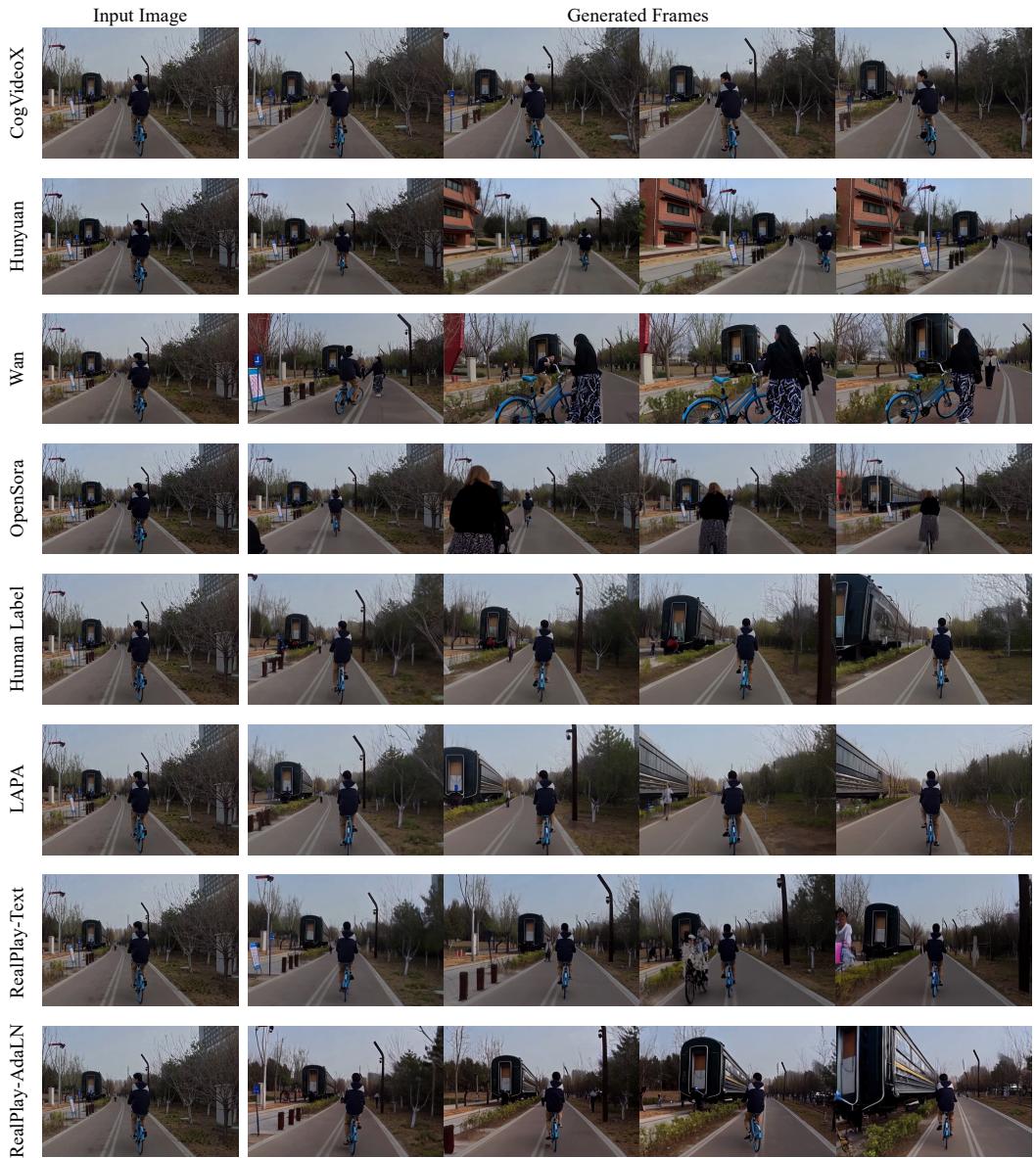


Figure 12: Qualitative comparison with baseline methods. The target action is **move left and then forward**.

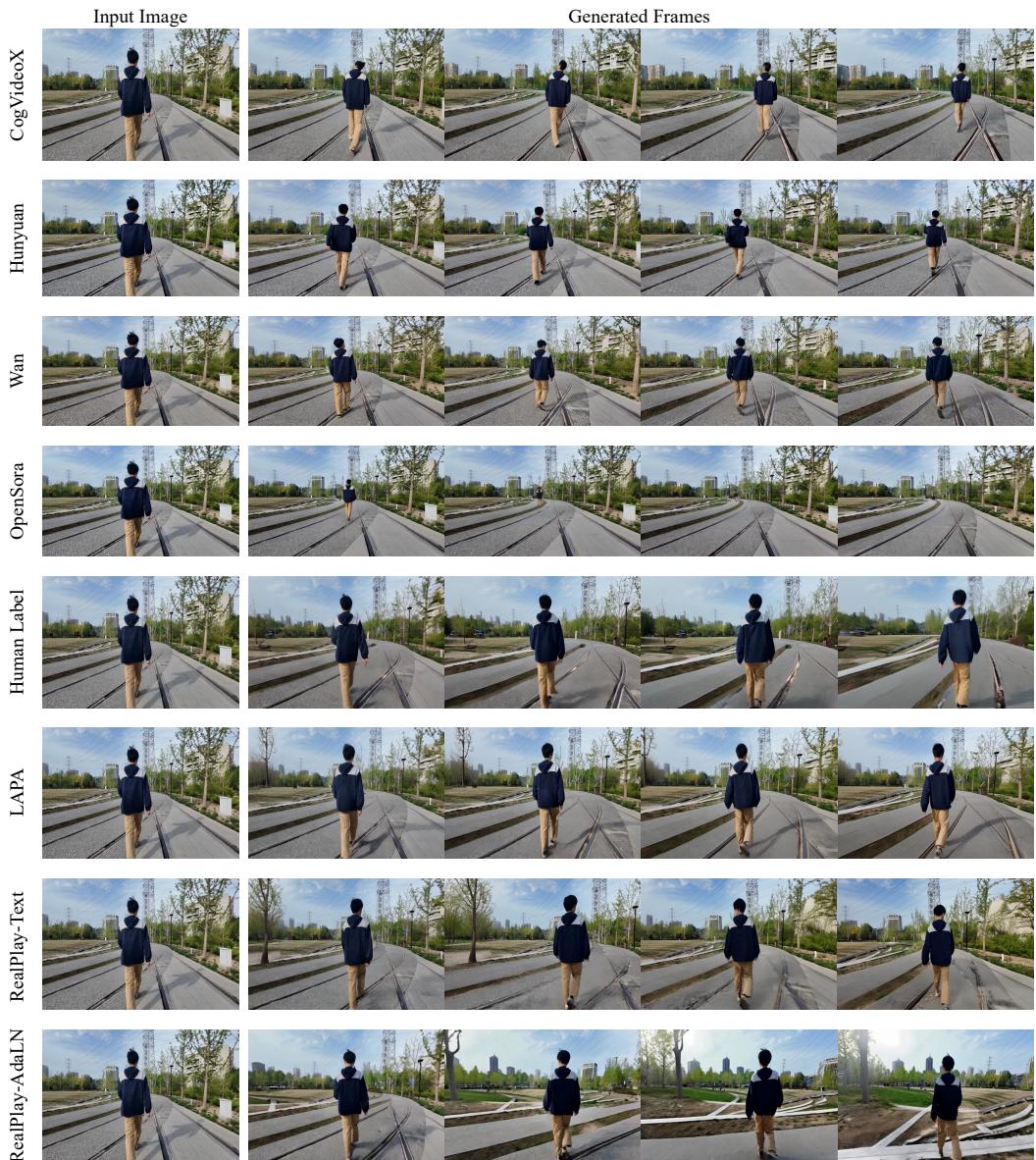


Figure 13: Qualitative comparison with baseline methods. The target action is **move left**.

References

- [1] E. Alonso, A. Jelley, V. Micheli, A. Kanervisto, A. J. Storkey, T. Pearce, and F. Fleuret. Diffusion for world modeling: Visual details matter in atari. *Advances in Neural Information Processing Systems*, 37:58757–58791, 2024.
- [2] J. Bruce, M. D. Dennis, A. Edwards, J. Parker-Holder, Y. Shi, E. Hughes, M. Lai, A. Mavalankar, R. Steigerwald, C. Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024.
- [3] H. Che, X. He, Q. Liu, C. Jin, and H. Chen. Gamegen-x: Interactive open-world game video generation. *arXiv preprint arXiv:2411.00769*, 2024.
- [4] B. Chen, D. Martí Monsó, Y. Du, M. Simchowitz, R. Tedrake, and V. Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion. *Advances in Neural Information Processing Systems*, 37:24081–24125, 2024.
- [5] G. Chen, D. Lin, J. Yang, C. Lin, J. Zhu, M. Fan, H. Zhang, S. Chen, Z. Chen, C. Ma, et al. Skyreels-v2: Infinite-length film generative model. *arXiv preprint arXiv:2504.13074*, 2025.
- [6] H. Chen, M. Xia, Y. He, Y. Zhang, X. Cun, S. Yang, J. Xing, Y. Liu, Q. Chen, X. Wang, et al. Videocrafter1: Open diffusion models for high-quality video generation. *arXiv preprint arXiv:2310.19512*, 2023.
- [7] H. Chen, Y. Zhang, X. Cun, M. Xia, X. Wang, C. Weng, and Y. Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7310–7320, 2024.
- [8] Decart, J. Quevedo, Q. McIntyre, S. Campbell, X. Chen, and R. Wachen. Oasis: A universe in a transformer. 2024. URL <https://oasis-model.github.io/>.
- [9] R. Feng, H. Zhang, Z. Yang, J. Xiao, Z. Shu, Z. Liu, A. Zheng, Y. Huang, Y. Liu, and H. Zhang. The matrix: Infinite-horizon world generation with real-time moving control. *arXiv preprint arXiv:2412.03568*, 2024.
- [10] K. Frans, D. Hafner, S. Levine, and P. Abbeel. One step diffusion via shortcut models. *arXiv preprint arXiv:2410.12557*, 2024.
- [11] K. Gao, J. Shi, H. Zhang, C. Wang, and J. Xiao. Vid-gpt: Introducing gpt-style autoregressive generation in video diffusion models. *arXiv preprint arXiv:2406.10981*, 2024.
- [12] Y. Gu, W. Mao, and M. Z. Shou. Long-context autoregressive video modeling with next-frame prediction. *arXiv preprint arXiv:2503.19325*, 2025.
- [13] Y. Guo, C. Yang, A. Rao, Z. Liang, Y. Wang, Y. Qiao, M. Agrawala, D. Lin, and B. Dai. Animated-iff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023.
- [14] R. Henschel, L. Khachatryan, H. Poghosyan, D. Hayrapetyan, V. Tadevosyan, Z. Wang, S. Navasardyan, and H. Shi. Streaming2v: Consistent, dynamic, and extendable long video generation from text. *arXiv preprint arXiv:2403.14773*, 2024.
- [15] J. Ho and T. Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [16] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [17] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M. Norouzi, and D. J. Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.
- [18] W. Hong, M. Ding, W. Zheng, X. Liu, and J. Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.
- [19] Z. Huang, Y. He, J. Yu, F. Zhang, C. Si, Y. Jiang, Y. Zhang, T. Wu, Q. Jin, N. 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.
- [20] Y. Jin, Z. Sun, N. Li, K. Xu, H. Jiang, N. Zhuang, Q. Huang, Y. Song, Y. Mu, and Z. Lin. Pyramidal flow matching for efficient video generative modeling. *arXiv preprint arXiv:2410.05954*, 2024.

- [21] J. Kim, J. Kang, J. Choi, and B. Han. Fifo-diffusion: Generating infinite videos from text without training. *arXiv preprint arXiv:2405.11473*, 2024.
- [22] W. Kong, Q. Tian, Z. Zhang, R. Min, Z. Dai, J. Zhou, J. Xiong, X. Li, B. Wu, J. Zhang, et al. Hunyanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024.
- [23] B. Lin, Y. Ge, X. Cheng, Z. Li, B. Zhu, S. Wang, X. He, Y. Ye, S. Yuan, L. Chen, et al. Open-sora plan: Open-source large video generation model. *arXiv preprint arXiv:2412.00131*, 2024.
- [24] F. Liu, W. Sun, H. Wang, Y. Wang, H. Sun, J. Ye, J. Zhang, and Y. Duan. Reconx: Reconstruct any scene from sparse views with video diffusion model. *arXiv preprint arXiv:2408.16767*, 2024.
- [25] G. Ma, H. Huang, K. Yan, L. Chen, N. Duan, S. Yin, C. Wan, R. Ming, X. Song, X. Chen, et al. Step-video-t2v technical report: The practice, challenges, and future of video foundation model. *arXiv preprint arXiv:2502.10248*, 2025.
- [26] K. Nan, R. Xie, P. Zhou, T. Fan, Z. Yang, Z. Chen, X. Li, J. Yang, and Y. Tai. Openvid-1m: A large-scale high-quality dataset for text-to-video generation. *arXiv preprint arXiv:2407.02371*, 2024.
- [27] J. Parker-Holder, P. Ball, J. Bruce, V. Dasagi, K. Holsheimer, C. Kaplanis, A. Moufarek, G. Scully, J. Shar, J. Shi, S. Spencer, J. Yung, M. Dennis, S. Kenjeyev, S. Long, V. Mnih, H. Chan, M. Gazeau, B. Li, F. Pardo, L. Wang, L. Zhang, F. Besse, T. Harley, A. Mitenkova, J. Wang, J. Clune, D. Hassabis, R. Hadsell, A. Bolton, S. Singh, and T. Rocktäschel. Genie 2: A large-scale foundation world model. 2024. URL <https://deepmind.google/discover/blog/genie-2-a-large-scale-foundation-world-model/>.
- [28] W. Peebles and S. Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023.
- [29] X. Peng, Z. Zheng, C. Shen, T. Young, X. Guo, B. Wang, H. Xu, H. Liu, M. Jiang, W. Li, et al. Open-sora 2.0: Training a commercial-level video generation model in \$200 k. *arXiv preprint arXiv:2503.09642*, 2025.
- [30] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- [31] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland, L. Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- [32] Sand-AI. Magi-1: Autoregressive video generation at scale, 2025. URL https://static.magi.world/static/files/MAGI_1.pdf.
- [33] A. Siarohin, S. Lathuilière, S. Tulyakov, E. Ricci, and N. Sebe. First order motion model for image animation. *Advances in neural information processing systems*, 32, 2019.
- [34] U. Singer, A. Polyak, T. Hayes, X. Yin, J. An, S. Zhang, Q. Hu, H. Yang, O. Ashual, O. Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022.
- [35] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [36] K. Song, B. Chen, M. Simchowitz, Y. Du, R. Tedrake, and V. Sitzmann. History-guided video diffusion. *arXiv preprint arXiv:2502.06764*, 2025.
- [37] W. Sun, S. Chen, F. Liu, Z. Chen, Y. Duan, J. Zhang, and Y. Wang. Dimensionx: Create any 3d and 4d scenes from a single image with controllable video diffusion. *arXiv preprint arXiv:2411.04928*, 2024.
- [38] Z. Teed and J. 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.
- [39] D. Vlevski, Y. Leviathan, M. Arar, and S. Fruchter. Diffusion models are real-time game engines. *arXiv preprint arXiv:2408.14837*, 2024.
- [40] A. Van Den Oord, O. Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.

- [41] A. Wang, B. Ai, B. Wen, C. Mao, C.-W. Xie, D. Chen, F. Yu, H. Zhao, J. Yang, J. Zeng, J. Wang, J. Zhang, J. Zhou, J. Wang, J. Chen, K. Zhu, K. Zhao, K. Yan, L. Huang, M. Feng, N. Zhang, P. Li, P. Wu, R. Chu, R. Feng, S. Zhang, S. Sun, T. Fang, T. Wang, T. Gui, T. Weng, T. Shen, W. Lin, W. Wang, W. Wang, W. Zhou, W. Wang, W. Shen, W. Yu, X. Shi, X. Huang, X. Xu, Y. Kou, Y. Lv, Y. Li, Y. Liu, Y. Wang, Y. Zhang, Y. Huang, Y. Li, Y. Wu, Y. Liu, Y. Pan, Y. Zheng, Y. Hong, Y. Shi, Y. Feng, Z. Jiang, Z. Han, Z.-F. Wu, and Z. Liu. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025.
- [42] C. Wu, J. Liang, X. Hu, Z. Gan, J. Wang, L. Wang, Z. Liu, Y. Fang, and N. Duan. Nuwa-infinity: Autoregressive over autoregressive generation for infinite visual synthesis. *arXiv preprint arXiv:2207.09814*, 2022.
- [43] T. Wu, Z. Fan, X. Liu, H.-T. Zheng, Y. Gong, J. Jiao, J. Li, J. Guo, N. Duan, W. Chen, et al. Ar-diffusion: Auto-regressive diffusion model for text generation. *Advances in Neural Information Processing Systems*, 36:39957–39974, 2023.
- [44] Z. Xiao, Y. Lan, Y. Zhou, W. Ouyang, S. Yang, Y. Zeng, and X. Pan. Worldmem: Long-term consistent world simulation with memory. *arXiv preprint arXiv:2504.12369*, 2025.
- [45] D. Xie, Z. Xu, Y. Hong, H. Tan, D. Liu, F. Liu, A. Kaufman, and Y. Zhou. Progressive autoregressive video diffusion models. *arXiv preprint arXiv:2410.08151*, 2024.
- [46] J. Xing, M. Xia, Y. Zhang, H. Chen, W. Yu, H. Liu, G. Liu, X. Wang, Y. Shan, and T.-T. Wong. Dynamicroft: Animating open-domain images with video diffusion priors. In *European Conference on Computer Vision*, pages 399–417. Springer, 2024.
- [47] M. Yang, J. Li, Z. Fang, S. Chen, Y. Yu, Q. Fu, W. Yang, and D. Ye. Playable game generation. *arXiv preprint arXiv:2412.00887*, 2024.
- [48] Z. Yang, J. Teng, W. Zheng, M. Ding, S. Huang, J. Xu, Y. Yang, W. Hong, X. Zhang, G. Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024.
- [49] S. Ye, J. Jang, B. Jeon, S. Joo, J. Yang, B. Peng, A. Mandlekar, R. Tan, Y.-W. Chao, B. Y. Lin, et al. Latent action pretraining from videos. *arXiv preprint arXiv:2410.11758*, 2024.
- [50] T. Yin, Q. Zhang, R. Zhang, W. T. Freeman, F. Durand, E. Shechtman, and X. Huang. From slow bidirectional to fast autoregressive video diffusion models. *arXiv preprint arXiv:2412.07772*, 2, 2024.
- [51] J. Yu, Y. Qin, X. Wang, P. Wan, D. Zhang, and X. Liu. Gamefactory: Creating new games with generative interactive videos. *arXiv preprint arXiv:2501.08325*, 2025.
- [52] Y. Zhang, C. Peng, B. Wang, P. Wang, Q. Zhu, Z. Gao, E. Li, Y. Liu, and Y. Zhou. Matrix-game: Interactive world foundation model. *arXiv*, 2025.
- [53] Z. Zheng, X. Peng, T. Yang, C. Shen, S. Li, H. Liu, Y. Zhou, T. Li, and Y. You. Open-sora: Democratizing efficient video production for all. *arXiv preprint arXiv:2412.20404*, 2024.