

GameFactory: Creating New Games with Generative Interactive Videos

Jiwen Yu^{1*†}

Yiran Qin^{1*}

Xintao Wang^{2‡}

Pengfei Wan²

Di Zhang²

Xihui Liu^{1‡}

¹ The University of Hong Kong

² Kuaishou Technology

<https://vvictoryuki.github.io/gamefactory/>

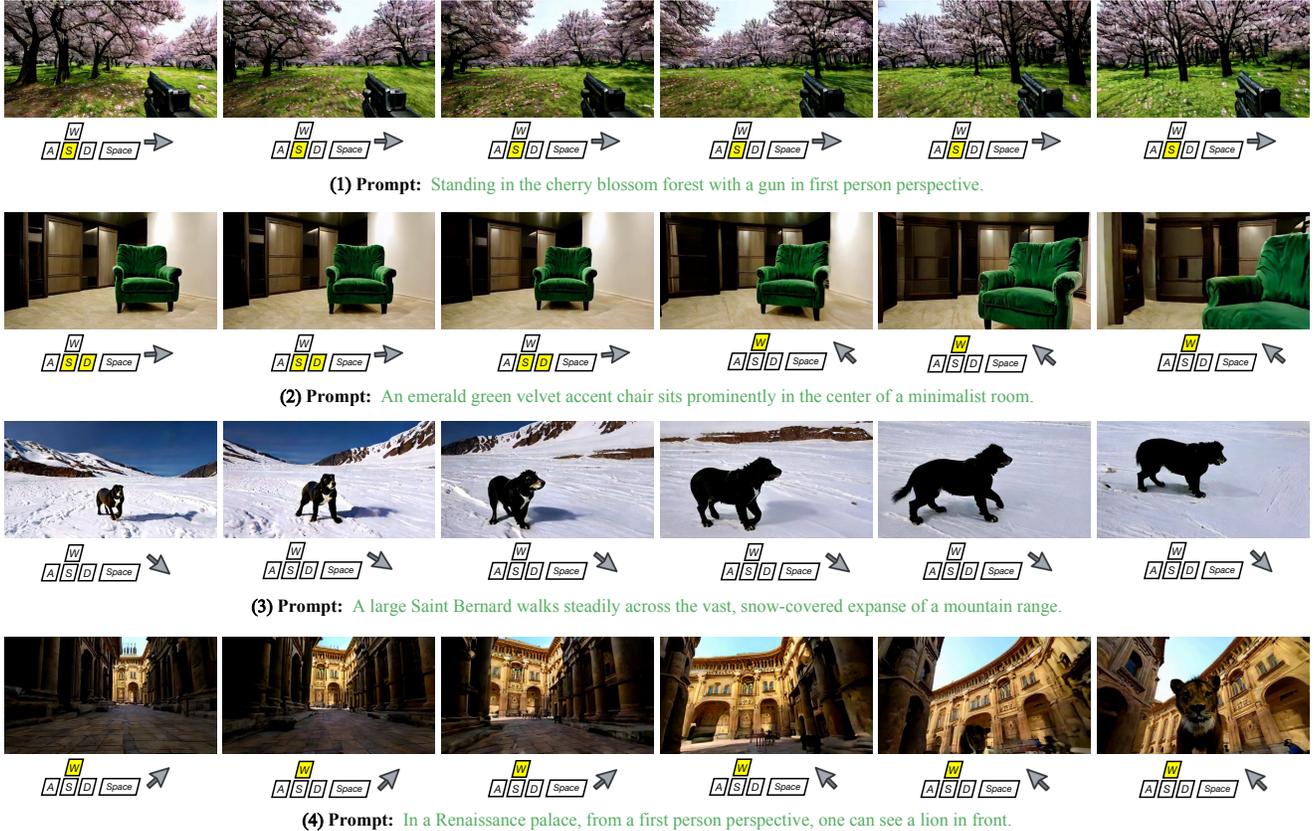


Figure 1. We propose GameFactory, a framework that leverages the powerful generative capabilities of pre-trained video models for the creation of new games. By learning action control from a *small-scale first-person Minecraft dataset*, this framework can transfer these control abilities to open-domain videos, ultimately allowing the creation of new games within open-domain scenes. As illustrated in (1)-(4), GameFactory supports action control across diverse, newly generated scenes in open domains, paving the way for the development of entirely new game experiences. The yellow buttons indicate **pressed keys**, and the arrows represent the direction of **mouse movements**.

Abstract

Generative game engines have the potential to revolutionize game development by autonomously creating new content and reducing manual workload. However, existing video-based game generation methods fail to address the critical challenge of scene generalization, limiting their applicability to existing games with fixed styles and scenes. In this paper, we present GameFactory, a framework focused on exploring scene generalization in game video gen-

eration. To enable the creation of entirely new and diverse games, we leverage pre-trained video diffusion models trained on open-domain video data. To bridge the domain gap between open-domain priors and small-scale game dataset, we propose a multi-phase training strategy that decouples game style learning from action control, preserving open-domain generalization while achieving action controllability. Using Minecraft as our data source, we release GF-Minecraft, a high-quality and diversity action-annotated video dataset for research. Furthermore, we extend our framework to enable autoregressive action-controllable game video generation, allowing the produc-

[‡]Corresponding authors. ^{*}Equal contribution. [†]Work done during an internship at KwaiVGI, Kuaishou Technology.

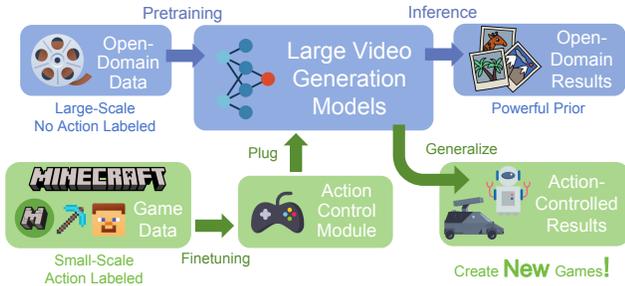


Figure 2. A schematic of our GameFactory creating new games based on pre-trained large video generation models. The upper blue section shows the generative capabilities of the pre-trained model in an open-domain, while the lower green section demonstrates how the action control module, learned from a small amount of game data, can be plugged in to create new games.

tion of unlimited-length interactive game videos. Experimental results demonstrate that GameFactory effectively generates open-domain, diverse, and action-controllable game videos, representing a significant step forward in AI-driven game generation. Our dataset and project page are publicly available at <https://vvictoryuki.github.io/gamefactory/>.

1. Introduction

Video diffusion models have demonstrated impressive capabilities in video generation and real-world physics simulation [24, 26, 45, 48]. This physics simulation ability can be applied to generate game physics, making them promising candidates for game engines [2, 11, 36]. These AI-driven engines have the potential to transform game development by significantly reducing manual work in the traditional game industry through automatic creation of game content. Given this promising direction, it is important to explore how to build a qualified generative game engine.

Generative game engines are typically implemented as video generation models with action controllability, enabling responses to user inputs like keyboard and mouse interactions. Scene generalization, as a crucial yet understudied area, enables the creation of new game scenes beyond existing games. Current research works [2, 3, 11, 36, 42] mainly focus on specific games such as DOOM [36, 42], Atari [2], CS:GO [2], Super Mario Bros [42], and Minecraft [11], or uses limited game-specific datasets [3]. This game-specific approach restricts models to generating content only for existing games, limiting their potential for creating new games. Thus, scene generalization remains a key challenge in advancing generative game engines.

To achieve scene generalization in game videos, collecting large-scale action-annotated video datasets would be the most direct approach. With sufficient data covering all possible scenarios, it could enable arbitrary game scene gen-

eration. However, such annotation is prohibitively expensive and impractical for open-domain scenarios. In contrast, open-domain videos without action annotations are abundant on the Internet and contain rich generative priors. Leveraging scene-generation priors trained on these large-scale open-domain datasets presents a more feasible path to scene generalization. As shown in Figure 2, using a video generation model pre-trained on **large-scale open-domain video data**, we can train an effective action control module with just **a small amount of action-annotated data**, achieving controllability while maintaining the powerful generative prior. We chose Minecraft [13] as our action data source for its unique advantages: diverse and complex action space with customizable inputs, frame-by-frame action annotation support, and cost-effective generation of high-quality action data without human bias.

We present GameFactory, a method for open-domain scene generalization in game video generation. GameFactory has several key designs. First, directly fine-tuning the pretrained model on Minecraft data not only imparts action control capabilities, but also risks embedding a distinct Minecraft style into the generated games, thereby compromising the model’s open-domain generalization. To address this issue, we propose a multiphase decoupled training strategy (Section 3.5) that separates the learning of the game style from the learning of action control. Second, we develop different control mechanisms specifically designed for the distinct nature of continuous mouse movements versus discrete keyboard inputs in the action control module (Section 3.3). Finally, we extend our model to achieve autoregressive long video generation (Section 3.6). Autoregressive generation is particularly desirable, as it enables the production of videos with unlimited length, better satisfying the practical requirements of game video generation.

In summary, our key contributions include:

- We propose GameFactory for open-domain scene generalization in game video generation. To overcome game-specific limitations, we leverage pre-trained video diffusion models with a multiphase training strategy that decouples game style from action control, enabling diverse interactive game environments beyond existing games.
- We introduce GF-Minecraft, a new action-annotated dataset for generative game engines. To address the challenge of human bias in action annotation, we curate a high-quality and diverse dataset from Minecraft.
- We develop an autoregressive generation mechanism that enables unlimited-length action-controllable game videos, meeting the critical requirements of continuous gameplay in practical game applications.

These advancements pave the way for AI-driven game engines with both scene generalization and practical utility.

2. Related Work

2.1. Video Diffusion Models

With the rise of diffusion models [19, 34, 35], significant progress has been made in visual generation, spanning both image [10, 28, 32] and video generation [6, 9, 24, 45, 48]. In particular, recent advancements in video generation have seen a shift from the U-Net architecture to Transformer-based architectures, enabling video diffusion models to generate highly realistic, longer-duration videos [26, 27]. The quality of these generated videos has led to the belief that video diffusion models are capable of understanding and simulating real-world physical rules [26, 40, 44], suggesting their potential use as world models in areas such as autonomous driving and robotics.

2.2. Controllable Video Generation

Text descriptions alone often provide limited control in text-to-video models, leading to ambiguous outputs. To enhance control, several methods have introduced additional control signals. For instance, approaches such as [16, 25, 38] incorporate images as control signals for the video generator, improving both video quality and temporal relationship modeling. Direct-a-Video [43] uses a camera embedder to adjust camera poses, but its reliance on only three camera parameters limits control to basic movements, like left panning. In contrast, MotionCtrl [37] and CameraCtrl [17] offer more complex and nuanced control over camera poses in generated videos.

2.3. Game Video Generation

Given the promising capabilities of video diffusion models, researchers have begun exploring their application in game generation [2, 3, 7, 11, 12, 14, 36, 42]. Genie [7] proposes a foundation model for playable world based on video generation. DIAMOND [2], GameNGen [36], Oasis [11] and PlayGen [42] leverage diffusion-based world modeling for specific games like Atari [2], CS:GO [2], DOOM [36, 42], Minecraft [11] and Super Mario Bro [42]. GameGenX [3] introduces OGameData for game video generation and control. However, these works, including some of our concurrent works [3, 11, 42], often suffer from overfitting to specific games or datasets, exhibiting limited scene generalization capabilities. Other concurrent works, such as Matrix [14] and Genie 2 [12], have discussed the issue of achieving control generalization. However, they still rely on collecting large amounts of high-cost action-labeled data to learn sufficient knowledge about game scenarios.

In contrast, our proposed GameFactory addresses scene generalization by utilizing pretrained video model priors and the easily-accessible and low-cost game dataset GF-Minecraft, with experiments validating its effectiveness.

3. Method

The key to creating new games using video generation models lies in scene generalization. However, it is impractical to fit all possible game scenarios merely by collecting large-scale gaming video datasets, as action annotations are expensive and the datasets can never truly achieve complete open-domain coverage. To address this challenge, we propose leveraging the open-domain generation prior of pre-trained models, combined with a small amount of action-annotated dataset to achieve effective action control with strong scene generalization capabilities.

Sec. 3.1 introduces preliminaries. Sec. 3.2 describes action-annotated data collection. Sec. 3.3 presents the action control module. Sec. 3.4 analyzes the challenges in maintaining the pre-trained model’s open-domain prior while achieving action control. Sec. 3.5 presents a multi-phase training strategy as the solution. Sec. 3.6 describes autoregressive generation for long video synthesis, which is essential for game video applications.

3.1. Preliminaries

We adopt a transformer-based latent video diffusion model [26, 27] as the backbone. Let \mathbf{X} represent a video sequence. To reduce modeling complexity, an Encoder $E(\cdot)$ compresses the video spatially and temporally into a latent representation $\mathbf{Z} = E(\mathbf{X})$. With temporal compression ratio r , a $(1 + rn)$ -frame video is compressed into $(1 + n)$ latent frames. Denoting the i -th frame as \mathbf{x}^i and i -th latent as \mathbf{z}^i , we have $\mathbf{X} = [\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^{rn}]$ and $\mathbf{Z} = [\mathbf{z}^0, \mathbf{z}^1, \dots, \mathbf{z}^n]$. To generate videos, we first train a noise predictor in the latent video diffusion model. During training, noise is added to the clean latent \mathbf{Z}_0 . Let \mathbf{Z}_t denote the noisy latent at timestep t . The noise addition process is formulated as:

$$\mathbf{Z}_t = \sqrt{\bar{\alpha}_t} \mathbf{Z}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1)$$

where $\sqrt{\bar{\alpha}_t}$ and $\sqrt{1 - \bar{\alpha}_t}$ are predefined scheduler parameters, and $\boldsymbol{\epsilon}$ is a random noise sampled from the standard normal distribution. The noise predictor is trained using the following loss function:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}[\|\boldsymbol{\epsilon}_\theta(\mathbf{Z}_t, \mathbf{p}, t) - \boldsymbol{\epsilon}\|_2^2], \quad (2)$$

where $\boldsymbol{\theta}$ denotes the model parameters and \mathbf{p} is the prompt input. During inference, with the trained noise predictor, we can sample clean latent \mathbf{Z}_0 from a noisy latent \mathbf{Z}_T . The predicted latent \mathbf{Z}_0 is then decoded back to video \mathbf{X} through $D(\cdot)$: $\mathbf{X} = D(\mathbf{Z}_0)$. When considering action control, the action $\mathbf{A} = [\mathbf{a}^1, \mathbf{a}^2, \dots, \mathbf{a}^{rn}]$, where \mathbf{a}^i represents the action taken at timestep $(i - 1)$ to transfer from \mathbf{x}^{i-1} to \mathbf{x}^i . The corresponding action-conditioned loss function is:

$$\mathcal{L}_a(\phi) = \mathbb{E}[\|\boldsymbol{\epsilon}_\phi(\mathbf{Z}_t, \mathbf{p}, \mathbf{A}, t) - \boldsymbol{\epsilon}\|_2^2], \quad (3)$$

where ϕ means parameters of the action control module.

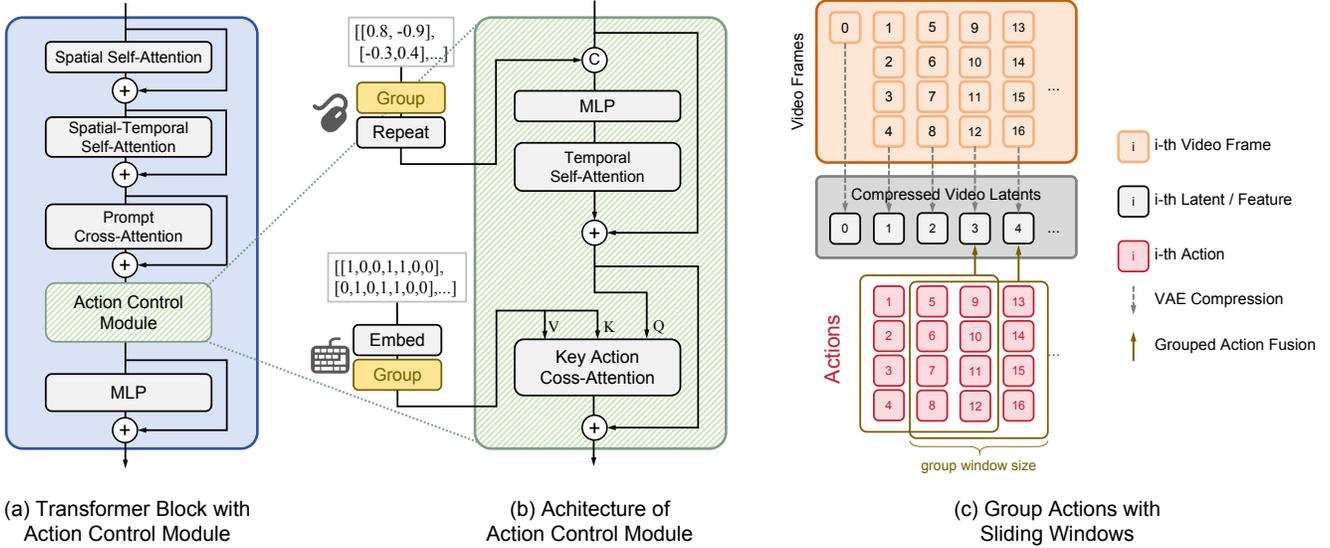


Figure 3. (a) Integration of Action Control Module into transformer blocks of the video diffusion model. (b) Different control mechanisms for continuous mouse signals and discrete keyboard signals (detailed analysis in Sec 3.3). (c) Due to temporal compression (compression ratio $r = 4$), the number of latent features differs from the number of actions, causing granularity mismatch during fusion. Grouping aligns these sequences for fusion. Additionally, the i -th latent feature can fuse with action groups within a previous window (window size $w = 3$), accounting for delayed action effects (e.g., ‘jump’ key affects several subsequent frames).

3.2. GF-Minecraft Dataset

For our action-controllable video generation model to simulate real game engines, the training data should satisfy three key requirements: (1) easily accessible with customizable action inputs to enable cost-effective data collection; (2) action sequences free from human bias, allowing extreme and low-probability action combinations to support arbitrary action inputs; (3) diverse game scenes with corresponding textual descriptions to learn scene-specific physical dynamics. As existing game video datasets fail to meet these criteria, we introduce our GF-Minecraft dataset, where ‘GF’ stands for our method GameFactory and ‘Minecraft’ refers to the game name. The advantages of our dataset are summarized as follows and its details can be found in the Appendix A.

Minecraft as an Accessible Data Source. We leverage Minecraft as our data collection platform for its comprehensive API that captures detailed environmental snapshots, enabling large-scale data collection with action annotations. The game also offers extensive scenes, navigable areas, diverse action space, and an open-world environment. By executing predefined action sequences, we collected 70 hours of gameplay video as GF-Minecraft Dataset.

Collecting Videos with Unbiased Action. Existing Minecraft datasets, such as VPT [4], are collected from real human gameplay, resulting in biased action distributions that favor common human behaviors. Models trained

on these datasets overlook rare action combinations, such as moving backward, jumping in place, or standing still with mouse movements. To eliminate such biases, we decompose keyboard and mouse inputs into atomic actions and ensure their balanced distribution. We also randomize the frame duration of each atomic action to avoid temporal bias.

Diverse Scenes with Textual Descriptions. To enhance dataset diversity, we captured videos across different scenes, weather conditions, and times of day. We segmented the videos and annotated them with textual descriptions using the pre-trained multimodal large language model MiniCPM [46].

3.3. Action Control Module

We incorporate action control modules into the transformer blocks of the video diffusion model to enable action-controllable generation. The architecture of the transformer block is shown in Fig. 3 (a), and the structure of the action control module is shown in Fig. 3 (b). Suppose that the input action includes both continuous mouse movement action \mathbf{M} and discrete keyboard action \mathbf{K} . The intermediate feature in the Transformer is denoted as \mathbf{F} :

$$\mathbf{M} = [\mathbf{m}^1, \mathbf{m}^2, \dots, \mathbf{m}^{rn}] \in \mathbb{R}^{rn \times d_1}, \quad (4)$$

$$\mathbf{K} = [\mathbf{k}^1, \mathbf{k}^2, \dots, \mathbf{k}^{rn}] \in \mathbb{R}^{rn \times d_2}, \quad (5)$$

$$\mathbf{F} = [\mathbf{f}^0, \mathbf{f}^1, \dots, \mathbf{f}^n] \in \mathbb{R}^{(n+1) \times l \times c}, \quad (6)$$

where \mathbf{m}^i represents the i -th mouse action, \mathbf{k}^i represents the i -th keyboard action, \mathbf{f}^i represents the feature of the i -th

latent frame. The total number of video frames is $(1 + rn)$, compressed to $(1 + n)$ latent frames, where r denotes the temporal compression ratio. d_1 and d_2 denote the dimension numbers of actions, l is the length of token sequence and c is the number of feature channels.

Grouping Actions with a Sliding Window. Due to the temporal compression ratio r , the number of actions (rn) differs from the number of features $(n + 1)$, creating a granularity mismatch for action-feature fusion. As shown in Fig. 3 (c), we address this by grouping actions using a sliding window of size w . For the i -th feature \mathbf{f}^i , we consider actions within $[\mathbf{a}^{r \times (i-w+1)}, \dots, \mathbf{a}^{ri}]$. This window design captures delayed action effects, such as how a jump command influences multiple subsequent frames. For out-of-range indices, boundary actions are used as padding. For mouse movement actions \mathbf{M} , the grouping operation can be represented as:

$$\mathbf{M}_{group} = \mathbf{group}(\mathbf{M}) \in \mathbb{R}^{(n+1) \times rw \times d_1}. \quad (7)$$

As for keyboard actions \mathbf{K} , we first learn the embedding of actions and add the positional encoding to get:

$$\mathbf{K}_{embed} = \mathbf{embed}(\mathbf{K}) \in \mathbb{R}^{rn \times c}. \quad (8)$$

After that, perform a grouping operation on the action embeddings, which can be represented as follows:

$$\mathbf{K}_{group} = \mathbf{group}(\mathbf{K}_{embed}) \in \mathbb{R}^{(n+1) \times rw \times c}. \quad (9)$$

Mouse Movements Control. To fuse the grouped mouse action \mathbf{M}_{group} with feature \mathbf{F} , we first reshape it from $\mathbb{R}^{(n+1) \times rw \times d_1}$ to $\mathbb{R}^{(n+1) \times 1 \times rwd_1}$. Then we repeat it in the dimension of the token sequence length to get:

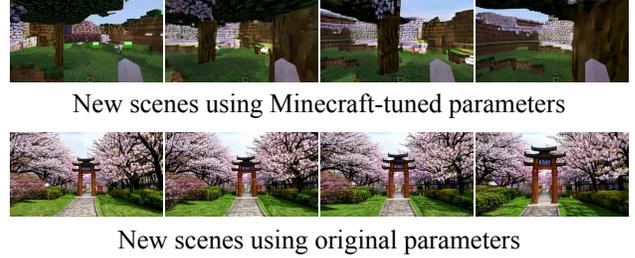
$$\mathbf{M}_{repeat} = \mathbf{repeat}(\mathbf{M}_{group}) \in \mathbb{R}^{(n+1) \times l \times rwd_1}, \quad (10)$$

and concatenate it with \mathbf{F} along with the channel dimension:

$$\mathbf{F}_{fused} = \mathbf{cat}([\mathbf{F}, \mathbf{M}_{repeat}]) \in \mathbb{R}^{(n+1) \times l \times (c+rwd_1)}. \quad (11)$$

After that, further learning on \mathbf{F}_{fused} is conducted through a layer of MLP and a layer of temporal self-attention.

Keyboard Actions Control. For discrete keyboard control, we perform a cross-attention calculation between the grouped action embeddings $\mathbf{K}_{group} \in \mathbb{R}^{(n+1) \times rw \times c}$ and $\mathbf{F} \in \mathbb{R}^{(n+1) \times l \times c}$, similar to the prompt cross-attention between text and \mathbf{F} . Specifically, \mathbf{K}_{group} serves as the key and value in the attention, while \mathbf{F} functions as the query.



Prompt: In a village, from a first person perspective, walking towards a torii gate with cherry blossoms in full bloom.

Figure 4. Style bias in video game generation. The model tuned on Minecraft data inherits its distinctive pixelated block style, creating a domain gap from the original parameters. This motivates decoupling action control learning from data style learning through specialized training strategies.

Action-Specific Control Module. Our experiments indicate concatenation works better for continuous actions (e.g., mouse movements) while cross-attention suits discrete actions (e.g., keyboard inputs). For continuous actions, concatenation preserves the magnitude information critical for control, whereas cross-attention’s similarity computations and normalization would diminish these numerical differences. For discrete actions, cross-attention’s effectiveness aligns with its successful applications in text prompt control, demonstrating its strength in categorical processing.

3.4. Leveraging Pre-trained Models for Action-Controlled Scene Generalization

Recent video generation models for games [2, 11, 36, 42], are limited to reproducing existing games rather than creating novel games, hindering the vision of a generative game engine that could reduce game development costs. While a large-scale open-domain game video dataset would be ideal, annotating such data is prohibitively expensive and achieving true open-domain coverage remains infeasible. We explore leveraging a pre-trained video generation model that can learn rich generative priors from unlabeled open-domain data. By fine-tuning an action control module using limited game data (e.g., our proposed GF-Minecraft Dataset) on this pre-trained model, we enable controlled game creation across open-domain scenes.

However, this approach faces a key challenge: achieving scene generalization without stylistic bias. While our model enables new game scene generation, the outputs tend to inherit the visual style of the training data (e.g., Minecraft’s pixelated blocks, Fig. 4). This style bias stems from the diffusion model’s loss function (Eq. 3), where learning action control naturally captures dataset-specific visual features that help minimize training loss. To address this, we next propose a multi-phase training strategy that decouples action control learning from visual style, enabling robust

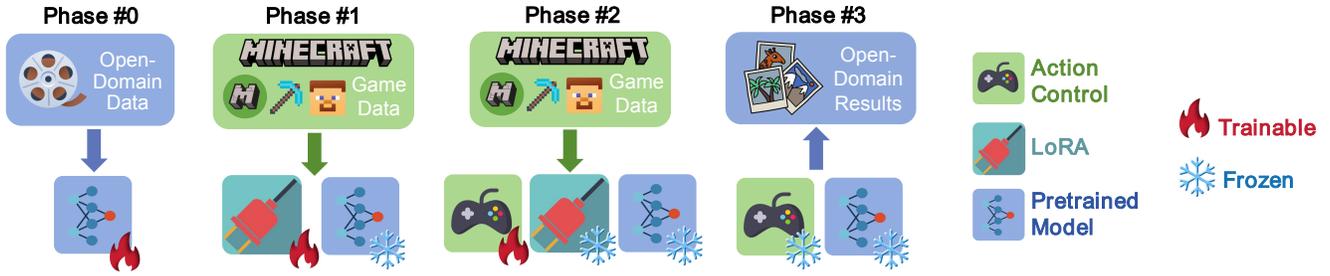


Figure 5. **Phase #0**: pretraining a video generation model on open-domain data. **Phase #1**: finetuning with LoRA for game video data. **Phase #2**: training the action control module while fixing other parameters. **Phase #3**: inference for action-controlled open-domain generation. To decouple style learning from action control, **Phase #1** learns game-specific style while **Phase #2** focuses on style-independent action control. This design preserves the open-domain capabilities from **Phase #0**, enabling generalization in **Phase #3**.

action control independent of specific data styles.

3.5. Multi-Phase Training Strategy

As illustrated in Fig. 5, our training process consists of **Phase #0** (model pretraining) and the following phases:

Phase #1: Tune LoRA to Fit Game Videos. We finetune the pre-trained video diffusion model using LoRA [21] to adapt it to specific game video while preserving most original parameters. This produces a model specialized for the target game domain. Better style adaptation in this phase allows the next phase to focus purely on action control, reducing style-control entanglement.

Phase #2: Tune Action Control Module. We freeze both pre-trained parameters and LoRA, only training the action control module with game videos and action signals. Since **Phase #1** has handled style adaptation through LoRA, the training loss now focuses less on style learning. This allows the model to concentrate on action control learning, as it becomes the main contributor to minimizing the diffusion loss. Such separation enables style-independent control that can generalize to open-domain scenarios.

Phase #3: Inference on Open Domain. During inference, we remove the LoRA weights for game style adaptation, keeping only the action control module parameters. Thanks to the decoupling in previous phases, the action control module can now work independently of specific game styles, enabling controlled game video generation across open-domain scenarios.

3.6. Autoregressive Generation of Long Action-Controllable Game Videos

Current action-controllable video generation methods are limited to fixed-length outputs, which is inadequate for practical game applications requiring continuous video

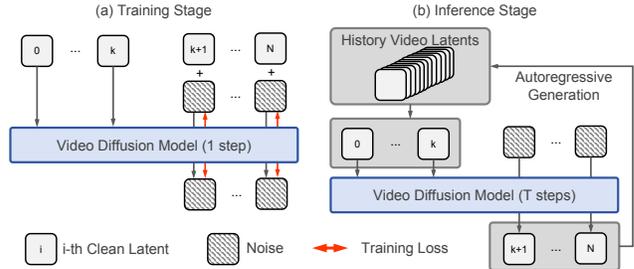


Figure 6. Illustration of autoregressive video generation. The frames from index 0 to k serve as conditional frames, while the remaining $N - k$ frames are for prediction, with k randomly selected. (a) Training stage: Loss computation and optimization focus only on the noise of predicted frames. (b) Inference stage: The model iteratively selects the latest $k + 1$ frames as conditions to generate $N - k$ new frames, enabling autoregressive generation.

streams. To address this, we develop an autoregressive approach that generates multiple frames per step based on previous outputs, enabling efficient long video generation.

One established approach for autoregressive video generation follows next-token prediction [23, 39, 47], but lacks accessible strong model parameters. While video diffusion transformer [24, 26, 27, 48] offers superior generation quality, it typically operates in a full-sequence manner. Inspired by Diffusion Forcing [8], which unlike standard diffusion models that require identical noise levels across frames, allows different noise levels where later frames with more noise can depend on earlier frames with less noise, we modify the video diffusion transformer to enable autoregressive generation. As shown in Fig. 6 (a), during training with $N + 1$ frame latents (from 0 to N), we randomly select the first $k + 1$ frames as conditions without adding noise, while adding noise only to the remaining $N - k$ frames for noise prediction training. Although the first $k + 1$ frames are input to the model, since they are assumed to be already generated, their predicted noise outputs are not utilized. To improve training efficiency, we focus on computing training losses only for the $N - k$ frames that require predic-

Control Module		Only-Key			Mouse-Small			Mouse-Large		
Key	Mouse	Flow-MSE↓	CLIP-Sim↑	FID↓	Flow-MSE↓	CLIP-Sim↑	FID↓	Flow-MSE↓	CLIP-Sim↑	FID↓
Cross Attention	Cross Attention	8.67	0.3313	107.13	20.46	0.3137	125.67	325.18	0.3103	167.37
Concat	Concat	22.37	0.3277	103.89	19.18	0.3159	133.42	258.93	0.3123	145.74
Cross Attention	Concat	7.79	0.3292	105.28	18.64	0.3184	127.84	249.54	0.3107	139.91

Table 1. Results of the ablation study on action control mechanisms. The findings indicate that an optimal approach for the action control module is to use cross-attention for discrete action control and concatenation for continuous action control.

tion. As for inference, as shown in Fig. 6 (b), after full-sequence generation of the first $N + 1$ frame latents, we can autoregressively generate new $N - k$ frames by selecting the most recent $k + 1$ frames from history video latents as conditions for each subsequent generation, and merge them into the history video latents. This process can be repeated to achieve infinite-length video generation. Unlike conventional next-frame generation methods [8, 11, 36], our approach supports **multi-frame generation in one step**, greatly reducing the time for long video generation.

Implementation Details. Since the model structure remains unchanged and only the noise levels vary across frames, the training configuration remains identical to full-sequence generation, including action control. We also add small noise perturbations (equivalent to step 15/1000) to conditional frames, which helps reduce error accumulation in long-term generation [11, 36]. The long video results for action control can be found on our project page.

4. Experiments

4.1. Implementation Details

Pretrained Model Setting. Our experiments are based on an internal 1B-sized transformer-based text-to-video diffusion model for research purpose, which is distilled from a larger pretrained video diffusion model and possesses strong generative priors in the open domain. The resolution of the game videos is 360×640 . The temporal compression rate of the VAE is $r = 4$.

Training and Inference Setting. Each phase of fine-tuning or training requires about two-four days of training on 8 A100 GPUs with a batch size of 64. The hyper parameters for LoRA finetuning can be referenced as $rank = 128$, with a learning rate of $1e-4$. The learning rate for training the action control module is $1e-5$. We only apply classifier free guidance [18] to the text prompt conditional input and use DDIM [33] sampling with 50 sampling steps.

Evaluation Setting. We retained 5% of our collected, segmented dataset as a test set, excluding it from training, and selected three subsets from it for ablation study: (1) **only-key**: contains only keyboard actions, designed to test

the model’s ability to follow discrete actions; (2) **mouse-small**: includes small-scale continuous mouse movements; and (3) **mouse-large**: includes large-scale continuous mouse movements. Additionally, for qualitative experiments, we support custom combinations of input actions, allowing us to test complex or rare action combinations. We use the following evaluation metrics: (1) **Flow-MSE**: calculates the optical flow of the generated video to reflect its dynamics, assessing action-following ability by measuring mean square error to the optical flow of the reference video; (2) **CLIP-Sim**: computes feature similarity in the CLIP [31] space to evaluate semantic relevance to the given text prompt; (3) **FID**: measures distribution differences between the generated videos and the reference, providing an assessment of generation quality.

4.2. Action Controllability

Ablation Study. We perform ablation studies on the control mechanisms for continuous mouse movement signals and discrete keyboard control signals, comparing two typical methods: cross-attention and concatenation-based injection. The results are presented in Tab. 1. To evaluate how well the model’s generated dynamics align with action controls, we use Flow-MSE as a metric. For discrete control signals, such as keyboard inputs, cross-attention outperforms concatenation, suggesting that category-based signal control benefits from similarity-based cross-attention, as is often used in text-based control. In contrast, for continuous mouse movement signals, concatenation is more effective than cross-attention. This may be due to cross-attention’s similarity computation, which tends to reduce the influence of the control signal’s magnitude, thereby affecting the final result. Additionally, the Flow-MSE values reveal that mouse movements have a larger impact on the visual output than keyboard inputs, particularly in the mouse-large test set where movement magnitude is greater. For the CLIP-Sim and FID metrics, there is little difference between the methods. This is primarily because our multi-phase training strategy decouples style learning into **Phase #1**, where style learning is consistent across all methods. As a result, metrics assessing semantic consistency and generation quality yield similar performance across different methods.

Demonstration of various action control in Minecraft.

In Fig. 7, we demonstrate the action control capabilities of

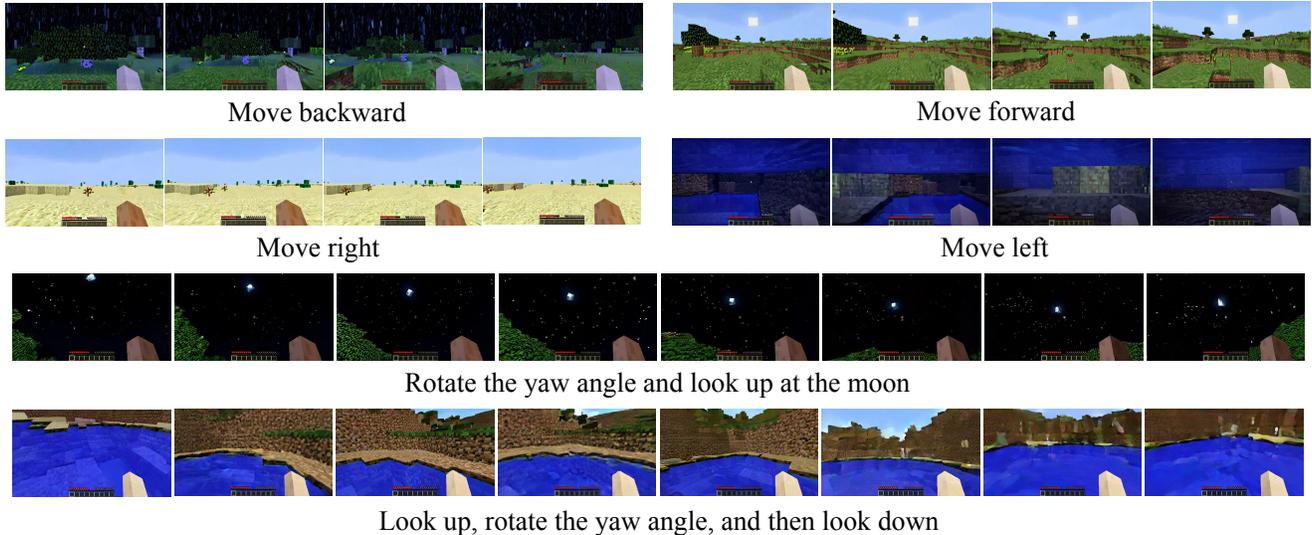


Figure 7. Demonstration of action control capabilities in the Minecraft domain. The model has successfully learned basic atomic actions (WASD keys) and mouse-based yaw and pitch controls. Additionally, it can combine these atomic actions to execute more complex movements. Note that the text below each video frame is a descriptive label of the content, not a text prompt provided to the model.

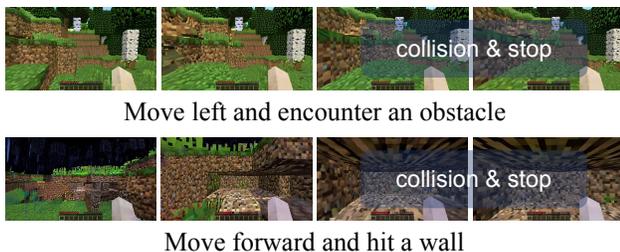


Figure 8. Demonstration of the learned response to collision, one of the most common interactions in Minecraft navigation. Note that the text below each video frame is a descriptive label of the content, not a text prompt provided to the model.

our proposed action control module within the Minecraft domain. The module successfully learns fundamental atomic actions, such as moving forward, backward, left, and right, as well as adjusting yaw and pitch angles. Additionally, the model learns to combine these actions to achieve more complex control. The diverse scenes in Fig. 7 reflect the scene variety in our dataset, highlighting its strength in capturing a wide range of scenarios. Notably, actions like moving backward, left, and right, although seemingly simple, are rare in typical human Minecraft gameplay. A dataset based on human gameplay would likely suffer from biased distributions, making it challenging to learn these less common actions effectively.

Interaction in generative Minecraft videos. For navigation agents, collision is one of the most critical physical interactions in a simulator. Since our data collection process

involves randomly generated scenes, our dataset naturally includes numerous examples of collisions. In these cases, even when action inputs are provided, the agent’s behavior should ideally remain stationary. Such corner cases can impact the learning of the action control module, especially when data volume is limited. However, during inference, we observed that our model has developed collision detection ability and provides appropriate interaction feedback, as shown in Fig. 8.

4.3. Scene Generalization

Create new games in generalized scenes. In Fig. 1, we showcase a variety of newly generated game videos across open-domain scenes. In theory, any scene within the generative capabilities of the pre-trained model could be included as content in these new games. These results inspire a vision of the future where generative game engines emerge as a new form of gaming, allowing players or game creators to generate and interact with anything they can imagine at minimal cost.

More inspiration from examples of racing games. In Fig. 9, we discuss an intriguing example. Since our collected Minecraft data is from a first-person perspective, the learned action space primarily generalizes to first-person scenes. Here, we experimented with a racing game scenario using a prompt for a car. Interestingly, we observed that the model’s learned yaw control for the mouse seamlessly generalized to steering control in the racing game. Additionally, certain directional controls, such as moving backward or sideways, were diminished, an adaptation that aligns well



Prompt: On a racing track, from a first person perspective, one can see holding a steering wheel.

Figure 9. Our model demonstrates the ability to generalize to a different game type, a racing game. Interestingly, the yaw control learned in Minecraft seamlessly transfers to steering control in the racing game, while unrelated actions, such as moving backward, left, or right, and pitch angle adjustments, automatically diminish.

with typical controls in a racing game, where these actions are rarely needed. This example not only highlights the strong generalization capabilities of our method but also raises the question: could there exist a larger, more versatile action space that encompasses a wider range of game controls, extending beyond first-person and racing games?

This example also leads us to wonder whether the racing game scenario could have applications in autonomous driving. If we were to collect an action dataset from an autonomous driving simulation environment, could a pre-trained model then generate unlimited open-domain autonomous driving data? Our exploration of scene generalization within generative game engines may hold valuable insights for other fields as well.

5. Potential of Generalizable World Model

We propose that the **GameFactory** we have developed is not merely a tool for creating new games but a **Generalizable World Model** with far-reaching implications. This model has the capability to generalize physical knowledge learned from small-scale labeled datasets to open-domain scenarios, addressing challenges in areas like autonomous driving [15, 20] and embodied AI [5, 29, 41, 49], which also face limitations due to the lack of large-scale action-labeled datasets. The generalizable world model has two key applications from different perspectives:

- **As a data producer:** It transfers knowledge from small labeled datasets to open-domain scenarios, enabling the generation of diverse unlimited action-annotated data that closely approximates real-world complexity.
- **As a simulator:** It provides an environment for directly training agents to perform real-world tasks [1, 22, 30], closely approximating real-world conditions. By enabling controlled and diverse scenario generation, including extreme situations that are difficult to capture in real-world data collection, it facilitates the development of policy models that are exposed to a wide range of environments and interactions, thereby improving their robustness and generalizability, and aiding in overcoming the challenges of sim-to-real transfer.

6. Conclusion

In this paper, we propose GameFactory, a framework that uses generative interactive videos to create new games, addressing significant gaps in existing research, particularly in scene generalization. The study of generative game engines still faces many challenges, including the design of diverse levels and gameplay, player feedback systems, in-game object manipulation, long-context memory, and real-time game generation. GameFactory marks our first effort in this field, and we aim to continue progressing toward the realization of a fully capable generative game engine.

References

- [1] Anurag Ajay, Seungwook Han, Yilun Du, Shuang Li, Abhi Gupta, Tommi Jaakkola, Josh Tenenbaum, Leslie Kaelbling, Akash Srivastava, and Pulkit Agrawal. Compositional foundation models for hierarchical planning. *Advances in Neural Information Processing Systems*, 36, 2024. 9
- [2] Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kanervisto, Amos Storkey, Tim Pearce, and François Fleuret. Diffusion for world modeling: Visual details matter in atari. In *Thirty-eighth Conference on Neural Information Processing Systems*, 2024. 2, 3, 5
- [3] Anonymous. Gamegen- \mathbb{X} : Interactive open-world game video generation. In *Submitted to The Thirteenth International Conference on Learning Representations*, 2024. under review. 2, 3
- [4] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampeiro, and Jeff Clune. Video pretraining (VPT): Learning to act by watching unlabeled online videos. In *Advances in Neural Information Processing Systems*, 2022. 4
- [5] Amir Bar, Gaoyue Zhou, Danny Tran, Trevor Darrell, and Yann LeCun. Navigation world models. *arXiv preprint arXiv:2412.03572*, 2024. 9
- [6] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 3
- [7] Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al.

- Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024. 3
- [8] Boyuan Chen, Diego Marti Monso, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann. Diffusion forcing: Next-token prediction meets full-sequence diffusion, 2024. 6, 7
- [9] Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models, 2024. 3
- [10] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- α : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 3
- [11] Etched Decart. Oasis: A universe in a transformer. <https://oasis-model.github.io/>, 2024. 2, 3, 5, 7
- [12] Google DeepMind. Genie 2: A large-scale foundation world model. <https://deepmind.google/discover/blog/genie-2-a-large-scale-foundation-world-model/>, 2024. 3
- [13] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. 2, 12
- [14] 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. 3
- [15] Shenyan Gao, Jiazhi Yang, Li Chen, Kashyap Chitta, Yihang Qiu, Andreas Geiger, Jun Zhang, and Hongyang Li. Vista: A generalizable driving world model with high fidelity and versatile controllability. *arXiv preprint arXiv:2405.17398*, 2024. 9
- [16] Yuwei Guo, Ceyuan Yang, Anyi Rao, Maneesh Agrawala, Dahua Lin, and Bo Dai. Sparsectrl: Adding sparse controls to text-to-video diffusion models. In *European Conference on Computer Vision*, pages 330–348. Springer, 2025. 3
- [17] 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
- [18] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 7
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 2020. 3
- [20] Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shotton, and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving. *arXiv preprint arXiv:2309.17080*, 2023. 9
- [21] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. 6
- [22] Jianguo Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d world. *arXiv preprint arXiv:2311.12871*, 2023. 9
- [23] Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Grant Schindler, Rachel Hornung, Vignesh Birodkar, Jimmy Yan, Ming-Chang Chiu, et al. Videopoet: A large language model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023. 6
- [24] PKU-Yuan Lab and Tuzhan AI etc. Open-sora-plan. <https://github.com/PKU-YuanGroup/Open-Sora-Plan>, 2024. 2, 3, 6
- [25] Haomiao Ni, Changhao Shi, Kai Li, Sharon X Huang, and Martin Renqiang Min. Conditional image-to-video generation with latent flow diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18444–18455, 2023. 3
- [26] OpenAI. Creating video from text. <https://openai.com/index/sora/>, 2024. 2, 3, 6
- [27] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023. 3, 6
- [28] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 3
- [29] Yiran Qin, Zhelun Shi, Jiwen Yu, Xijun Wang, Enshen Zhou, Lijun Li, Zhenfei Yin, Xihui Liu, Lu Sheng, Jing Shao, et al. Worldsimbench: Towards video generation models as world simulators. *arXiv preprint arXiv:2410.18072*, 2024. 9
- [30] Yiran Qin, Enshen Zhou, Qichang Liu, Zhenfei Yin, Lu Sheng, Ruimao Zhang, Yu Qiao, and Jing Shao. Mp5: A multi-modal open-ended embodied system in minecraft via active perception. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16307–16316, 2024. 9
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 2021. 7
- [32] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022. 3
- [33] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv:2010.02502*, 2020. 7
- [34] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 2019. 3

- [35] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. International Conference on Learning Representations, 2021. 3
- [36] Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-time game engines. arXiv preprint arXiv:2408.14837, 2024. 2, 3, 5, 7
- [37] Zhouxia Wang, Ziyang Yuan, Xintao Wang, Yaowei Li, Tianshui Chen, Menghan Xia, Ping Luo, and Ying Shan. Motionctrl: A unified and flexible motion controller for video generation. In ACM SIGGRAPH 2024 Conference Papers, 2024. 3, 13
- [38] Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Xintao Wang, Tien-Tsin Wong, and Ying Shan. Dynamicrafter: Animating open-domain images with video diffusion priors, 2023. 3
- [39] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using vq-vae and transformers. arXiv preprint arXiv:2104.10157, 2021. 6
- [40] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. arXiv preprint arXiv:2310.06114, 2023. 3
- [41] Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. arXiv preprint arXiv:2310.06114, 2023. 9
- [42] Mingyu Yang, Junyou Li, Zhongbin Fang, Sheng Chen, Yangbin Yu, Qiang Fu, Wei Yang, and Deheng Ye. Playable game generation. arXiv preprint arXiv:2412.00887, 2024. 2, 3, 5
- [43] Shiyuan Yang, Liang Hou, Haibin Huang, Chongyang Ma, Pengfei Wan, Di Zhang, Xiaodong Chen, and Jing Liao. Direct-a-video: Customized video generation with user-directed camera movement and object motion. In ACM SIGGRAPH 2024 Conference Papers, pages 1–12, 2024. 3
- [44] Sherry Yang, Jacob C Walker, Jack Parker-Holder, Yilun Du, Jake Bruce, Andre Barreto, Pieter Abbeel, and Dale Schuurmans. Position: Video as the new language for real-world decision making. In Proceedings of the 41st International Conference on Machine Learning, 2024. 3
- [45] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. arXiv preprint arXiv:2408.06072, 2024. 2, 3
- [46] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. arXiv preprint arXiv:2408.01800, 2024. 4, 12
- [47] Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. Language model beats diffusion—tokenizer is key to visual generation. arXiv preprint arXiv:2310.05737, 2023. 6
- [48] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all. <https://github.com/hpcaitech/OpenSora>, 2024. 2, 3, 6
- [49] Fangqi Zhu, Hongtao Wu, Song Guo, Yuxiao Liu, Chilam Cheang, and Tao Kong. Irasim: Learning interactive real-robot action simulators. arXiv preprint arXiv:2406.14540, 2024. 9

Additional results of action control in Minecraft and open-domain scenarios can be found in our project page: <https://vvictoryuki.github.io/gamefactory/>.

Appendix A. Details of GF-Minecraft Dataset

A.1. Basic Information

For data collection, we use Minedojo [13] to obtain Minecraft snapshots which contain a wide array of details within the agent’s surroundings. To mitigate biases introduced by human player habits, we uniformly sampled the frequency and duration of each action and randomly combined them to enhance action generalization (e.g., pressing two keys simultaneously or combining keyboard and mouse operations). Subsequently, the agent was connected to the MineDojo simulation environment to execute the action sequences and collect corresponding observations.

To further enhance the diversity and generalization of the generated game scenes, we preconfigured three biomes (forest, plains, desert), three weather conditions (clear, rain, thunder), and six different times of day (“Starting of a day”, “Noon, sun is at its peak”, “Sunset”, “Beginning of night”, “Midnight, moon is at its peak”, “Beginning of sunrise”). This process resulted in the generation of 2,000 video clips with action annotations, each containing 2,000 frames.

A.2. Data PreProcessing

After obtaining the full video, we preprocess it in two steps: (1) **Slicing**: Randomly sample consecutive sequences of k frames from the full video to form new video clips. Sampling continues until the total number of sampled frames reaches n times the total frame count of the original video. In our experimental setup, $k = 81$, which is slightly larger than the generation length of the video diffusion model (77), and $n = 3$; (2) **Text Annotation**: We use the open-source multimodal large language model MiniCPM-V [46] to annotate the sliced video clips with text. The corresponding prompts and annotation examples are shown in Fig. 10.



Q: Describe the scenery and objects in the Minecraft video.

A: The video showcases a **desert biome** in the game Minecraft, characterized by **sandy terrain** and **sparse vegetation**. The environment includes **sand blocks** forming the ground and **cacti** scattered around as the only **plant life** visible. In the background, there are rock formations that add to the arid **landscape's ruggedness**. A notable feature is a **green, pixelated entity** resembling a **zombie-like figure** with outstretched arms, which remains stationary throughout the video. The sky transitions from **orange hues** at the horizon to **darker shades** above, suggesting either **sunrise or sunset**. The player's HUD (Heads-Up Display) is consistently present, indicating health and inventory status, but no actions such as mining or crafting are depicted within these frames.

Figure 10. An example of video clip annotation, where words describing scenes and objects are highlighted in red and bolded.

A.3. Details of Action Space

We use the part of the action space of Minedojo [13] which encompasses nearly all actions available to human players, including keypresses, mouse movements, and clicks. We used keypresses and mouse movements as our control signal. The specific binary actions used in our setup are listed in Table 2. Interface₁ to Interface₅ represent different MineDojo interfaces, where mutually exclusive actions (e.g., moving forward and backward) are assigned to the same interface. Mouse movements are represented by the offset of the mouse pointer from the center of the game region. For each frame in the video, we calculate the cumulative offset relative to the first frame, and this absolute offset is used as input to the model.

Appendix B. Supplementary Comparative Experimental Results

B.1. Qualitative Results of Ablation Study on Action Control Module

Based on the experimental results shown in Table 1, it can be observed that cross attention performs better in terms of action control capability for keyboard inputs, while concatenation performs relatively poorly. In our experiments, we even found

Table 2. **Details of Action Space.** The term *Control Signal* refers to the raw input signals utilized for training purposes, while the *Action Interface* represents the corresponding interface in the MineDojo platform that maps these input signals to actionable commands.

Behavior	Control Signal	Action Interface
forward	W key	Interface ₁
back	S key	Interface ₁
left	A key	Interface ₂
right	D key	Interface ₂
jump	space key	Interface ₃
sneak	shift key	Interface ₃
sprint	ctrl key	Interface ₃
vertical perspective movement	mouse movement(yaw)	Interface ₄
horizontal perspective movement	mouse movement(pitch)	Interface ₅

that concatenation sometimes struggles to follow discrete keyboard input signals, which is consistent with the significant differences observed in the metrics. Related qualitative results are shown in Fig. 11.

Additionally, the metric results in Table 1 indicate that concatenation has certain advantages over cross attention in modeling continuous mouse movements. This is somewhat similar to some works on camera pose control [37]. However, upon a closer comparison of the qualitative results, we found it challenging to distinguish the differences between them visually. This is reasonable because, as shown in the quantitative results in Table 1, the difference between concatenation and cross attention is relatively small, even for large-scale variations in mouse movements.

B.2. Evaluation on Multi-Phase Training Strategy for Scene Generalization

To evaluate the scene generalization capability of GameFactory in open-domain scenarios, we used the action sequences from the test set along with open-domain prompts and generated open-domain results using a model trained with multi-phase training. Similarly, we used the action sequences and prompts annotated from the test set to generate Minecraft videos with the same action controls using the model trained with multi-phase training. We then compared the Flow-MSE of both results against the reference videos. As shown in Tab. 3, the differences between the two are minor, indicating that the action control capability learned on Minecraft data has been transferred to open-domain videos to a certain extent.

Method	Domain	Flow-MSE↓	Domain-Sim↑	CLIP-Sim↑	FID↓
Multi-Phase Training	Minecraft	43.48	-	-	-
Multi-Phase Training	Open-Domain	54.13	0.7565	0.3181	121.18
One-Phase Training	Open-Domain	76.02	0.7345	0.3111	167.79

Table 3. Quantitative results of evaluation on scene generalization.

Next, we compared the generalization capability in open-domain scenarios between multi-phase training and one-phase training (i.e., directly training action control on the original parameters) after replacing the text-to-video parameters with the original parameters. The Domain-Sim metric is calculated by extracting CLIP features from the video frames and comparing their similarity in CLIP space to the CLIP features of corresponding video frames generated by the original model with the same prompt. This metric reflects the domain difference between the generated open-domain action-controlled videos and those generated by the original model. According to the results in Tab. 3, multi-phase training demonstrates better action control capability and achieves higher domain similarity with the original model. Additionally, the generated videos exhibit slightly better quality and alignment with the prompts. The qualitative results are shown in Fig. 12, where it can be observed that one-phase training results in degraded action control capability in open-domain scenarios. Moreover, the generated frames display distortion, leading to an overall decline in video quality.

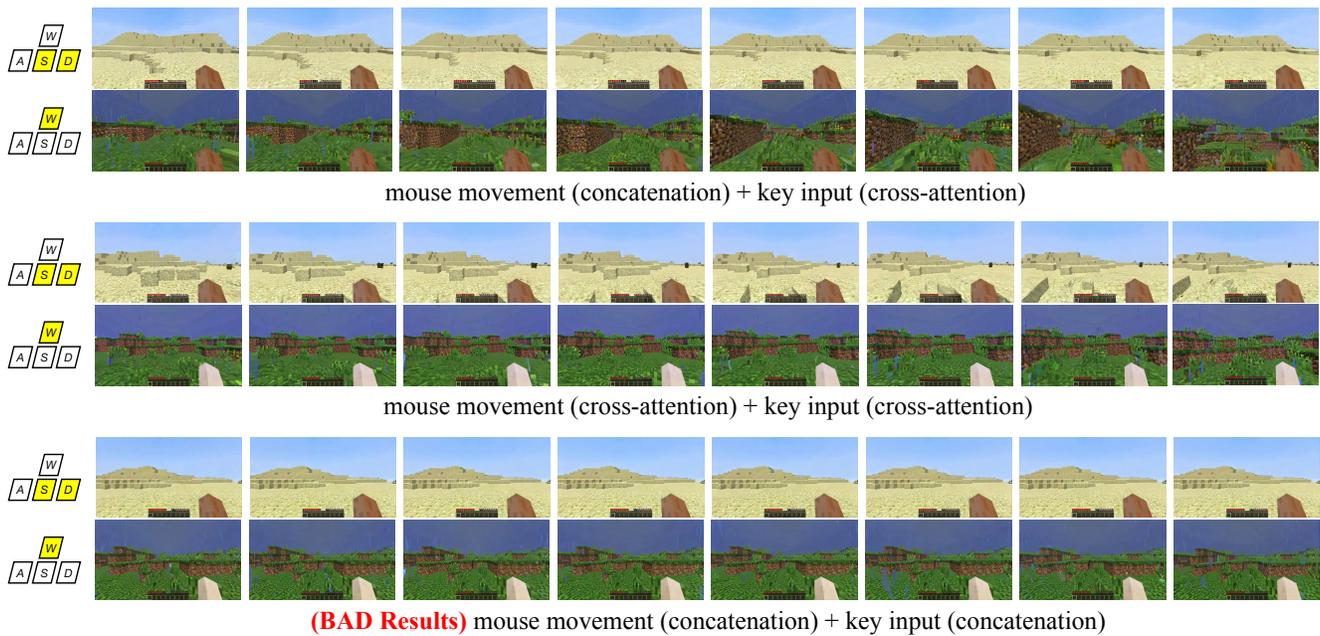
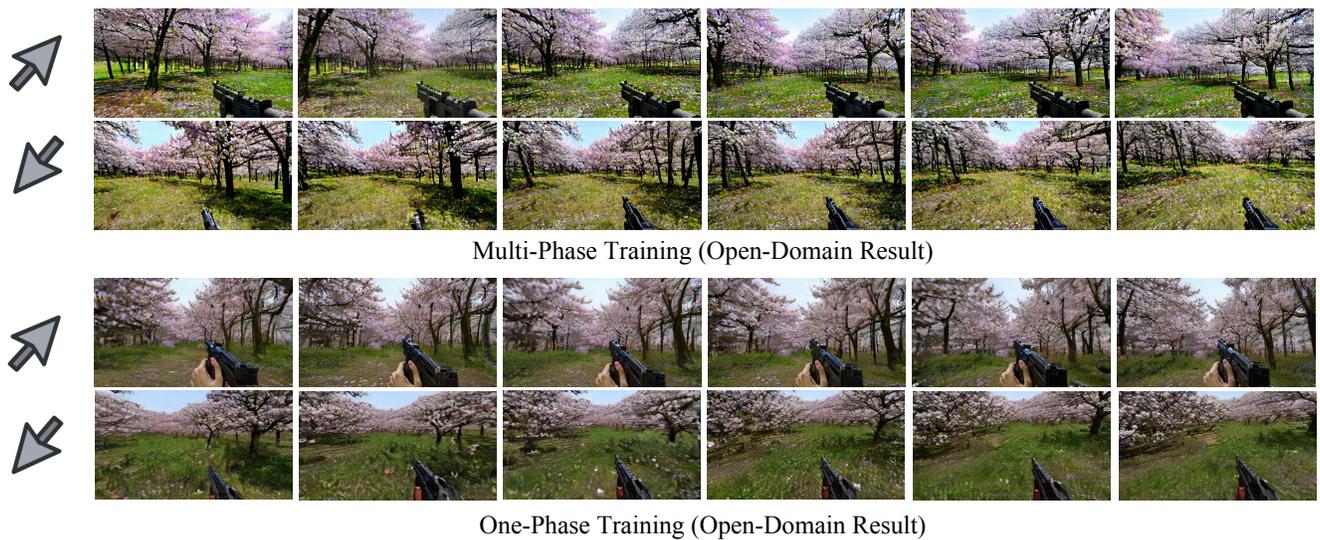


Figure 11. Qualitative comparison of key input control performance. It can be observed that cross-attention significantly outperforms concatenation in handling discrete key input signals, while concatenation may fail to respond to the key input. The yellow buttons indicate pressed keys.



Prompt: Standing in the cherry blossom forest with a gun in first person perspective.

Figure 12. Qualitative comparison of multi-phase training with one-phase training for scene generalization. The arrows represent the direction of mouse movements.