

# MineDreamer: Learning to Follow Instructions via Chain-of-Imagination for Simulated-World Control

Enshen Zhou<sup>1,2\*</sup>, Yiran Qin<sup>1,3\*</sup>,  
Zhenfei Yin<sup>1,4</sup>, Yuzhou Huang<sup>3</sup>, Ruimao Zhang<sup>3†</sup>, Lu Sheng<sup>2†</sup>,  
Yu Qiao<sup>1</sup>, Jing Shao<sup>1‡</sup>

<sup>1</sup> Shanghai Artificial Intelligence Laboratory

<sup>2</sup> Beihang University

<sup>3</sup> The Chinese University of Hong Kong, Shenzhen (CUHK-Shenzhen)

<sup>4</sup> The University of Sydney

zhouenshen@buaa.edu.cn yiranqin@link.cuhk.edu.cn

<https://sites.google.com/view/minedreamer/main>

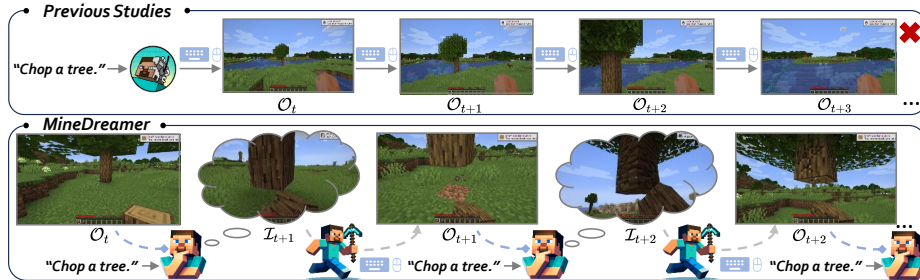
**Abstract.** It is a long-lasting goal to design a generalist-embodied agent that can follow diverse instructions in human-like ways. However, existing approaches often fail to steadily follow instructions due to difficulties in understanding abstract and sequential natural language instructions. To this end, we introduce *MineDreamer*, an open-ended embodied agent built upon the challenging Minecraft simulator with an innovative paradigm that enhances instruction-following ability in low-level control signal generation. Specifically, *MineDreamer* is developed on top of recent advances in Multimodal Large Language Models (MLLMs) and diffusion models, and we employ a Chain-of-Imagination (CoI) mechanism to envision the step-by-step process of executing instructions and translating imaginations into more precise visual prompts tailored to the current state; subsequently, the agent generates keyboard-and-mouse actions to efficiently achieve these imaginations, steadily following the instructions at each step. Extensive experiments demonstrate that *MineDreamer* follows single and multi-step instructions steadily, significantly outperforming the best generalist agent baseline and nearly doubling its performance. Moreover, qualitative analysis of the agent’s imaginative ability reveals its generalization and comprehension of the open world.

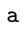
**Keywords:** Chain-of-Imagination · multimodal large language model · instruction following · low-level control

## 1 Introduction

One of the core objectives of current embodied intelligence is to develop a generalist low-level control agent that can follow diverse instructions to solve endless open-world embodied tasks [4, 5, 8, 42, 57]. Recent studies [4, 5, 8, 42] successfully unlock the instruction-following ability of foundation models [3, 12, 15] in

\* Equal contribution † Corresponding author ‡ Project leader



**Fig. 1: Comparison between *MineDreamer* and previous studies.** In “Chop a tree”  task, *MineDreamer* employs a Chain-of-Imagination mechanism, where it imagines step by step what to do next tailored to the current state. Imaginations contain environmental understanding and physical rules (e.g., perspective-based size changes). These can serve as more precise visual prompts to steadily guide the agent in generating actions to achieve these imaginations as effectively as possible at each step. Previous approaches have seen a tree, but missed the opportunity to chop it down.

the sequential decision-making domain [8, 11, 33, 57, 63, 69, 73]. However, these methods [5, 42] struggle to enable agents to follow textual instructions steadily, due to the: (1) Many textual instructions are abstract for low-level control and models struggle to effectively understand. They should be transformed into more effective prompts that consider how to execute instructions based on the current state. Hence, simple textual instructions cannot provide a precise demonstration of the desired behavior. (2) Many textual instructions are sequential, and executing them may require considering the current state and breaking down the task into multiple stages for step-by-step completion. Therefore, steady action generation driven by single-text instructions often fails.

To address the above issues, this work aims to explore how to unlock the situation-aware reasoning ability for a pre-trained decision-making foundation model. We introduce a simple yet effective mechanism called Chain-of-Imagination (CoI), which enables the agent to imagine and act upon the next stage step by step according to the instructions. Our method is motivated by two ideas: (1) When solving complex problems, humans often envision the goal of the next stage based on the current state. If we can break down the sequential instructions into multiple stages according to the current state, step by step, we can enable agents to follow instructions steadily. (2) Inspired by prompt tuning [34, 78, 79], if we can provide visual prompts containing physical rules and environmental understanding for each imagined step, tailored to optimally describe the desired behavior in the current state, which are more intuitive and efficient than task instructions, we can better guide the foundation model in predicting actions.

To this end, we propose *MineDreamer* within Minecraft, which generates a series of “imagined” sub-steps based on the textual instructions and current state. These visual sub-steps are then fed into a pre-trained decision-making foundation model to generate low-level control actions aimed at achieving the sub-steps. Specifically, *MineDreamer* comprises three modules: (1) An Imaginator, a diffusion model enhanced by a Multimodal Large Language Model (MLLM), can better generate imaginations that contain the physical rules and environmental

understanding. **(2)** A Prompt Generator, the bridge between Imaginator and PolicyNet, can convert future imaginations into latent visual prompts that offer more logical and precise demonstrations of the desired behavior. **(3)** A PolicyNet, a foundation model, can use latent prompts as guidance to predict actions for agents in an open-world environment.

Notably, as shown in Fig. 1, *MineDreamer* leverages a Chain-of-Imagination mechanism through multi-turn interaction between the Imaginator and the PolicyNet and cyclically generates latent visual prompts that better align with the current state to guide the PolicyNet in following instructions steadily in action generation. This mechanism represents an attempt to implement “self multi-turn interaction” in the sequential decision-making domain. Training an Imaginator in an open-world environment to envision the image of the next step requires extensive data. We employ the *Goal Drift Collection* method to gather a large amount of egocentric embodied data, which helps the Imaginator to understand how to achieve the instruction sequentially and how to achieve it repeatedly.

Our main contributions are as follows:

- We introduce the Chain-of-Imagination(CoI) method, which introduces “self multi-turn interaction” to the sequential decision-making domain and enables the agent to follow human instructions steadily in action generation.
- We propose the *Goal Drift Collection* method and an MLLM-enhanced diffusion model that can generate imaginations adhering to physical rules and environmental understanding, providing more precise visual prompts relevant to the current state and instructions.
- Leveraging these methods, we create an embodied agent in Minecraft named *MineDreamer* that has achieved nearly double the performance of the best generalist agent baseline in executing single and multi-step instructions steadily.

## 2 Related Work

### 2.1 Build Instruction-Following Agents in Minecraft

Research on generalist agents in Minecraft’s complex and dynamic environment is increasingly popular in AI. Despite the exploration of Large Language Models [7, 14, 48, 54, 65, 66] as high-level task planners that guide agents in executing long-horizon tasks [25, 50, 70–72, 81] like Voyager [70] and MP5 [50], we still require lower-level controllers [3, 8, 20, 27, 42] to execute the generated plans. In the sequential decision-making domain, DreamerV3 [27] trains agents using a world model, while VPT [3] builds a large foundational model to generate actions by learning from extensive video data. However, neither can follow instructions. GROOT [8] is developed to follow video instructions but fails to follow text instructions. STEVE-1 [42], an evolution of VPT [3], is built for text instructions but struggles to understand natural language prompts, despite extensive prompt engineering. Therefore, we create *MineDreamer*, which, leveraging the Chain-of-Imagination mechanism, generates more precise visual prompts step-by-step, enabling it to follow instructions steadily in action generation.

## 2.2 Conditioned Diffusion Models in Embodied Scenario

With the development of the text-to-image diffusion model [18, 30, 46, 55, 58, 60], the instruction-based diffusion methods [6, 9, 21, 23, 29, 32, 35, 67, 76] have recently marked considerable progress in generative tasks, especially in embodied scenarios. UniPi [19] and HiP [1] integrate video diffusion with inverse dynamics to generate robot control signals for specific tasks. SkillDiffuser [41] applies interpretable hierarchical planning via skill abstractions in diffusion-based task execution. While existing methods can only handle embodied tasks limited to fixed environments, the emergence of Multimodal Large Language Models (MLLMs) [13, 22, 43, 49, 62, 74, 75, 80] has showcased superior reasoning and perceptual abilities in open-world environment. Inspired by this, we create an MLLM-enhanced diffusion model, focusing on the model’s understanding of physics rules and environmental understanding, and its ability to create high-quality egocentric images for guiding low-level action generation.

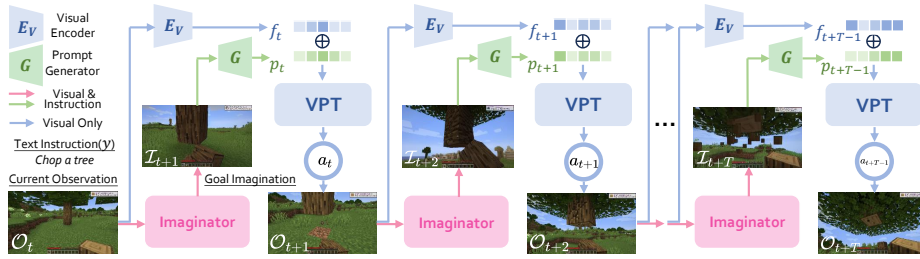
## 3 Method

In this section, we first provide an overview (Sec. 3.1) of our *MineDreamer*, including its mechanisms and features. Next, we introduce the purpose and workflow of the Chain-of-Imagination (CoI) mechanism (Sec. 3.2) regarding Fig. 2. To implement CoI and collect extensive embodied data to train Imaginator, we elaborate on the dataset construction (Sec. 3.3), including *Goal Drift Collection* method. Finally, we provide the necessary details of each part, including Imaginator (Sec. 3.4), Prompt Generator, and PolicyNet (Sec. 3.5).

### 3.1 Overview

Our *MineDreamer* comprises three modules, *i.e.*, Imaginator, Prompt Generator, and PolicyNet. Our objective is to empower agents, especially foundation models in the sequential decision-making domain, to follow human instructions steadily and act accordingly. The Imaginator is a parameter-efficiently fine-tuned diffusion model specific to Minecraft utilizing the visual reasoning ability of a Multimodal Large Language Model (MLLM). The Prompt Generator reconstructs latent visual prompts from the current observations, future imaginations, and instructions. PolicyNet is the existing Video Pretraining (VPT) [3] model, trained on 70k hours of Minecraft gameplay.

**Why future goal imagination?** Given a pre-trained model that can predict actions, the intuitive approach is to input the current state and instructions to guide it directly. So why the future goal imagination? In practice, we find that future goal imagination proves more interpretable for humans, easing debugging, and improving interaction and safety assessment [40, 51, 56, 77]. Furthermore, images yield flexible, explicit representations, facilitating natural language goal decomposition into clearer stages by learned physical rules and environmental understanding, helping the low-level control model “plan” what to do now.



**Fig. 2: The Overview of Chain-of-Imagination.** The Imaginator imagines a goal imagination based on the instruction and current observation. The Prompt Generator transforms this into a precise visual prompt, considering both the instruction and observed image. The Visual Encoder encodes the current observation, integrates it with this prompt, and inputs this into VPT. VPT then determines the agent’s next action, leading to a new observation, and the cycle continues. Note that VPT’s input is historical observations, so the figure cannot fully represent the autoregressive process. More details about VPT as PolicyNet can be found in Sec. 3.5.

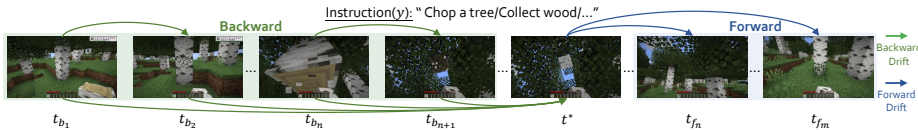
**Why can *MineDreamer* follow instructions more steadily?** Firstly, *MineDreamer* employs a Chain-of-Imagination (CoI) mechanism for incremental goal achievement via self-multi-turn interactions, enabling the agent to appropriately respond to the current state. In addition, with the help of this mechanism, the Prompt Generator crafts logical latent visual prompts that provide clear demonstrations of desired behaviors, ensuring that the agent steadily follows instructions. Furthermore, the enhanced Imaginator not only comprehends open-ended visual concepts, enabling it to imagine images of novel instructions it has never seen before but also ensures these images adhere to physical rules and environmental understanding, thereby sharpening the precision of prompts. Thus, *MineDreamer* can follow instructions steadily in an open-world environment.

### 3.2 Chain-of-Imagination

Chain-of-Imagination (CoI) enables the agent to envision the steps needed to achieve a goal iteratively. As shown in Fig. 2, it is an example to demonstrate how CoI works. First, the Imaginator takes in the user’s instructions  $y$  and current observations  $\mathcal{O}_t$  and imagines a future image  $\mathcal{I}_{t+1}$  depicting a moment within the process of completing the given instruction  $y$ , which is closely related to the current observation  $\mathcal{O}_t$ . Next, the Prompt Generator progressively creates a more precise latent visual prompt  $p_t$  in awareness of the current observation  $\mathcal{O}_t$ , instruction  $y$  and future imagination  $\mathcal{I}_{t+1}$ , aligning with the visual input space of the Video Pretraining (VPT) [3] model. The Visual Encoder then processes  $\mathcal{O}_t$  into a representation  $f_t$ , which is combined with  $p_t$  and fed into VPT [3]. Finally, VPT [3] progressively predicts an action (*i.e.*, keyboard and mouse) from the observation history, interacts with the environment, gathers a new observation  $\mathcal{O}_{t+1}$ , and repeats the cycle later.

### 3.3 Datasets

We train the Imaginator with the Goal Drift Dataset, which includes 500k triplets (current observation, future goal imagination, instruction) from the OpenAI Contractor Gameplay Dataset [3], using the *Goal Drift Collection* method.



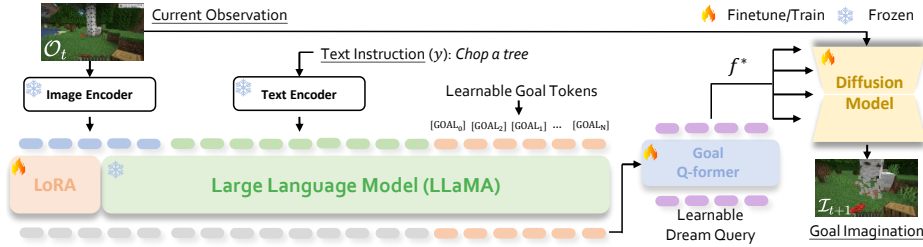
**Fig. 3: Goal Drift Collection.** For each timestamp  $t^*$ , we form many triplets comprising (current observation, goal imagination, instruction) associated with the game event-related instructions completed by contractors. Each pair of linked images forms a training triplet with its instruction for the Imaginator in this figure.

**OpenAI Contractor Gameplay Dataset.** OpenAI Contractor Gameplay Dataset [3] is created by hiring human contractors to play Minecraft and complete tasks like house building. Game events, like “mine\_block”, noting the type of block broken, are logged with timestamps. These timestamps ( $t^*$ ) provide precise progress tracking and align with completed event-related instructions.

**Goal Drift Collection.** The Gameplay Dataset allows us to construct numerous embodied data by using specific event-related instructions achieved at each timestamp  $t^*$ . Yet, directly pairing images from these timestamps  $t^*$  as future goal imaginations  $\mathcal{O}_{t^*}$  with images from a fixed timestep  $T$  earlier as current observations  $\mathcal{O}_{t^*-T}$ , along with instruction  $y$ , could lead to certain problems: **(1) Goal Illusion: The Imaginator edits the observation to depict the completed instruction.** Training the Imaginator on such data may reduce it to an image editor, as it generates imaginations without regard to the environment because all goal imaginations in the dataset represent the moment when instruction is completed. For instance, given the instruction “Break dirt”  $\blacksquare$  while facing the sky, the Imaginator may unrealistically insert a broken dirt block  $\blacksquare$  into the sky. **(2) Imagination Stagnation: The Imaginator fails to conceive repeated task completion.** The Imaginator is trained to envision the instructions’ fulfillment once, not recognizing the need for repetition, as all current observations precede the achievement of instructions. For instance, given “Chop a tree”  $\blacklozenge$ , after cutting the uppermost wood  $\blacksquare$  by looking up, the agent will not look down for more trees  $\blacklozenge$ , impeding continuous task performance.

To address the aforementioned issues, we propose the *Goal Drift Collection* method to gather Goal Drift Dataset. From the Gameplay Dataset, we form many triplets (current observation, goal imagination, instruction) at each timestamp  $t^*$ , all associated with the same event-related instructions  $y$  completed by the contractors. Fig. 3 shows that a pair of linked images with instructions  $y$  constitutes a training triplet. Our approach has both Backward Drift, which helps the model understand the step-by-step completion of tasks to mitigate Goal Illusion, and Forward Drift, which enables the model to learn how to accomplish instructions repeatedly to reduce Imagination Stagnation. The details of collecting three kinds of data samples corresponding to each  $t^*$  are as follows:

1. Backward Drift 1: We set  $t_{b_1}$  as  $t^*$  backward by fixed  $T_b$  time steps and then select  $m - 2$  random timestamps between  $t_{b_1}$  and  $t^*$  to form the sequence  $t_{b_1}, \dots, t_{b_m}$ , where  $t^*$  is  $t_{b_m}$ . At each time step, the current and next observations are paired as the current observations and goal imagination, respectively, which can form  $m - 1$  samples.



**Fig. 4: The Overall Framework of Imaginator.** For the goal understanding, we add  $k$  [GOAL] tokens to the end of instruction  $y$  and input them with current observation  $\mathcal{O}_t$  into LLaVA [43]. Then LLaVA [43] generates hidden states for the [GOAL] tokens, which the Q-Former processes to produce the feature  $f^*$ . Subsequently, the image encoder  $\mathbf{E}_v$  combines its output with  $f^*$  in the diffusion models for instruction-based future goal imagination generation.

2. Backward Drift 2: In  $t_{b_1}, \dots, t_{b_m}$ , the observations at each timestamp except for  $t_{b_m}$  are used as the current observations, and the observation at  $t^*$  serve as the goal imagination, which can form  $m - 1$  samples.
3. Forward Drift: We set  $t_{f_m}$  as  $t^*$  forward by fixed  $T_f$  time steps and randomly select  $m - 2$  timestamps between  $t^*$  and  $t_{f_m}$ , where  $t^*$  is  $t_{f_1}$ . The observation at  $t^*$  serves as the current observation, and the observations at future timestamps serve as the goal imaginations, which can form  $m - 1$  samples.

For more details about the dataset and collection method, please check Supp. B.

### 3.4 Imaginator

Inspired by prompt tuning [34, 78, 79], we introduce Imaginator, an MLLM-enhanced diffusion model that imagines step by step what to do next based on the current state and instruction, enabling the creation of more precise visual prompts for improved low-level control demonstrations of the desired behavior. Imaginator’s training data utilizes the Goal Drift Dataset from Sec 3.3, consisting of (current observation, goal imagination, instruction) triplets.

**Goal Understanding via Task Instruction Following.** Given a current observation  $\mathcal{O}_t$  and a textual instruction  $y$ , the Imaginator generates a future goal imagination  $\mathcal{I}_{t+1}$  for the PromptGenerator’s visual prompt. In Fig. 4, current observation  $\mathcal{O}_t$  is encoded by a frozen image encoder  $\mathbf{E}_v$  into  $\mathbf{E}_v(\mathcal{O}_t)$ , textual instruction  $y$  is tokenized into  $(x_1, \dots, x_T)$ , they are sent to the LLM together. Imaginator now can acquire a goal imagination of the instruction intention but are limited to the language modality. Inspired by GILL [37], we bridge the language-vision modalities gap by extending the LLM’s vocabulary with  $k$  Learnable Goal Tokens  $[\text{GOAL}_1], \dots, [\text{GOAL}_k]$ , appending them to instruction  $y$ . Specifically, a trainable matrix  $\mathbf{E}_g$ , representing these [GOAL] embeddings, is added to the LLM’s embedding matrix. We aim to minimize the negative log-likelihood of predicting the next [GOAL] token given previously generated [GOAL] tokens:

$$\mathcal{L}_{\text{LLM}} = - \sum_{i=1}^k \log p_{\{\theta_L \cup \theta_t \cup \mathbf{E}_g\}}([\text{GOAL}_i] \mid \mathbf{E}_v(\mathcal{O}_t), x_1, \dots, x_T, [\text{GOAL}_1], \dots, [\text{GOAL}_{i-1}]) \quad (1)$$

We add LoRA [31] parameters  $\theta_l$  into the LLM’s self-attention projection layers for efficient fine-tuning while keeping all LLM parameters  $\theta_L$  frozen. During training, only the LoRA [31] parameters  $\theta_l$  and the Learnable Goal Tokens  $\mathbf{E}_g$  are updated. The hidden states  $h_{[\text{GOAL}]}$  corresponding to  $\mathbf{E}_g$  tokens are used to generate imaginations in the following module.

**Goal Imagination Generation via Latent Imagination.** To address the disparity between the LLM’s hidden states and the CLIP [53] text encoder’s feature spaces, we must transform the LLM’s sequential goal tokens into semantically relevant representations for guiding goal imagination generation. Inspired by BLIP2 [39] and InstructBLIP [16], we employ a Goal Q-Former  $\mathcal{Q}$  with several Learnable Dream Query, to derive the goal imagination representation  $f^*$ :

$$f^* = \mathcal{Q}(h_{[\text{GOAL}]}) \quad (2)$$

To enhance goal imagination with representation  $f^*$  to guide imagination generation, we utilize a latent diffusion model combining a variational autoencoder (VAE) [36] for latent space denoising diffusion. Drawing from Instruct-Pix2Pix’s [6] latent diffusion approach, a cornerstone in instruction-based image editing, our model introduces noise to the latent encoding  $z = \mathcal{E}(\mathcal{I}_{t+1})$  of the goal imagination  $\mathcal{I}_{t+1}$  through encoder  $\mathcal{E}$ , yielding a noisy latent  $z_s$  across timesteps  $s \in S$ . A U-Net [59]  $\epsilon_\delta$  is trained to estimate this noise, conditional on the current observation  $c_o = \mathcal{E}(O_t)$  and text instruction  $c_T$ , by merging  $c_o$  with  $z_s$ . The specific process can be formulated as follows:

$$\mathcal{L}_{\text{dream}} = \mathbb{E}_{\mathcal{E}(\mathcal{I}_{t+1}), \mathcal{E}(O_t), c_T, \epsilon \sim \mathcal{N}(0,1), s} [\|\epsilon - \epsilon_\delta(s, \text{concat}[z_s, \mathcal{E}(O_t)] + f^*)\|_2^2] \quad (3)$$

where  $\epsilon$  is unscaled noise,  $s$  is the sampling step,  $z_s$  is latent noise at step  $s$ ,  $\mathcal{E}(O_t)$  is the current observation condition, and  $c_T$  is the text instruction condition. The `concat` corresponds to the concatenation operation.

### 3.5 Prompt Generator and PolicyNet

To transform goal imaginations into precise latent visual prompts that the PolicyNet can understand, we require a Prompt Generator to serve as the bridge between the Imaginator and the PolicyNet. Inspired by STEVE-1 [42], our prompt generator is a conditional variational autoencoder (CVAE) [36, 64] model trained on the Goal Drift subset dataset. It encodes the current observations, goal imaginations, and instructions by MineCLIP [20] to produce three embeddings. These embeddings are then reconstructed into a latent visual embedding within the MineCLIP [20] visual space and a linear layer then projects it into the visual input space of our PolicyNet.

In our PolicyNet, we utilize the architecture of the existing model named VPT [3] and the training parameters of STEVE-1 [42]. Specifically, as shown in Fig. 2, we first process the current observation with a Visual Encoder (*i.e.*, ResNet [28]) of VPT [3] and get representation  $f_t$ . After adding it with the latent visual prompts  $p_t$  generated by the Prompt Generator, the sum result  $o_t$  is then fed into the PolicyNet. PolicyNet, whose backbone is Transformer-XL [17], processes the current input representations  $o_t$  and autoregressively predicts the next action  $a_t$ . We can describe the process where the Prompt Generator creates



latent visual prompts  $p_t$  and PolicyNet predicts the next action  $a_t$  based on them and historical observations using the following simple notation:

$$p_t \leftarrow \mathcal{G}(\mathcal{O}_t, \mathcal{I}_{t+1}, y), \quad f_t \leftarrow \mathcal{V}(\mathcal{O}_t), \quad o_t \leftarrow f_t + p_t, \quad a_t \leftarrow \mathcal{T}(o_{t-T}, \dots, o_t) \quad (4)$$

where  $\mathcal{G}$  is PromptGenerator,  $\mathcal{V}$  is VisualEncoder, and  $\mathcal{T}$  is TransformerXL [17].

## 4 Experiments

### 4.1 Experimental Setup

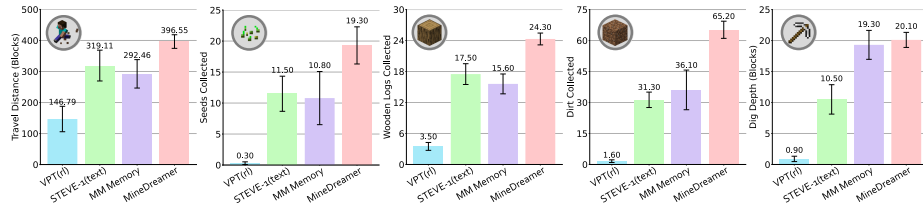
**Training Process.** The training process of Imaginator is divided into three main stages. In the first stage, the MLLM is aligned with the CLIP [54] text encoder [53] using the QFormer [39]. In the second stage, we apply Instruct-Pix2Pix [6] to warm up the weights for the diffusion model in Minecraft. In the third stage, we optimize Imaginator in an end-to-end manner. To be specific, the weights of LLaVA [43] are frozen and LoRA [31] is added for efficient fine-tuning. For the diffusion model, we directly use the weights pre-trained in the second stage as the initial weights in Imaginator. The CVAE [36, 64] within the Prompt Generator features a Gaussian prior and a Gaussian posterior, with its encoder and decoder, parameterized as three-layer MLPs, each with 512 hidden units and layer normalization [2], similar to the architecture of STEVE-1’s [42] prior. More training details can be found in Supp. C.

**Training Datasets.** In the first stage of Imaginator, we use the extensive corpus CC12M [10], and our Goal Drift Dataset is used in the second and third stages. We follow STEVE-1’s [42] approach for CVAE [36, 64] training, curating a subset of approximately 10k quadruplets from the Goal Drift Dataset for our test tasks. This subset includes current observations, goal imaginations, and instructions that match the Goal Drift Dataset. We use the MineCLIP [20] video encoder to transform the goal imagination and the previous 16 frames into a visual prompt embedding, which acts as the ground truth. More details can be found in Supp. B.

**Environment Setting.** We employ MineRL [26] as the Minecraft simulation. The observation space is limited to RGB images, and the action space is confined to keyboard and mouse controls, which are consistent with human interaction. For more details about the simulator, please check Supp. A.

**Baseline.** We compare *MineDreamer* with three baseline:

1. VPT [3], a foundation model pretrained on 70k hours gameplay. Here, we select the VPT(rl), which is finetuned by reinforcement learning on the original VPT [3] foundation model but **cannot follow instructions**.
2. STEVE-1 [42], an instruction-following agent finetuned from VPT(rl). Here, we select STEVE-1(text), which uses a simple prior to aligning the text with the visual space, **without considering the current observation**.
3. Multi-Modal Memory, a substitute for the Imaginator and Prompt Generator in *MineDreamer*, efficiently searches through extensive instruction-video pairs to find the most relevant video as a visual prompt based on the given instruction and the current observation, which effectively **leverages the current observation and incorporates a CoI mechanism**.



**Fig. 5: Performance on Programmatic Evaluation.** *MineDreamer* surpasses the unconditional VPT [3], the text-conditioned STEVE-1 [42] that ignores current state, and the Multi-Modal Memory that utilizes current state with a CoI mechanism.

For more details about the baseline, please check Supp. D.1.

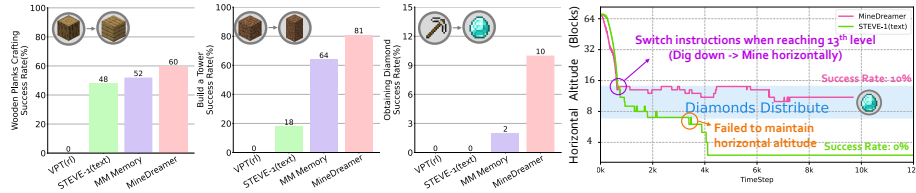
**Evaluation.** We utilize STEVE-1’s [42] *early-game evaluation suite*, which comprises two evaluations: (1) Programmatic Evaluation, a quantitative evaluation used to evaluate an agent’s ability to execute **single-step** instruction steadily. We track the states provided by the simulator to calculate metrics (*e.g.*, wooden log collection, travel distance). (2) Command-Switching Evaluation, a quantitative evaluation designed to assess whether the agent can successfully execute **multi-step** instructions in sequence to complete long-horizon tasks (*e.g.*, obtaining diamond). We use the success rate as the metric for evaluation. More evaluation details can be found in Supp. D.2 and Supp. D.3.

## 4.2 Performance on Textual Instructions Control

**Programmatic Evaluation.** We quantitatively evaluate all agents on 5 tasks and plot the programmatic metric performances (mean and 95% confidence intervals). Each task runs 10 trials with distinct environment seeds, limiting 3,000 frames (*i.e.*, 2.5 minutes of gameplay) which are consistent with STEVE-1 [42]. Unlike STEVE-1 [42], we condition all agents with the most suitable biome.

Fig. 5 compares the performance of our *MineDreamer* with the unconditional VPT [3], the text-conditioned STEVE-1 [42] and *MineDreamer* using Multi-Modal Memory. With appropriate text instructions, *MineDreamer* significantly outperforms the unconditional VPT [3], collecting  $64\times$  more seeds,  $7\times$  more wood,  $41\times$  more dirt, traveling  $2.7\times$  further, and digging  $22\times$  deeper. It also surpasses the STEVE-1 [42], collecting  $1.7\times$  more seeds,  $1.4\times$  more wood,  $2.1\times$  more dirt, traveling  $1.2\times$  further, and digging  $1.9\times$  deeper. Compared to Multi-Modal Memory, *MineDreamer* collects  $1.8\times$  more seeds,  $1.5\times$  more wood,  $1.8\times$  more dirt, travels  $1.3\times$  further, and digs  $1.1\times$  deeper. This demonstrates that our CoI mechanism, which breaks down instructions into multiple stages and executes them step by step, leads to steadier instruction following compared to STEVE-1 [42] which uses direct text instruction guidance. Unlike Multi-Modal Memory, which also features the CoI mechanism, our method generates future imaginations that closely resemble the current state at each stage, resulting in providing more precise visual prompts of the desired behavior, thus enhancing the stability of action generation.

We also observe an interesting phenomenon: while Multi-Modal Memory, using the CoI mechanism and current observations, outperforms unconditional



**Fig. 6: Performance on Command-Switching Evaluation.** (Left) *MineDreamer* swiftly adapts to instructions and follows them steadily, achieving a higher success rate than the unconditional VPT [3], the text-conditioned STEVE-1 [42], and the Multi-Modal Memory with CoI mechanism. (Right) *MineDreamer* can dig down  $\downarrow$  to a depth of 13 and steadily mine horizontally  $\rightarrow$  to obtain diamonds  $\diamond$  with an average success rate of 10%, while STEVE-1 [42] struggles to maintain a consistent altitude.

VPT [3], it sometimes underperforms compared to STEVE-1 [42]. Upon reviewing the recorded videos and the results of memory retrieval, we find that due to the vast diversity of open-world environments, the videos retrieved by Multi-Modal Memory still exhibit slight differences from the current state. This discrepancy misguides the PolicyNet in predicting agent actions, indicating that the CoI’s effectiveness hinges on the relevancy and precision of future imaginations or visual prompts to the current state.

**Command-Switching Evaluation for Long-Horizon Tasks.** In this part, we explore agents’ ability to solve long-horizon tasks that require executing multi-step instructions in sequence, including (1) collect wood  $\log$  and then craft planks  $\text{plank}$ , (2) gather dirt  $\text{dirt}$  and then build a tower  $\text{tower}$  and (3) dig down  $\downarrow$  and then mine horizontally  $\rightarrow$  for diamonds  $\diamond$ , each with 50 trials. Tasks 1 and 2 limits 3,000 frames (*i.e.*, 2.5 minutes of gameplay), with instructions changing at 1,500 and 2,000 frames. Task 3 limits 12,000 frames (*i.e.*, 10 minutes of gameplay), switching instructions upon reaching the 13th floor, as diamonds  $\diamond$  are commonly found between the 7th and 14th floors.

In Fig. 6 (Left), *MineDreamer* consistently surpasses VPT [3] and STEVE-1 [42] in Command-Switching tasks. VPT’s [3] inability to follow instructions leads to a complete failure in executing sequential instructions, as evidenced by a 0% success rate in the evaluation. Although STEVE-1 [42] occasionally completes Command-Switching tasks, it underperforms compared to *MineDreamer*. For instance, in the Obtain diamond  $\diamond$  task, STEVE-1’s [42] success rate is 0%, while Multi-Modal Memory’s success rate is 2%, notably lower than *MineDreamer*’s 10%. As shown in Fig. 6 (Right), we reconstruct an instance where two agents act in the same environment based on the simulator records. Initially, both *MineDreamer* and STEVE-1 [42] rapidly dig down  $\downarrow$  to the target depth and then mine horizontally  $\rightarrow$  to obtain diamonds  $\diamond$ . Compared to STEVE-1 [42], *MineDreamer* can consistently maintain the specified horizontal level over an extended period and successfully obtains diamonds  $\diamond$  around the 10k steps in this instance. While STEVE-1 [42] manages to maintain its specified horizontal level for a long time, it ultimately fails to do so and becomes stuck in the bedrock layer (*i.e.*, the agent cannot break any block), resulting in a 0% success rate. This demonstrates that, even when instructions are switched rapidly, the

CoI mechanism can still drive the agent to generate future goal imaginations that align with the current state. Visual prompts generated from these imaginations enable the agent to quickly adapt its actions to correspond with the new instructions while steadily following the instructions in action generation.




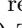
**Fig. 7:** Qualitative Comparison of Goal Imagination Generation. When compared to InstructPix2Pix [6] that have undergone further fine-tuning on our Goal Drift Dataset, our approach demonstrates superior goal imagination capabilities in embodied scenarios. See Sec. 4.3 for a more detailed analysis.

### 4.3 Qualitative Results of Imaginator

We compare Imaginator with the existing state-of-the-art instruction-based image editing model, namely InstructPix2Pix [6]. Given this model has been trained on specific datasets, its performance would inevitably be suboptimal if directly applied to the Minecraft domain. To facilitate a fair comparison, we fine-tune InstructPix2Pix [6] using the same training set employed by the Imaginator and assess the performance of the fine-tuned models in addressing tasks in Minecraft. Fig 7 shows qualitative results in the evaluation set, our methodology exhibits enhanced abilities in Goal Imagination Generation within intricate scenarios.

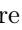
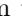
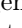
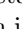
The first comparison shows that the Imaginator adeptly captures the agent’s perspective shift as it advances, whereas InstructPix2Pix [6] struggles to generate images in alignment with the provided instructions. In the second instance, the Imaginator specifically visualizes the region with felled trees 🌲, contrasting with InstructPix2Pix [6], which yields an image markedly divergent from the existing observation background. The third comparison highlights the Imaginator’s ability to depict enhanced visibility following torch placement, in contrast to InstructPix2Pix [6], which merely adds torches without the associated increase in illumination. These observations suggest that in scenarios requiring instruction reasoning and goal understanding, a simple CLIP [54] text encoder may struggle to guide the diffusion model to generate reasonable goal imagination.

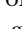


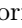
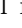



**Fig. 8: The Generalizability of *MineDreamer*.** (Left) Despite excluding data involving ‘Dirt’  or ‘Dig’  from Goal Drift Dataset and retraining, Imaginator can still generate relatively high-quality imaginations aligned with the instruction’s concept. (Right) The retrained Imaginator remains operational with the CoI mechanism and can handle unseen instructions while largely preserving its previous performance.



However, the MLLM can fully utilize its powerful reasoning ability, vast environmental knowledge, and intrinsic physical rules to correctly understand the goal and generate goal imagination. More visual results can be found in Supp. F and Supp. G.

#### 4.4 Discussion on Generalization

In this part, we will explore the generalizability of *MineDreamer*, as the agent’s ability to generalize is key to its behavior in the open world where environments are complex and instructions vary widely. Since STEVE-1 [42] has shown its prior ability to map text to visual prompts effectively, and our Prompt Generator is built upon it, we will now concentrate on the generalizability of our Imaginator and the entire agent. At first, we exclude data related to the words ‘Dirt’  or ‘Dig’  from the Goal Drift Dataset and retrain the model. Then, we observe the images generated in response to the instruction “Collect dirt”  based on the current state and the quantity of dirt  collected by the agent.



As shown in Fig. 8, we find that even after completely removing the concepts of ‘Dirt’  or ‘Dig’ , Imaginator is still able to generate goal imaginations of relatively good quality aligned with the instruction’s concept (*i.e.*, **agent points towards the dirt  and attempt to break it**), which can still guide the PolicyNet to follow instructions. The resulting collection of dirt  is about 70% of the original amount, which shows that the Imaginator can respond to unseen novel instructions while largely maintaining its previous performance. We attribute this to three key factors: (1) The MLLM within Imaginator has the relevant environmental knowledge to map the text ‘Dirt’  to its corresponding element in Minecraft images, recognizing its visual counterpart; (2) training data for related tasks, such as “Collect seeds” , enables the MLLM

**Table 1:** We study the impact of dataset collection methods on agent performance. Values in parentheses represent 95% confidence intervals.

Instruction	Fixed Timestep Backward	Only Backward Drift	Only Forward Drift	Normal
“Chop a tree” 	7.60(3.84, 11.36)	10.10(2.82, 5.58)	4.20(2.82, 5.58)	<b>24.30(21.71, 26.89)</b>
“Collect dirt” 	38.60(21.97, 55.23)	30.30(20.71, 39.89)	18.10(6.74, 29.46)	<b>65.20(55.81, 74.59)</b>

to comprehend the meaning of action ‘Collect’ in Minecraft task; **(3)** The pre-trained Diffusion model can generalize to the Minecraft domain and generate goal imaginations by leveraging the MLLM’s latent representations for understanding textual semantics mentioned above from the instructions.

**Table 2:** We study the impact of the Chain-of-Imagination and diffusion model ability on agent performance. Values in parentheses represent 95% confidence intervals.

Instruction	wo CoI	Random Noise	Instruct-Pix2Pix	Normal
“Chop a tree” 	18.70(15.26, 22.14)	2.70(0.85, 4.55)	22.90(20.17, 25.63)	<b>24.30(21.71, 26.89)</b>
“Collect dirt” 	53.50(36.93, 70.07)	10.90(3.95, 17.85)	59.50(54.00, 65.00)	<b>65.20(55.81, 74.59)</b>

#### 4.5 What Contributes to Performance

**Dataset Collection Method.** In Tab. 1, we study the impact on agent performance by training with datasets of equal size collected using fixed Backward timesteps, only Backward Drift, only Forward Drift, and normal *Goal Drift Dataset Collection*. Although data collected using the first three methods can enable the agent to follow instructions, the Imaginator is affected by Goal Illusion and Imagination Stagnation, which are discussed in Sec. 3.3. This results in the Imaginator’s inability to envision the step-by-step process of completing the instruction and how to steadily complete the instruction multiple times.


**Chain-of-Imagination.** In Tab. 2, we explore the effect of the CoI mechanism on agent performance, where “wo-CoI” denotes the scenario where the agent generates the goal imagination and visual prompt only at the beginning and remains unchanged thereafter. Compared to normal performance, “wo-CoI” achieves about 77%. This is because the visual prompts generated at the beginning become less capable of providing precise demonstrations of the desired behavior in later stages, resulting in hindering the ability to guide the agent step by step more steadily.

**Diffusion Model Ability.** In Tab. 2, we explore the impact of diffusion model ability on performance. Using “random noise” as a goal imagination results in vague visual prompts, which drastically reduce performance to merely 10% of its original level. The performance of InstructPix2Pix [6] and our MLLM-enhanced diffusion model are comparable; however, by leveraging MLLM, our generated images adhere more closely to physical rules and environmental knowledge, as shown in Fig. 7. Additionally, as discussed in Sec. 4.2, we discover that the CoI mechanism demands a certain quality of goal imagination, suggesting that the stronger the Imaginator, the better it can guide agents to follow instructions.

More ablation studies can be found in Supp. E.

## 5 Conclusion and Limitation

In this paper, we introduce an innovative paradigm for enhancing the instruction-following ability of agents in simulated-world control. We prove that by employing a Chain-of-Imagination mechanism to envision the step-by-step process of executing instructions, and translating imaginations into precise visual prompts

tailored to the current state and instruction, can significantly help the foundation model follow instructions steadily in action generation. Our Agent, *MineDreamer* in Minecraft, showcases its strong instruction-following ability. Furthermore, we show its potential as a high-level planner’s downstream controller in the challenging “Obtain diamond”  task. We believe this novel paradigm will inspire future research and generalize to other domains and open-world environments.

**Limitation.** Firstly, generating high-quality imagination can take seconds, slowing down frequent-use scenarios. Speed enhancements via distillation [61] and quantization [24] may mitigate this. Secondly, the Imaginator may produce unrealistic hallucinations. Integrating world knowledge via methods such as RAG [38] or reducing MLLM hallucinations [45] could mitigate this.

## References

1. Ajay, A., Han, S., Du, Y., Li, S., Gupta, A., Jaakkola, T., Tenenbaum, J., Kaelbling, L., Srivastava, A., Agrawal, P.: Compositional foundation models for hierarchical planning. *Advances in Neural Information Processing Systems* **36** (2024) [4](#)
2. Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization. *arXiv preprint arXiv:1607.06450* (2016) [9](#)
3. Baker, B., Akkaya, I., Zhokov, P., Huizinga, J., Tang, J., Ecoffet, A., Houghton, B., Sampedro, R., Clune, J.: Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems* **35**, 24639–24654 (2022) [1](#), [3](#), [4](#), [5](#), [6](#), [8](#), [9](#), [10](#), [11](#), [2](#), [7](#), [12](#)
4. Brohan, A., Brown, N., Carbajal, J., Chebotar, Y., Chen, X., Chormanski, K., Ding, T., Driess, D., Dubey, A., Finn, C., et al.: Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818* (2023) [1](#)
5. Brohan, A., Brown, N., Carbajal, J., Chebotar, Y., Dabis, J., Finn, C., Gopalakrishnan, K., Hausman, K., Herzog, A., Hsu, J., et al.: Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817* (2022) [1](#), [2](#)
6. Brooks, T., Holynski, A., Efros, A.A.: Instructpix2pix: Learning to follow image editing instructions. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 18392–18402 (2023) [4](#), [8](#), [9](#), [12](#), [14](#), [6](#), [16](#), [17](#)
7. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. *Advances in neural information processing systems* **33**, 1877–1901 (2020) [3](#)
8. Cai, S., Zhang, B., Wang, Z., Ma, X., Liu, A., Liang, Y.: Groot: Learning to follow instructions by watching gameplay videos. *arXiv preprint arXiv:2310.08235* (2023) [1](#), [2](#), [3](#), [11](#)
9. Cao, M., Wang, X., Qi, Z., Shan, Y., Qie, X., Zheng, Y.: Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. *arXiv preprint arXiv:2304.08465* (2023) [4](#)
10. Changpinyo, S., Sharma, P., Ding, N., Soricut, R.: Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 3558–3568 (2021) [9](#)

11. Chen, L., Lu, K., Rajeswaran, A., Lee, K., Grover, A., Laskin, M., Abbeel, P., Srinivas, A., Mordatch, I.: Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems* **34**, 15084–15097 (2021) [2](#)
12. Chen, X., Djolonga, J., Padlewski, P., Mustafa, B., Changpinyo, S., Wu, J., Ruiz, C.R., Goodman, S., Wang, X., Tay, Y., et al.: Pali-x: On scaling up a multilingual vision and language model. *arXiv preprint arXiv:2305.18565* (2023) [1](#)
13. Chen, Z., Wang, Z., Wang, Z., Liu, H., Yin, Z., Liu, S., Sheng, L., Ouyang, W., Qiao, Y., Shao, J.: Octavius: Mitigating task interference in mlms via moe. *arXiv preprint arXiv:2311.02684* (2023) [4](#)
14. Chiang, W.L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J.E., et al.: Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023) (2023) [3](#)
15. Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H.W., Sutton, C., Gehrmann, S., et al.: Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research* **24**(240), 1–113 (2023) [1](#)
16. Dai, W., Li, J., Li, D., Tiong, A.M.H., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: Instructblip: Towards general-purpose vision-language models with instruction tuning (2023) [8](#)
17. Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q.V., Salakhutdinov, R.: Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860* (2019) [8](#), [9](#), [12](#)
18. Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. *Advances in neural information processing systems* **34**, 8780–8794 (2021) [4](#)
19. Du, Y., Yang, S., Dai, B., Dai, H., Nachum, O., Tenenbaum, J., Schuurmans, D., Abbeel, P.: Learning universal policies via text-guided video generation. *Advances in Neural Information Processing Systems* **36** (2024) [4](#)
20. Fan, L., Wang, G., Jiang, Y., Mandlekar, A., Yang, Y., Zhu, H., Tang, A., Huang, D.A., Zhu, Y., Anandkumar, A.: Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems* **35**, 18343–18362 (2022) [3](#), [8](#), [9](#), [2](#), [5](#), [7](#), [14](#)
21. Fu, T.J., Hu, W., Du, X., Wang, W.Y., Yang, Y., Gan, Z.: Guiding instruction-based image editing via multimodal large language models. *arXiv preprint arXiv:2309.17102* (2023) [4](#)
22. Gao, P., Han, J., Zhang, R., Lin, Z., Geng, S., Zhou, A., Zhang, W., Lu, P., He, C., Yue, X., et al.: Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010* (2023) [4](#)
23. Geng, Z., Yang, B., Hang, T., Li, C., Gu, S., Zhang, T., Bao, J., Zhang, Z., Hu, H., Chen, D., et al.: Instructdiffusion: A generalist modeling interface for vision tasks. *arXiv preprint arXiv:2309.03895* (2023) [4](#)
24. Gholami, A., Kim, S., Dong, Z., Yao, Z., Mahoney, M.W., Keutzer, K.: A survey of quantization methods for efficient neural network inference (2021) [15](#)
25. Gong, R., Huang, Q., Ma, X., Vo, H., Durante, Z., Noda, Y., Zheng, Z., Zhu, S.C., Terzopoulos, D., Fei-Fei, L., et al.: Mindagent: Emergent gaming interaction. *arXiv preprint arXiv:2309.09971* (2023) [3](#)
26. Guss, W.H., Houghton, B., Topin, N., Wang, P., Codel, C., Veloso, M., Salakhutdinov, R.: Minerl: a large-scale dataset of minecraft demonstrations. In: *Proceedings of the 28th International Joint Conference on Artificial Intelligence*. pp. 2442–2448 (2019) [9](#), [1](#), [3](#), [4](#), [10](#), [11](#)



27. Hafner, D., Pasukonis, J., Ba, J., Lillicrap, T.: Mastering diverse domains through world models. arXiv preprint arXiv:2301.04104 (2023) [3](#)
28. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016) [8](#)
29. Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., Cohen-Or, D.: Prompt-to-prompt image editing with cross attention control. arXiv preprint arXiv:2208.01626 (2022) [4](#)
30. Ho, J., Salimans, T.: Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598 (2022) [4](#), [11](#), [12](#), [13](#)
31. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models. arXiv preprint arXiv:2106.09685 (2021) [8](#), [9](#), [6](#)
32. Huang, Y., Xie, L., Wang, X., Yuan, Z., Cun, X., Ge, Y., Zhou, J., Dong, C., Huang, R., Zhang, R., et al.: Smartedit: Exploring complex instruction-based image editing with multimodal large language models. arXiv preprint arXiv:2312.06739 (2023) [4](#)
33. Janner, M., Li, Q., Levine, S.: Offline reinforcement learning as one big sequence modeling problem. Advances in neural information processing systems **34**, 1273–1286 (2021) [2](#)
34. Jia, M., Tang, L., Chen, B.C., Cardie, C., Belongie, S., Hariharan, B., Lim, S.N.: Visual prompt tuning. In: European Conference on Computer Vision. pp. 709–727. Springer (2022) [2](#), [7](#)
35. Kawar, B., Zada, S., Lang, O., Tov, O., Chang, H., Dekel, T., Mosseri, I., Irani, M.: Imagic: Text-based real image editing with diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6007–6017 (2023) [4](#)
36. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114 (2013) [8](#), [9](#), [5](#), [7](#), [14](#)
37. Koh, J.Y., Fried, D., Salakhutdinov, R.R.: Generating images with multimodal language models. Advances in Neural Information Processing Systems **36** (2024) [7](#)
38. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.t., Rocktäschel, T., et al.: Retrieval-augmented generation for knowledge-intensive nlp tasks. Advances in Neural Information Processing Systems **33**, 9459–9474 (2020) [15](#)
39. Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597 (2023) [8](#), [9](#), [6](#)
40. Li, L., Dong, B., Wang, R., Hu, X., Zuo, W., Lin, D., Qiao, Y., Shao, J.: Salad-bench: A hierarchical and comprehensive safety benchmark for large language models. arXiv preprint arXiv:2402.05044 (2024) [4](#)
41. Liang, Z., Mu, Y., Ma, H., Tomizuka, M., Ding, M., Luo, P.: Skilldiffuser: Interpretable hierarchical planning via skill abstractions in diffusion-based task execution. arXiv preprint arXiv:2312.11598 (2023) [4](#)
42. Lifshitz, S., Paster, K., Chan, H., Ba, J., McIlraith, S.: Steve-1: A generative model for text-to-behavior in minecraft. arXiv preprint arXiv:2306.00937 (2023) [1](#), [2](#), [3](#), [8](#), [9](#), [10](#), [11](#), [13](#), [5](#), [7](#), [12](#)
43. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. arXiv preprint arXiv:2304.08485 (2023) [4](#), [7](#), [9](#), [6](#)
44. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017) [6](#)

45. Minervini, P., et al.: awesome-hallucination-detection. <https://github.com/EdinburghNLP/awesome-hallucination-detection> (2014) 15
46. Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., Chen, M.: Glide: Towards photorealistic image generation and editing with text-guided diffusion models. arXiv preprint arXiv:2112.10741 (2021) 4
47. OpenAI: Gpt-4v(ision) system card (2023), <https://openai.com/research/gpt-4v-system-card> 5
48. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al.: Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* **35**, 27730–27744 (2022) 3
49. Peng, Z., Wang, W., Dong, L., Hao, Y., Huang, S., Ma, S., Wei, F.: Kosmos-2: Grounding multimodal large language models to the world. arXiv preprint arXiv:2306.14824 (2023) 4
50. Qin, Y., Zhou, E., Liu, Q., Yin, Z., Sheng, L., Zhang, R., Qiao, Y., Shao, J.: Mp5: A multi-modal open-ended embodied system in minecraft via active perception. arXiv preprint arXiv:2312.07472 (2023) 3
51. Qu, Y., Shen, X., He, X., Backes, M., Zannettou, S., Zhang, Y.: Unsafe diffusion: On the generation of unsafe images and hateful memes from text-to-image models. In: *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*. pp. 3403–3417 (2023) 4
52. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: *ICML (2021)* 9
53. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: *International conference on machine learning*. pp. 8748–8763. PMLR (2021) 8, 9, 6
54. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al.: Language models are unsupervised multitask learners. *OpenAI blog* **1**(8), 9 (2019) 3, 9, 12
55. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 **1**(2), 3 (2022) 4
56. Rando, J., Paleka, D., Lindner, D., Heim, L., Tramèr, F.: Red-teaming the stable diffusion safety filter. arXiv preprint arXiv:2210.04610 (2022) 4
57. Reed, S., Zolna, K., Parisotto, E., Colmenarejo, S.G., Novikov, A., Barth-Maron, G., Gimenez, M., Sulsky, Y., Kay, J., Springenberg, J.T., et al.: A generalist agent. arXiv preprint arXiv:2205.06175 (2022) 1, 2
58. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 10684–10695 (2022) 4
59. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*. pp. 234–241. Springer (2015) 8
60. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* **35**, 36479–36494 (2022) 4
61. Salimans, T., Ho, J.: Progressive distillation for fast sampling of diffusion models. arXiv preprint arXiv:2202.00512 (2022) 15

62. Shi, Z., Wang, Z., Fan, H., Yin, Z., Sheng, L., Qiao, Y., Shao, J.: Chef: A comprehensive evaluation framework for standardized assessment of multimodal large language models. arXiv preprint arXiv:2311.02692 (2023) [4](#)
63. Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al.: Mastering the game of go with deep neural networks and tree search. *nature* **529**(7587), 484–489 (2016) [2](#)
64. Sohn, K., Lee, H., Yan, X.: Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems* **28** (2015) [8](#), [9](#), [5](#), [7](#), [14](#)
65. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023) [3](#)
66. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al.: Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023) [3](#)
67. Tumanyan, N., Geyer, M., Bagon, S., Dekel, T.: Plug-and-play diffusion features for text-driven image-to-image translation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 1921–1930 (2023) [4](#)
68. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017) [6](#)
69. Vinyals, O., Babuschkin, I., Czarnecki, W.M., Mathieu, M., Dudzik, A., Chung, J., Choi, D.H., Powell, R., Ewalds, T., Georgiev, P., et al.: Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature* **575**(7782), 350–354 (2019) [2](#)
70. Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., Anandkumar, A.: Voyager: An open-ended embodied agent with large language models. arXiv preprint arXiv:2305.16291 (2023) [3](#)
71. Wang, Z., Cai, S., Liu, A., Jin, Y., Hou, J., Zhang, B., Lin, H., He, Z., Zheng, Z., Yang, Y., et al.: Jarvis-1: Open-world multi-task agents with memory-augmented multimodal language models. arXiv preprint arXiv:2311.05997 (2023) [3](#)
72. Wang, Z., Cai, S., Liu, A., Ma, X., Liang, Y.: Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. arXiv preprint arXiv:2302.01560 (2023) [3](#)
73. Wen, M., Lin, R., Wang, H., Yang, Y., Wen, Y., Mai, L., Wang, J., Zhang, H., Zhang, W.: Large sequence models for sequential decision-making: a survey. *Frontiers of Computer Science* **17**(6), 176349 (2023) [2](#)
74. Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y., et al.: mplug-owl: Modularization empowers large language models with multimodality. arXiv preprint arXiv:2304.14178 (2023) [4](#)
75. Yin, Z., Wang, J., Cao, J., Shi, Z., Liu, D., Li, M., Sheng, L., Bai, L., Huang, X., Wang, Z., et al.: Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. arXiv preprint arXiv:2306.06687 (2023) [4](#)
76. Zhang, K., Mo, L., Chen, W., Sun, H., Su, Y.: Magicbrush: A manually annotated dataset for instruction-guided image editing. arXiv preprint arXiv:2306.10012 (2023) [4](#)
77. Zhang, Z., Zhang, Y., Li, L., Gao, H., Wang, L., Lu, H., Zhao, F., Qiao, Y., Shao, J.: Psysafe: A comprehensive framework for psychological-based attack, defense, and evaluation of multi-agent system safety. arXiv preprint arXiv:2401.11880 (2024) [4](#)

78. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Conditional prompt learning for vision-language models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16816–16825 (2022) [2](#), [7](#)
79. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. *International Journal of Computer Vision* **130**(9), 2337–2348 (2022) [2](#), [7](#)
80. Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M.: Minigt-4: Enhancing vision-language understanding with advanced large language models. arXiv preprint arXiv:2304.10592 (2023) [4](#)
81. Zhu, X., Chen, Y., Tian, H., Tao, C., Su, W., Yang, C., Huang, G., Li, B., Lu, L., Wang, X., et al.: Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory. arXiv preprint arXiv:2305.17144 (2023) [3](#)

# MineDreamer: Learning to Follow Instructions via Chain-of-Imagination for Simulated-World Control

## Supplementary Material

The supplementary document is organized as follows:

- Sec. **A**: Environment Setting, like observation and action space.
- Sec. **B**: Dataset composition and collection.
- Sec. **C**: Implementation Details, like training details.
- Sec. **D**: Experiment Details, like baseline and evaluation details.
- Sec. **E**: More Ablation Studies about *MineDreamer*.
- Sec. **F**: More Visual Results about Imagination in *MineDreamer*.
- Sec. **G**: Demo videos about *MineDreamer*.

## A Minecraft Environment

Minecraft is a widely popular sandbox game that offers players the freedom to build and explore their worlds without limits, which also extends to AI agents as well. Within the game, AI agents encounter situations that closely mirror real-world challenges, requiring them to make decisions and solve endless tasks in an open-world setting. Consequently, Minecraft is an ideal platform for AI evaluation and stands as an exemplary benchmark for AI testing, due to its vast freedom and open nature. With the help of Minecraft, AI researchers can more easily simulate a wide variety of complex and dynamic environments and tasks, allowing them to conduct experiments that enhance the practical and applicable value of AI technologies.

We use MineRL [26] v1.0, which corresponds to Minecraft 1.16.5, as our simulation platform, ensuring an environment that is consistent with those used by VPT [3] and STEVE-1 [42]. In this version of MineRL [26], a significant advancement over its predecessor (*i.e.*, MineRL v0.4.4), lies in the simulation environment. The environment now enables AI agents to interact in a manner entirely consistent with human players, eschewing primitive actions or script-based APIs. This approach presents a more complex and challenging scenario for AI research. More specifically, AI agents experience the environment as humans do, solely through egocentric RGB images, devoid of any privileged in-game information. Additionally, their interactions with the environment are restricted to low-level keyboard and mouse actions. Consequently, AI agents trained in this version of MineRL [26] (*i.e.*, MineRL v1.0) resemble embodied agents capable of performing various tasks in an open-world environment, demonstrating a higher degree of generalization. Furthermore, the abundance of gaming videos available on the internet (*e.g.*, YouTube), provides AI researchers with the opportunity to

---

<https://github.com/minerllabs/minerl/releases/tag/v1.0>

harness these vast datasets for extensive pre-training, enabling the development of a foundation model in the sequential decision-making domain.

### A.1 Observation Space

Our observation space aligns with that of human players, comprising simply the raw pixels from Minecraft. This includes the hotbar, health indicators, player hands, equipped items, and the game environment itself. Specifically, the simulator produces RGB images with a resolution of 640x360. When the agent takes action within the environment, the simulator renders the player’s first-person perspective with a field of view of 70 degrees. If the agent opens the inventory, the simulator will render the GUI interface along with the mouse cursor.

Notably, we do not employ privileged information such as voxels and lidar information available in MineDojo [20], which could be provided to the agent. During actual inference, the PolicyNet of *MineDreamer* **only accepts the raw RGB pixels observations as input** that the agent can obtain from the environment and generates text-conditioned low-level action controls based on these observations, which are consistent with those used in VPT [3] and STEVE-1 [42].

### A.2 Action Space

As shown in Tab. 3, our action space encompasses a vast array of actions that are consistent with those of human players (*i.e.*, keyboard and mouse), including keypresses, mouse movements, and clicks. Excluding the “chat” action, which serves to initialize the agent with pre-defined conditions, more details can be found in Supp. A.3. Keyboard presses and mouse clicks are binary functional actions (*e.g.*, “Forward”, “Back”, “Left”, “Right” and *etc.*). Beyond these binary input options, our action space also has mouse cursor movements. While the GUI is closed (*i.e.*, activated by pressing “E” for the GUI inventory) and remains inactive, the mouse’s horizontal and vertical movements direct the agent’s yaw and pitch. Conversely, with GUI open, the same movements are re-purposed to navigate the cursor across the display.

It is noteworthy that we have not employed structured APIs such as “craft” and “smelt” as seen in MineDojo [20], which replace the need for precise mouse movements that are necessary for interacting with the inventory for certain tasks, effectively turning these operations into GUI functional binary actions. During actual inference, our *MineDreamer*’s PolicyNet **only outputs keyboard and mouse actions** to dictate the agent’s movements, aligning these actions with those utilized in VPT [3] and STEVE-1 [42].

### A.3 Environment Settings and Rules


In our experiments, the agent’s *initial position* at the start of the game, as well as the *seed* used to generate the environment, are completely random. This introduces an element of unpredictability and variety into the experimental setup, ensuring that the agent will encounter a wide range of scenarios and challenges.


**Table 3: Action Space utilized in the MineRL [26] simulator.** The action space primarily consists of 14 keyboard and mouse operations, with detailed descriptions sourced from the Minecraft wiki (<https://minecraft.fandom.com/wiki/Controls>).

Index	Action	Human Action	Description
1	Forward	key W	Move forward.
2	Back	key S	Move backward.
3	Left	key A	Strafe left.
4	Right	key D	Strafe right.
5	Inventory	key E	Open or close GUI inventory.
6	Drop	key Q	Drop a single item from the stack of items the player is currently holding.
7	Jump	key Space	Jump. When in the water, it keeps the player afloat.
8	Sneak	key left Shift	Move slowly in the current direction of movement.
9	Sprint	key left Ctrl	Move fast in the current direction of movement.
10	Attack	left Mouse Button	Destroy blocks (hold down); Attack entity (click once); Pick up the stack of items or place the stack of items in the GUI (click once)
11	Use	right mouse Button	Place the item being held or interact with the block that the player is currently looking at.
12	Hotbar.[1-9]	keys 1 - 9	Switch the appropriate hotbar cell.
13	Yaw	move Mouse X	Turning; aiming; camera movement.Ranging from -180 to +180.
14	Pitch	move Mouse Y	Turning; aiming; camera movement.Ranging from -180 to +180.

To better evaluate the agent’s ability to follow textual instructions for action prediction and its ability to rapidly adapt its behavior based on instructions, we have modified MineRL [26] to enable “chat” action operations. This allows for the swift initialization of the agent with predefined conditions through instructions. Specifically, for Programmatic Evaluation, we ensure that each experiment for all agents is conducted with the same seed and within the biome most conducive to completing the current instruction; across multiple experiments, different seeds are used. For Command-Switching Evaluation for Long-Horizon Tasks, all agents are placed in the same seed and biome optimal for the current instruction as well. In addition, the following rules are applied as aids:

- `/difficulty peaceful`: Set the difficulty of the environment to peaceful mode.
- `/gamerule doDaylightCycle false`: Set the environment to daytime forever.
- `/gamerule keep inventory true`: Set agent to not drop items upon death.


Specifically, for the task of “Obtain diamonds” , we add two additional rules on top of the aforementioned ones as assistance:

- `/effect give @a night_vision 99999 250 true`: Help the agent see more clearly in extremely dark environments (*e.g.*, at night or underground).
- `/give @p minecraft:diamond_pickaxe`: Provided the agent with a diamond pickaxe , enabling it to break almost all blocks and mine all ores within Minecraft.

For details regarding the most suitable biome used in the experiments, please check Supp. D.2 and Supp. D.3.

## B Dataset Details

### B.1 OpenAI Contractor Gameplay Dataset

All our raw data are based on the contractor dataset, which consists of offline trajectory data in Minecraft used for training VPT [3]. This dataset is created by hiring human contractors to play Minecraft and complete predetermined tasks, and it includes video (*i.e.*, image sequences), along with corresponding action sequences and metadata. OpenAI releases six subsets of contractor data: 6.x, 7.x, 8.x, 9.x, 10.x, and the MineRL BASALT 2022 dataset. Our Goal Drift Dataset ultimately selects three of these subsets as our raw data, including 8.x (house building from scratch), 10.x ( ) , and the FindCave dataset from the MineRL BASALT 2022 dataset. For each video, there is an associated metadata file that not only records the contractor’s actions for every frame but also documents events triggered by the contractor within the simulator; the specific events are detailed in Tab. 4.

**Table 4: The detailed event name and description in MineRL [26] simulator.** The simulator records the names of events that occur as well as related information, including quantities. We can use these events to collect a large amount of data for completing event-related instruction tasks with clarity.

Event Name	Description
mine_block	The moment the agent breaks a block, the type of block is recorded.
craft_item	The moment the agent crafts items, the type and number of items are recorded.
use_item	The moment the agent uses or places items, the type of item is recorded.
kill_entity	The moment the agent kills an entity, the type of entity is recorded.
break_item	The moment the tool of the agent is broken, the type of tool is recorded.
pick_up	The moment the agent picks up items, the type and number of items are recorded.

### B.2 Event Selection

In constructing our dataset, we opt to **select events directly from the MineRL simulator and supplement them with manually annotated events**. Specifically, to train the Imaginator within the constraints of limited resources, we focus on the following types of events: “mine\_block”, “craft\_item”, “use\_item”, “kill\_entity” and a manually defined event named “easy\_action”. Details of the specific items selected for each event can be found in Tab. 5. The simulator’s built-in events have a clearly defined completion time  $t^*$ , while manually annotated events are marked with a manually labeled completion time.

---

<https://github.com/openai/Video-Pre-Training>



**Table 5: Details of the specific items selected for each event.** We select four built-in events from the simulator, along with a manually defined event called “easy\_action”. The built-in events have a clearly defined completion moment, while the collection of the “easy\_action” event is manually annotated.

Event Name	mine_block	craft_item	use_item	kill_entity	easy_action
Detail items	Wooden Log	wooden planks	torch	sheep	Go explore
	Grass				Dig down
	Dirt				Look at the sky
	Grass Block				Go Swimming
	Sand				Stay underwater
	Snow				Build a tower
	Stone				Mine horizontally
	Coal Ore				
	Iron Ore				
	Redstone Ore				
	Diamond Ore				

### B.3 Dataset Collection

After obtaining the completion times  $t^*$  for all events, we employ `gpt-4-turbo` [47] to generate corresponding event-related instructions. Specifically, we provide `gpt-4-turbo` [47] with the event’s name, description, and detailed items, and prompt it to generate multiple distinct simple instructions. These instructions include specific actions, while others mention the items to be obtained upon completing the action. For instance, for “Grass” in the “mine\_block” event, `gpt-4-turbo` [47] would generate instructions like “break grass”, “break tall grass”, “gather seeds”, and “collect seeds”. After gathering instructions for all events, we apply the *Goal Drift Collection* method described in Sec. 3.3 of the main paper to conduct backward and forward drift on the completion times  $t^*$  of event-related instructions. For each pair (current observation, goal imagination), there are many instructions created by `gpt-4-turbo` [47] to describe that event. This process results in a substantial collection of triplets (current observation, goal imagination, instruction), which serve as training data for the Imaginator, forming what we call the Goal Drift Dataset. The final Goal Drift Dataset contains approximately 500,000 triplets (current observation, goal imagination, instruction), with about 400,000 of these triplets derived from events built into the simulator.

We follow the method used in STEVE-1 [42] for training the CVAE [36, 64] and collect a subset of approximately 10,000 quadruplets from the Goal Drift Dataset for the events we need to test subsequently. This subset consists of quadruplets where the current observation, goal imagination, and instruction are consistent as conditions with the Goal Drift Dataset. Additionally, there is a visual prompt embedding that serves as ground truth. This embedding is derived from a video composed of the goal imagination and the preceding 16 frames, processed through the MineCLIP [20] video encoder.

## C Implementation Details

### C.1 Imaginator

The training process of Imaginator is divided into three main stages. In the first stage, the MLLM is aligned with the CLIP [53] text encoder using the QFormer [39]. In the second stage, we apply InstructPix2Pix [6] to warm up the weights for the diffusion model in Minecraft. In the third stage, we optimize Imaginator in an end-to-end manner. To be specific, the weights of LLaVA [43] are frozen and LoRA [31] is added for efficient fine-tuning. For the diffusion model, we directly use the weights pre-trained in the second stage as the initial weights in Imaginator.

For the Large Language Model with visual input (e.g., LLaVA [43]), we choose LLaVA-1.1-7b [43] as the base model. During training, the weights of LLaVA are frozen and we add LoRA for efficient fine-tuning. We expand the original LLM vocabulary with 32 new tokens. The QFormer is composed of 6 transformer [68] layers and 77 learnable query tokens. We use the AdamW optimizer [44] in all three stages. In the initial stage of training, we configure the learning rate and weight decay parameters at  $2e-4$  and 0, respectively. The training targets for this stage encompass a dual-objective framework, comprising the Mean Squared Error (MSE) loss between the outputs of LLaVA [43] and the CLIP [53] text encoder, alongside the language model loss. Both losses are assigned equal weights of 1. The training setting in the second is the same as InstructPix2Pix [6]. In the final stage, the settings for the learning rate, weight decay, and warm-up ratio are adjusted to  $1e-5$ , 0, and 0.001, respectively. During this phase, the loss function is diffusion loss.

**Table 6: The Hyperparameters of Imaginator.**

Hyperparameter Name	Value
base_model	LLaVA [43]
input_image_size	$256 \times 256$
expand_vocabulary_num	32
transformer_layers_num	6
QFormer_learnable_query_num	77
optimizer	AdamW [44]
learning_rate_initial_stage	$2e-4$
weight_decay_initial_stage	0
learning_rate_final_stage	$1e-5$
weight_decay_final_stage	0
warm-up_ratio_final_stage	0.001
n_iterations_initial_stage	5000
n_iterations_final_stage	10000

## C.2 Prompt Generator

Our Prompt Generator is mainly a conditional variational autoencoder (CVAE) [36, 64] with a Gaussian prior and a Gaussian posterior similar to STEVE-1 [42]. Both the encoder and decoder of CVAE [36, 64] are parameterized as three-layer MLPs with 512 hidden units and layer normalization. It encodes the current observations, goal imaginations, and instructions then reconstructs a latent visual embedding, and uses a linear layer to project this embedding into the visual input space of our PolicyNet as the final visual prompt.

It is noteworthy that instead of using raw pixel images and natural language instructions directly as conditions to generate pixel-level videos depicting the execution of an instruction from the current observation to the imagined target, we opt to perform reconstruction within the visual space of MineCLIP [20], where MineCLIP [20] is a pre-trained CLIP model that employs a contrastive objective on pairs of Minecraft videos and associated transcripts from the web. Specifically, the process of generating prompts by the Prompt Generator mainly involves three steps. First, we stack the current observation and the goal imagination 16 times each to create two static 16-frame videos. These are then processed through MineCLIP [20]’s video encoder to obtain two visual embeddings. Concurrently, the instruction is encoded into a text embedding using MineCLIP [20]’s text encoder. This ensures that all embeddings are encoded within the MineCLIP [20] space. We then train a CVAE [36, 64] using the ELBO loss, which reconstructs a latent visual embedding from the previous three embeddings. This representation is a video embedding that describes the process within the MineCLIP [20] visual space. This representation is a video embedding that captures the process within the MineCLIP [20] visual space. The ground truth for this is mentioned in Supp. B.3 and is derived from the goal imagination and the preceding 16 frames, which have been processed through the MineCLIP [20] video encoder. In the end, we use a linear layer to project the latent visual embedding into the visual input space of the PolicyNet as the final visual prompt. For each event to be evaluated subsequently, we train a CVAE [36, 64] on the dataset, specifically for 150 epochs with early stopping on a small validation set. Notably, the parameters of the MineCLIP [20] within Prompt Generator remain unchanged, as do the parameters of the linear layer that maps MineCLIP [20]’s visual space to the visual input space of PolicyNet, whose parameters come from STEVE-1 [42]. The hyperparameters used during the training are listed in the following Tab. 7.

## D Experiment Details

In this section, we first detail the three baselines we select. We then separately present the Programmatic Evaluation details and the Command-Switching Evaluation for Long-Horizon Tasks details.

### D.1 Baseline Details

**Video Pretraining (VPT)** [3] is the first foundation model in the Minecraft domain, pre-trained on 70k hours of gameplay by Baker et al. [3]. Its archi-

**Table 7: The Hyperparameters of CVAE [36,64] within Prompt Generator.**

Hyperparameter Name	Value
architecture	MLP
visual_prompt_dim	512
text_dim	512
current_img_dim	512
goal_img_dim	512
hidden_layers	3
batch_size	256
learning_rate	1e-4
$\beta$	0.001
n_epochs	150

ecture primarily consists of two parts: ImpalaCNN and TransformerXL [17]. VPT [3] has three variants: VPT(fd), VPT(bc), and VPT(rl), representing the vanilla foundation model, the behavior cloning fine-tuned model, and the RL fine-tuned model, respectively. Specifically, they initially pre-trained on a large corpus of YouTube videos using a behavior cloning algorithm to obtain VPT(fd), which is capable of free exploration within the environment. This model gains a fundamental understanding of the environment and acquires some environmental knowledge. To enhance the agent’s capability in completing early-game tasks (*e.g.*, “Collect wood” 🌲 and “Craft wooden planks” 🪵), they collect an “Early-Game” video dataset and fine-tune the VPT(fd) to obtain VPT(bc). This model performs well in early-game tasks but struggles with long-horizon tasks, such as obtaining diamonds 💎. Building on VPT(bc), they employ online reinforcement learning with carefully designed rewards to fine-tune the model, enabling it to complete the task of obtaining diamonds 💎 from scratch, ultimately resulting in the creation of VPT(rl). Hence, it is noteworthy that **all three variants of VPT [3] are unable to follow instructions**; they must first be fine-tuned on downstream tasks before they can be completed. Despite their extensive environmental knowledge, this knowledge cannot be unlocked by instruction-following capabilities. In our experiments, we use VPT(rl) because it initially seeks out trees 🌲 and gathers wood 🪵, a critical step in the pathway to obtaining diamonds 💎. When set in the appropriate biome, VPT(rl) explores further 🏠 and collects more wood 🪵 compared to VPT(fd) and VPT(bc).

**STEVE-1** [42] is a Minecraft agent that can follow both textual and visual instructions, built upon MineCLIP [20] and VPT [3]. Drawing from the paradigms of instruction tuning in large language models and multimodal large language models, it successfully unlocks the instruction-following abilities of the foundation model (*i.e.*, VPT [3]) in the domain of decision-making. STEVE-1 [42] comes in two variants, STEVE-1(visual) and STEVE-1(text). The training process is divided into two steps. The first step involves training a policy conditioned on future video as visual instructions using the packed hindsight relabeling method. Specifically, they utilize the OpenAI Contractor Gameplay Dataset to fine-tune

VPT(rl) to follow visual instructions, resulting in STEVE-1(visual). The second step is to train a model that can map text instructions to visual instructions. Inspired by UnCLIP, they trained a Conditional Variational Autoencoder (CVAE) [36, 64] on a dataset of video-text pairs they collected, thus obtaining STEVE-1(text) which can follow text instructions. It is important to note that the visual or **textual instruction variants of STEVE-1 [42] do not consider the current observation and remain unchanged throughout the task, serving as an initial guide without adapting to environmental changes.**


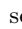

**Multi-Modal Memory** serves as a substitute for the Imaginator and Prompt Generator in the *MineDreamer* framework, essentially functioning by supplying PolicyNet with video prompts that best align with the current observations and textual instructions, similar to the approach of STEVE-1 (visual). We construct a multi-modal memory comprised of numerous video-text pairs. This memory is specifically built upon the triplets (current observation, goal imagination, instruction) from the Goal Drift Dataset. By tracing back 16 frames from the timestamp of the goal imagination, we create a 16-frame video segment, resulting in a revised triplet format: (current observation, goal imagination video, instruction). Each event, whether from the MineRL [26] environment or manually defined, contains 1,000 pairs. The retrieval process is as follows: First, we encode the current instruction and all instructions in the multi-modal memory using the OpenCLIP [52] text encoder to obtain embeddings. We then compare these embeddings using cosine similarity. Next, within the memory corresponding to the text instruction with the highest similarity, we find the match where the current observation and the memory’s observation, once encoded through the OpenCLIP [52] Image encoder, have the highest cosine similarity in their embeddings. Finally, the video from the final retrieval result is then encoded using the MineCLIP [20] video encoder, and the resulting visual embedding is used as the final visual prompt. Therefore, **Multi-Modal Memory leverages the current observation and also utilizes the Chain-of-Imagination (CoI) mechanism.**





## D.2 Programmatic Evaluation Details

In this part, we will elaborate on the selection of experimental tasks for Programmatic Evaluation, the methodology for calculating evaluation metrics, and the specific details of the experimental setup.






For the Programmatic Evaluation, we evaluate the agents on five **single-step** instruction tasks derived from the *early-game evaluation suite* proposed in Table 3 of the STEVE-1 [42] appendix. The purpose of this evaluation is to quantitatively measure an agent’s ability to follow instructions with minimal human intervention. Specifically, we calculate the programmatic evaluation metrics by monitoring the state of the MineRL [26] environment during each evaluation episode. Consistent with VPT [3] and STEVE-1 [42], we compute multiple pro-

grammatic metrics, including travel distance, dig depth, and early-game item collection. The calculation is as follows:

1. **Travel Distance (Blocks)**: The agent’s maximum horizontal displacement, in the X-Z plane, is measured from the initial spawn point.
2. **Dig Depth (Blocks)**: The agent’s maximum vertical (Y-axis) displacement is measured from its initial spawn point.
3. **Early-Game Inventory Counts**: The maximum number of log , seed , and dirt  items seen in the agent’s inventory during the episode.

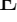
We test all agents on these five **single-step** instruction tasks, with each task running 10 episodes of 3000 timesteps (*i.e.*, 2.5 minutes of gameplay). Each episode used a unique environmental seed, yet all agents were tested under the same seed for consistency. It is important to note a key difference in our experimental setup compared to STEVE-1 [42]: **for each task, we initialize the agents in the biome most conducive to task completion** to enhance the reliability of our evaluation metrics. For instance, in the “Chop a tree”  task, all agents are spawned in a forest biome, rather than a plain, to avoid the added randomness of searching for trees  before chopping them. Due to a limited computational budget, we do not generate goal imaginations for every frame within an episode. In MineRL [26], an agent can perform only one mouse or keyboard action per frame, and for tasks such as breaking a block of dirt , it requires approximately 25 frames of consistently holding down the left mouse button. Therefore, **we decide to imagine a goal imagination and translate it to a visual prompt every 25 frames** ultimately, which then guides the action generation for the following 25 frames (*i.e.*, the visual prompt  $p_t$  will not change for the next 25 frames). This interval is chosen because, aside from the “Chop a tree”  task, the other four tasks can be achieved within 25 frames (*i.e.*, just over 1 second of gameplay), thereby necessitating a new round of imagination to guide subsequent actions. The detailed settings for the Programmatic Evaluation can be found in Tab. 8.











**Table 8: The detailed settings for the Programmatic Evaluation.**

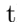



Id	Text Instruction	Biome	Time Limit	Imagination Interval	Metric
1	go explore 	Plains	3000 Frames	25 Frames	Travel Distance (Blocks)
2	collect seeds 	Plains			Seeds Collected
3	chop a tree 	Forest			Wooden Logs Collected
4	collect dirt 	Plains			Dirt Collected
5	dig down 	Plains			Dig Depth (Blocks)

### D.3 Command-Switching Evaluation Details

In this part, we will also detail the selection of experimental tasks for Command-Switching Evaluation for Long-Horizon Tasks, the calculation methods for evaluation metrics, and the specific details of the experimental setup.

The Command-Switching Evaluation for Long-Horizon Tasks comprises three **multi-step** instructions tasks sourced from the *early-game evaluation suite* of STEVE-1 [42], except the “Obtain diamonds”  task which originates from GROOT [8], designed to steadily follow video instructions. These tasks aim to evaluate an agent’s ability to swiftly adapt to new instructions following an instruction switch, a critical capability for a downstream controller operating under an LLM-based high-level planner. We employ success rate as the performance metric, also by monitoring the MineRL [26] environment state throughout each evaluation episode. The criteria for determining success across the three different tasks are as follows:

1. **collect wood  and then craft planks **: Success is defined as successfully crafting at least one wooden log  into four wooden planks  within the given time frame.
2. **gather dirt  and then build a tower **: Success is defined as successfully building a tower  with a height of at least 7 blocks within the given time frame.
3. **dig down  and then mine horizontally **: Success is obtaining at least one diamond  within the given time frame.

For these three **multi-step** instructions tasks, we run 50 episodes of testing per task. The time limit for the first two tasks is set at 3000 frames (*i.e.*, 2.5 minutes of gameplay), consistent with STEVE-1 [42], while the final task has an episode time limit of 12,000 frames (*i.e.*, 10 minutes of gameplay), aligning with what is mentioned in the main paper of GROOT [8]. Each episode utilizes a unique environmental seed to ensure variability; however, all agents are tested with the same seed for consistency across episodes. It is important to note that our experimental setup differs from that of STEVE-1 [42] in that **we initialize the agents in the biome most conducive to task completion** for each task. Specifically, as mentioned in Supp. A.3, we utilize the “chat” action to initialize the agent. For the “Obtain diamonds”  task, we equip the agent with night vision and a diamond pickaxe , which is consistent with the description provided in the main paper of GROOT [8]. **Considering that STEVE-1 [42] may not be explicitly trained on the “mine horizontally”  instruction, we augment STEVE-1 [42]’s prior original training data with the corresponding text-video pairs from the Goal Drift Dataset and retrain the prior.** This ensures that the updated prior can map the textual instruction “mine horizontally”  to the associated visual instructions. The detailed settings for the Command-Switching Evaluation for Long-Horizon Tasks experiment can be found in Tab. 9.

## E More Ablation Studies

In this section, we introduce additional ablation studies to explore various contributors to performance. This includes the use of Classifier-Free Guidance [30] during inference, the selection of Drift Lengths from the Goal Drift Dataset, and the generation strategies for Visual Prompts. We employ the same experimental

**Table 9: The detailed settings for the Command-Switching Evaluation.**

Id	Text Instruction	Biome	Switch Condition	Time Limit	Imagination Interval
1	chop a tree 🌳 craft wooden planks 🪵	Forest	Reach 1500 Frames	3000 Frames	25 Frames
2	collect dirt 🪨 build a tower 🏰	Plains	Reach 2000 Frames	3000 Frames	25 Frames
3	dig down 🛖 mine horizontally 🛖	Plains	Reach 13th floors	12000 Frames	25 Frames

settings as our Programmatic Evaluation, compare the performance of different ablations, and plot the results, showing both the mean and 95% confidence intervals of the programmatic metrics.

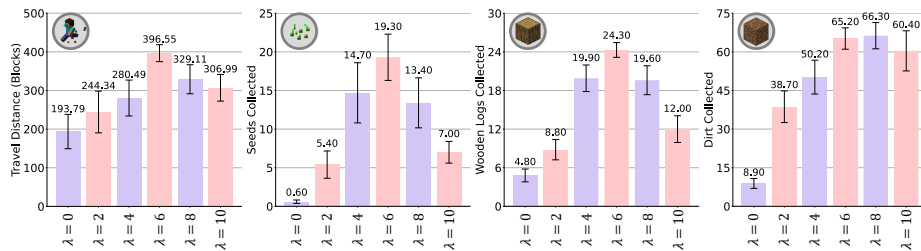
### E.1 Classifier-Free Guidance During Inference

Given that VPT [3] is a foundation model obtained through behavior cloning from extensive video demonstrations without instruction guidance during training, this may lead to a smoother behavior distribution learned by VPT [3]. Consequently, even after fine-tuning VPT [3] for instruction-following abilities, when provided with direct instruction as a condition, it tends to act based on its previously learned knowledge from behavior cloning. It fails to steadily follow the instructions given, similar to the observation in Appendix I of Baker et al. [3]. We believe that this bias arises inherently from the training process of VPT [3]. Inspired by STEVE-1 [42], we employ classifier-free guidance [30] to mitigate this bias as much as possible in the action logits space before sampling the action. Specifically, for each inference, we perform two computations of logits through PolicyNet: one with visual prompt guidance and the other without. At each timestep, **we subtract a certain proportion of the action logits from the unconditioned PolicyNet from those predicted by the visual prompt-conditioned PolicyNet**. The equation for computing logits is directly borrowed from STEVE-1 [42].

$$f_t \leftarrow \mathcal{V}(\mathcal{O}_t), \quad o_t \leftarrow f_t + p_t, \quad \text{logits} \leftarrow (1 + \lambda) \underbrace{\mathcal{T}_\theta(o_{t-T}, \dots, o_t)}_{\text{conditional logits}} - \lambda \underbrace{\mathcal{T}_\theta(f_{t-T}, \dots, f_t)}_{\text{unconditional logits}} \quad (1)$$

where  $\mathcal{V}$  is the VisualEncoder and  $\mathcal{T}$  is the TransformerXL [17],  $\mathcal{O}_t$  is the current observation,  $p_t$  is the visual prompt,  $\lambda$  is the trade-off parameter between the visual prompt conditioned logits and unconditioned logits. By setting a suitable value for  $\lambda$ , we can encourage PolicyNet to follow the instructions in action generation more steadily. Fig. 9 illustrates how choosing different values affects the agent’s performance in Programmatic Evaluation. When the value of  $\lambda$  is less than 6, performance improves with an increase in  $\lambda$ , indicating that classifier-free guidance [30] can significantly reduce the bias introduced by prior behavior. The agent performs optimally when  $\lambda$  is 6 to 8; beyond this range, the performance begins to decline. This decrease is due to excessive guidance disrupting the agent’s original understanding and knowledge of the environment, impeding its ability to act normally. Ultimately, we opt for a value of  $\lambda$  equal to 6.





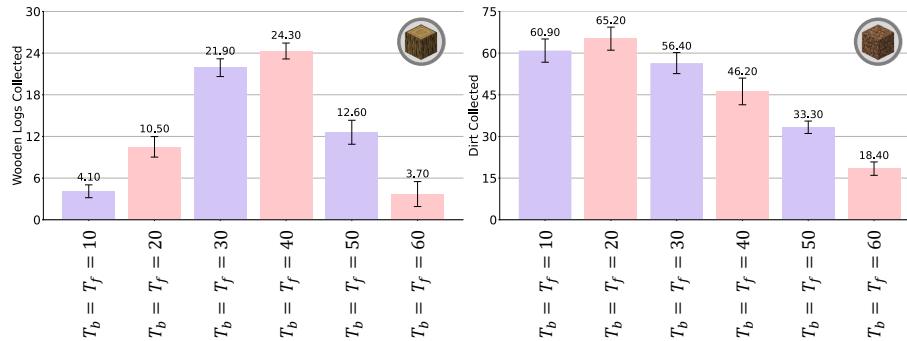
**Fig. 9: The impact of different values of condition scale  $\lambda$  on the performance of the agent by using classifier-free guidance.** Selecting the optimal parameter  $\lambda$  to balance between visual prompt-conditioned and unconditioned settings can significantly enhance agent performance, consistently improving its ability to follow instructions. By using the best  $\lambda$  ( $\lambda = 6$ ), *MineDreamer* when significantly outperforms *MineDreamer* when  $\lambda = 0$  (no guidance), collecting  $32\times$  more seeds 🌱,  $5\times$  more wood 🪵,  $7.3\times$  more dirt 🟫, travelling  $2\times$  further 🚶.

After utilising classifier-free guidance [30] during inference, *MineDreamer* when  $\lambda = 6$  significantly outperforms *MineDreamer* when  $\lambda = 0$  (no guidance), collecting  $32\times$  more seeds 🌱,  $5\times$  more wood 🪵,  $7.3\times$  more dirt 🟫, travelling  $2\times$  further 🚶. Therefore, selecting an appropriate value for parameter  $\lambda$  to balance the trade-off between visual prompt-conditioned and unconditioned logits can significantly enhance the agent’s performance and steadily improve its ability to follow instructions in action generation. Although this technique trick is effective during inference, it still needs to find the best hyperparameter in practice. In the future, eliminating biased behaviors directly from the fine-tuning process training would be meaningful.

## E.2 Selection of Drift Lengths

During Goal Drift Dataset collection, we utilize fixed values for  $T_b$  and  $T_f$  to address the challenges of “Goal Illusion” and “Imagination Stagnation” by performing backward and forward drifts around the event occurrence moment  $t^*$ . The specific algorithmic procedure is detailed in Sec. 3.3 of the main paper.

It is noteworthy that we observe an inconsistency in the optimal Drift Length for each event. As illustrated in Fig. 10, the optimal drift length for “wood” 🪵 is approximately 40, while for “dirt” 🟫, it is around 20. We find that the best drift length correlates with the amount of time required to complete the instruction task once. Also as shown in Fig. 10, selecting appropriate values for  $T_b$  and  $T_f$  can effectively address the “Goal Illusion” and “Imagination Stagnation” challenges mentioned in Sec. 3.3 of the main paper, enhancing the agent’s ability to steadily follow instructions in action generation. Although this method is effective, it does require the cumbersome task of selecting the right length for each event. We believe that mitigating or eliminating the interference caused by varying drift lengths during training or fine-tuning in the future will be meaningful.



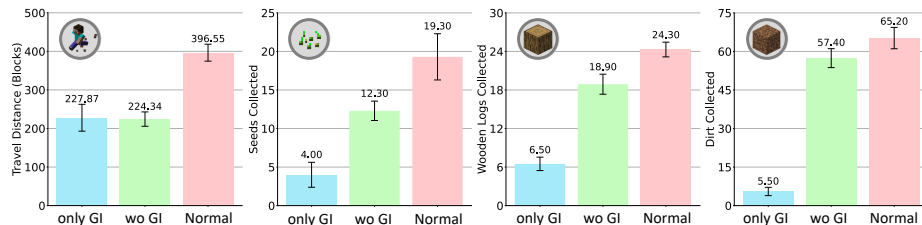
**Fig. 10: The influence of different goal drift lengths on agent performance.** For each event, there is an optimal goal drift length that is correlated with the duration of an instruction task completion one time. Employing the appropriate goal drift length can address the dual challenges of “Goal Illusion” and “Imagination Stagnation”, thereby enhancing the agent’s ability to steadily follow instructions in action generation.

### E.3 Generation Strategies for Visual Prompts

When the agent acts in the simulator, the Imaginator first creates a goal imagination of the next stage to complete the given instruction based on the current observation and instruction. Then, the Prompt Generator creates a visual prompt from this goal imagination, integrating the current observation and instruction. This part investigates strategies for generating visual prompts. We consider two variants:

1. Unlike current methods, we can synthesize the imagination into a 16-frame video by simply stacking it 16 times, and we encode this video with the MineCLIP [20] video encoder to project it into the MineCLIP [20] space and align it with the PolicyNet using a linear layer. More specifically, we bypass the reconstruction step mentioned in Supp. C.2 and directly transform the only goal imagination into the required visual prompt for the MineCLIP [20] visual space.
2. In contrast to current methods, we eliminate the Imaginator and retrain a Prompt Generator to directly reconstruct visual prompts from current observations and instructions. Specifically, we retrain a CVAE [36, 64] without using the goal of Imagination as a guiding condition for visual prompt generation.

From Fig. 11, it is evident that our approach enables the agent to follow instructions more steadily, as our visual prompt provides a more precise demonstration of the desired behaviour customized to the current environment. One drawback of using goal imagination stacked into a 16-frame video as a visual prompt is that the depicted behavior resembles a static state. This can confuse PolicyNet, making it unclear whether to remain stationary or to achieve the state represented in the video. A limitation of reconstructing the current observation and instruction into a visual prompt is that the CVAE [36, 64]’s ability to model



**Fig. 11: The impact of different visual prompt generation strategies on the performance.** “only GI” refers to bypassing the CVAE reconstruction phase in Prompt Generator and directly stacking the goal imagination into a static 16-frame video as the visual prompt. “wo GI” indicates that the CVAE reconstructs the visual prompt without using goal imagination as a condition, thus skipping the imagination phase of the Imaginator.

future spatiotemporal aspects is subpar. Without relying on goal imagination, it struggles to accurately reconstruct the demonstration of the desired behaviour. This occasionally results in misleading the agent, preventing it from steadily following instructions during action generation.

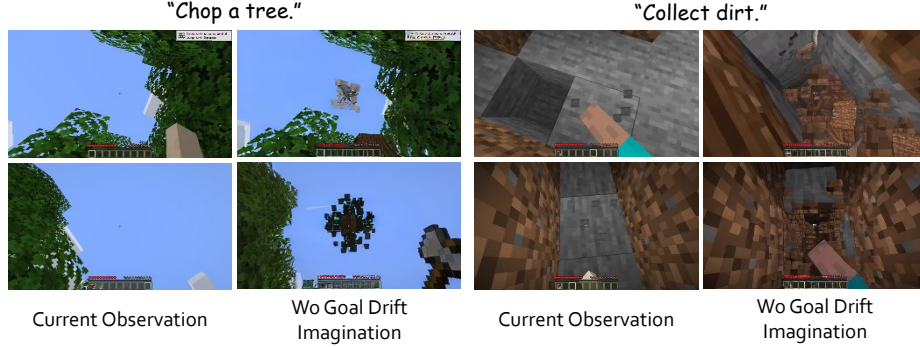
## F More Visual Results

Imagination visual results are all hosted on our [project webpage](#).

### F.1 Imagination Visual Results without Goal Drift

To evaluate the efficacy of the Goal Drift data collection method, we carry out experiments comparing various data collection approaches. Fig. 12 illustrates the imagination generated by the Imaginator trained on data collected without any goal drift. Due to the absence of backward drift, all imaginations generated by the Imaginator correspond to the moment when the event-related instructions are completed. Consequently, this leads to the phenomenon of “Goal Illusion”, where the Imaginator edits the current observation to depict the completed instruction. For the instruction “Chop a tree” 🌳, when the agent faces the sky, the Imaginator may unrealistically insert a broken wooden log 🪵 into the sky. For the instruction “Collect dirt” 🪨, even though the agent is pointing at a stone 🪨, the Imaginator still imagines dirt and shatters it, resulting in the agent eventually attempting to break the stone 🪨. Fig. 13 shows the imaginations generated by the Imaginator trained on data collected without forward drift. Because there is no forward drift, all imaginations generated by the Imaginator represent moments before the completion of event-related instructions. This results in the phenomenon of “Imagination Stagnation”, where the Imaginator fails to conceive repeated task completion. For the instruction “Chop a tree” 🌳, after cutting the uppermost wood 🪵 by looking up, the agent will not look down for more trees 🌳, which impedes continuous task performance. In contrast, an

Imaginator trained with data collected including forward drift is able to understand that the agent should now look down to find other trees 🌳 to continue the task.



**Fig. 12: Imagination Visual Results without Goal Drift.** Due to the absence of goal drift, the imaginations generated by the Imaginator are all related to the moment of event-related instruction completion, leading to the phenomenon known as “Goal Illusion”, where the Imaginator edits the current observation to represent the executed instruction. In the figure depicted, the agent inserts broken wooden blocks 🪵 into the sky and, facing a stone 🪨, imagines itself breaking dirt 🟫.



## F.2 Imagination Visual Results on Evaluation Set

We compare Imaginator with the existing state-of-the-art instruction-based image editing model, namely InstructPix2Pix [6]. Given this model has been trained on specific datasets, its performance would inevitably be suboptimal if directly applied to the Minecraft domain. To facilitate a fair comparison, we fine-tune InstructPix2Pix [6] using the same training set employed by the Imaginator and evaluate the performance of the fine-tuned models in addressing tasks in Minecraft. Fig 14 shows qualitative results in the evaluation set, our methodology exhibits enhanced abilities in Goal Imagination Generation within intricate scenarios.

## F.3 Imagination Visual Results During Agent Solving Tasks

We visualize the agent’s imagination during task execution alongside the next observation in Fig. 15 and Fig. 16 to evaluate the Imaginator’s generalization capability in open scenarios. It is observed that the Imaginator is capable of generating high-quality visualizations that closely align with the current scene in an open environment, thereby guiding the subsequent PolicyNet to autoregressively predict the next action steadily.



**Fig. 13: Imagination Visual Results without Forward Drift.** Due to the lack of forward drift, the imaginations produced by the Imaginator are all from moments prior to the completion of event-related instructions, resulting in a phenomenon called “Imagination Stagnation”. This means the Imaginator fails to anticipate the outcomes of repeated tasks. For example, in the figure provided, after the agent cuts the uppermost wood  by looking up, it will not look down for more trees  to continue the task.

#### F.4 User Studies

To further evaluate *MineDreamer*’s efficacy, we conduct a user study. Specifically, we randomly select 15 images from the evaluation set, representing a wide range of tasks and scenarios within Minecraft. For each image, we generate results using both InstructPix2Pix [6] and *MineDreamer*, then randomly shuffle the order of these results. As noted in Sec. 4.3 of the main paper, InstructPix2Pix [6] is fine-tuned on the same dataset as *MineDreamer*. This process yield 15 sets of images in a shuffled sequence. Participants are asked to independently identify the two superior images for each set: the first being the one that best matches the given instructions (named **Instruct-Alignment**), and the second being the image that most closely mirrors real-world appearances, including perspective and physical laws (named **Image Quality**). A total of 25 individuals participate in the study. The findings, illustrated in Fig. 17, reveal that over 69.40% of participants find *MineDreamer*’s outputs to be more aligned with the instructions, and more than 70.31% favor the results produced by *MineDreamer* for their realism. These outcomes further underscore *MineDreamer*’s instruction following ability and generalization ability.

## G Demo Videos

Demo videos are all hosted on our [project webpage](#).

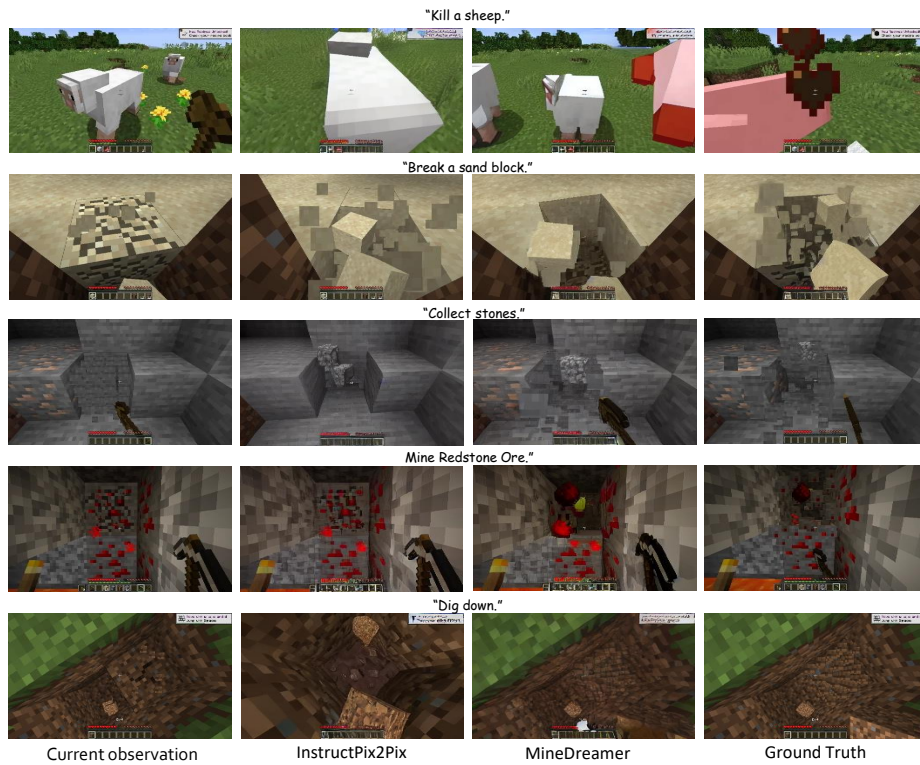


Fig. 14: Imagination visual results on Goal Drift Evaluation Set.



**Fig. 15:** Imagination visual results during agent solving tasks.

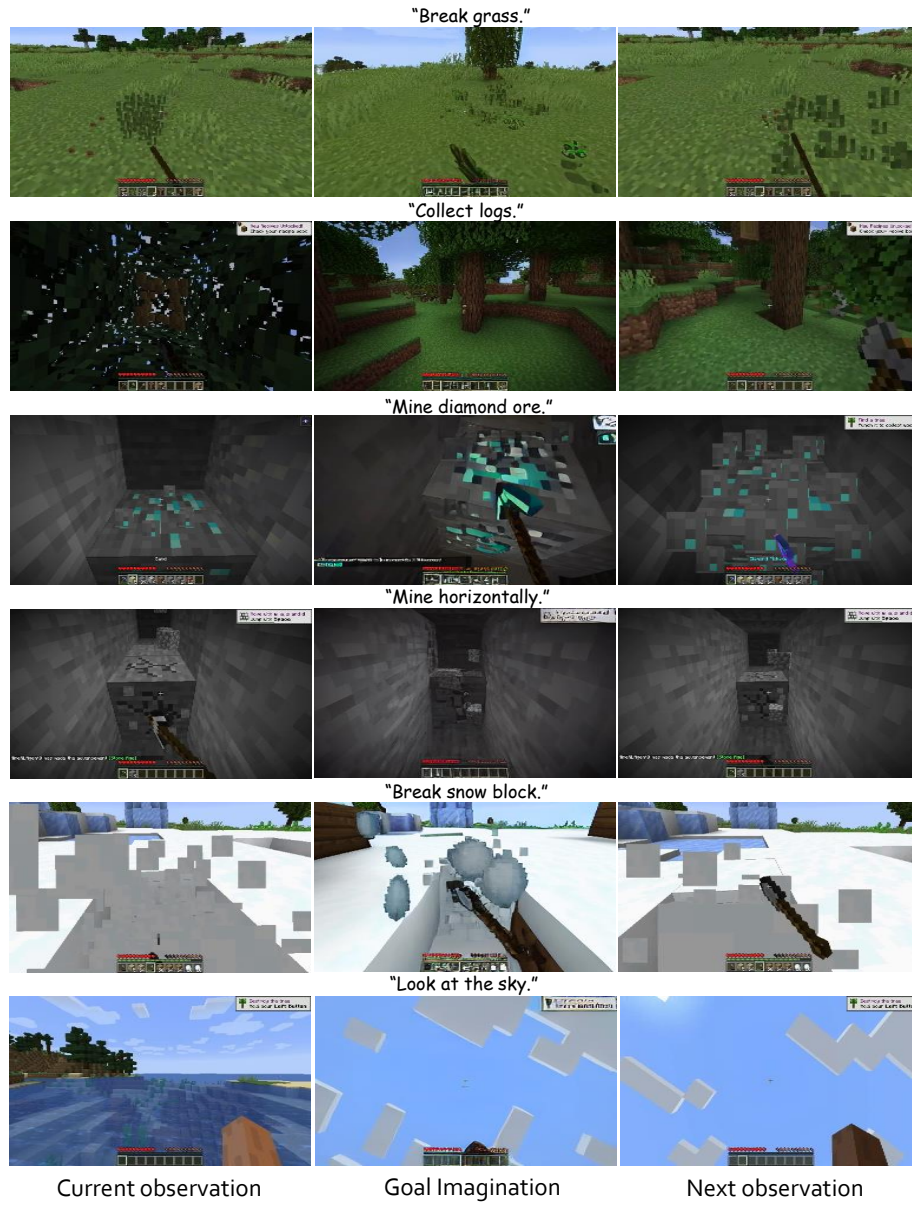
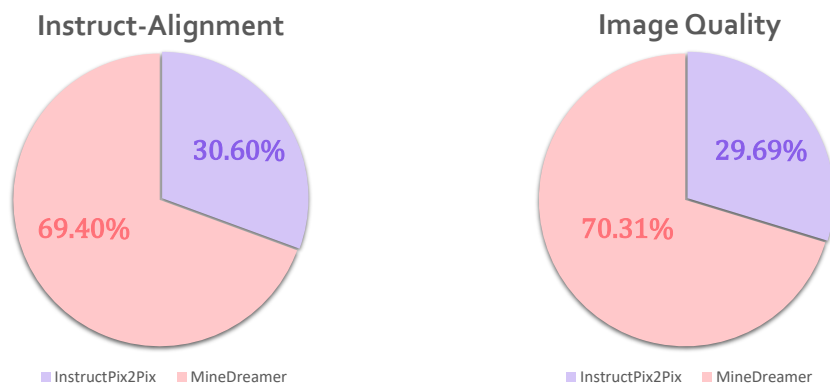


Fig. 16: Imagination visual results during agent solving tasks.





**Fig. 17: The results of user studies**, comparing the results generated by Instruct-Pix2Pix and *MineDreamer*. Based on the results from both the Instruction Alignment and Image Quality perspectives, *MineDreamer* demonstrates superior effectiveness.

### G.1 Programmatic Evaluation

We demonstrate videos of the four tasks from the Programmatic Evaluation on the aforementioned anonymous project webpage. Of course, you can also view the demo videos for the respective tasks by directly accessing the video URLs.

- “Go explore” 🧑: <https://youtu.be/UdG0ckoGRCY>
- “Collect seeds” 🌱: [https://youtu.be/TFchu\\_YBiuI](https://youtu.be/TFchu_YBiuI)
- “Chop a tree” 🌳: [https://youtu.be/Sx\\_NKjq5DTA](https://youtu.be/Sx_NKjq5DTA)
- “Collect dirt” 🪨: <https://youtu.be/7TOR0SOFaB8>

### G.2 Command-Switching Evaluation

We demonstrate videos of the three tasks from the Command-Switching Evaluation on the anonymous project webpage mentioned above. Of course, you can also view the demo videos for the respective tasks by accessing the video URLs.

- “Chop a tree 🌳 to Craft planks 🪵”: [https://youtu.be/YtY2M\\_Hi70E](https://youtu.be/YtY2M_Hi70E)
- “Gather dirt 🪨 to Build a tower 🏰”: <https://youtu.be/Zy2t2RpeNtQ>
- “Obtain Diamond” 💎: <https://youtu.be/hThbWh0q5EE>