

# Cultivating Game Sense for Yourself: Making VLMs Gaming Experts

Anonymous ACL submission

## Abstract

Developing agents capable of fluid gameplay in first/third-person games without API access remains a critical challenge in Artificial General Intelligence (AGI). Recent efforts leverage Vision Language Models (VLMs) as direct controllers, frequently pausing the game to analyze screens and plan action through language reasoning. However, this inefficient paradigm fundamentally restricts agents to basic and non-fluent interactions: relying on isolated VLM reasoning for each action makes it impossible to handle tasks requiring high reactivity (e.g., FPS shooting) or dynamic adaptability (e.g., ACT combat). To handle this, we propose a paradigm shift in gameplay agent design: instead of directly controlling gameplay, VLM develops specialized execution modules tailored for tasks like shooting and combat. These modules handle real-time game interactions, elevating VLM to a high-level developer. Building upon this paradigm, we introduce GameSense, a gameplay agent framework where VLM develops task-specific game sense modules by observing task execution and leveraging vision tools and neural network training pipelines. These modules encapsulate action-feedback logic, ranging from direct action rules to neural network-based decisions. Experiments demonstrate that our framework is the first to achieve fluent gameplay in diverse genres, including ACT, FPS, and Flappy Bird, setting a new benchmark for game-playing agents.

## 1 Introduction

Developing agents that fluidly play first/third-person games without API access remains a critical challenge in AGI, where complexity mirrors real-world embodied tasks (Lu et al., 2024; Wang et al., 2024). Agents must navigate diverse tasks, ranging from combat encounters to environmental navigation, while executing precise real-time actions (Hu et al., 2024). Traditional reinforcement learning (RL) approaches struggle to handle such a broad

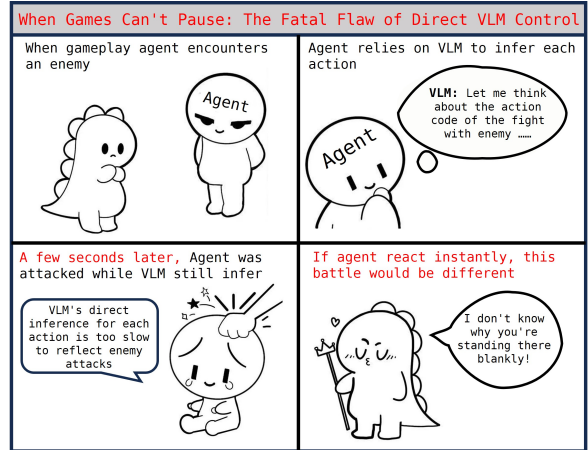


Figure 1: The ‘thinking time’ of direct VLM control becomes a critical vulnerability in real-time games, highlighting the need for a paradigm shift on VLM use: from direct controller to execution module developer

spectrum of demands due to their limited task generalization (de Woillemont et al., 2022; Justesen et al., 2019). Recently, the emergence of Vision Language Models (VLMs) has opened new possibilities in this domain. With their strengths in visual understanding and decision-making, VLMs interact with games purely through visual understanding of game screens. This ability offers a promising direction for developing non-API-dependent gameplay agents (Tan et al., 2024; Liu et al., 2024a; Wang et al., 2023b).

Recent VLM-based approaches leverage VLMs as direct game controllers through a pause-and-plan paradigm (Tan et al., 2024; Chen et al., 2024): the agent periodically pauses gameplay, using VLM and vision tools (e.g., OCR, segmentation) together to analyze game screens, plan actions and then directly output key-mouse command to control game. However, this paradigm suffers from fundamental limitations: (1) it heavily depends on the game’s support to pause at any moment, which disrupts the gameplay flow and limits its applicability to

a narrow range of games that support such interruptions; (2) Requiring VLM reasoning for every action makes it unsuitable for tasks demanding high reactivity (e.g., FPS shooting); (3) VLM outputs simple key-mouse control commands without real-time interactive logic for game environments, making it hard solve tasks demanding dynamic adaptation (e.g., action game combat). These limitations reflect a **fundamental mismatch**: VLMs excel at time-consuming deliberate reasoning (scene understanding and planning) but struggle with rapid, continuous game interactions requiring millisecond-level responses (shown in Figure 1).

We observe that most human game actions rarely rely on deliberate reasoning, but rather flows from quick-fire game sense - a set of trained reflexes and patterns developed through practice. This observation suggests a fundamental paradigm shift: Unlike using VLMs to directly control every game actions, we should elevate them to develop task-specific execution modules that can handle real-time interactions autonomously. These specialized modules, developed by VLM, solve specific tasks requiring rapid reactions or frequent environmental interactions. This paradigm shift bridges the VLM’s reasoning with real-time gameplay demands, enabling more versatile game agents.

Based on this new paradigm, we present GameSense, a framework that empowers VLMs to develop and optimize task-specific execution modules, termed Game Sense Modules (GSMs). GameSense equips VLMs with essential tools, including vision tools and neural network training pipelines, to create GSMs tailored for diverse gameplay tasks. These modules can range from simple action-feedback loops (e.g., combat patterns based on HP bar monitoring) to complex, learned behaviors (e.g., boss fight strategies optimized through RL). These modules are seamlessly integrated into the gameplay loop: when the agent identifies a specific task, it activates the corresponding module and refines it based on execution feedback. By shifting VLMs’ role from direct controller to the developer of GSMs, GameSense achieves efficient execution and promotes continuous improvement in gameplay performance.

Experiments demonstrate that GameSense is the first agent to achieve fluent gameplay in diverse game genres. In ACT/FPS games, our framework achieves the highest success rates in combat tasks, while achieving the highest exploration scores without gameplay pausing. In contrast, existing VLM-

based methods either fail to complete such tasks or rely heavily on frequent gameplay pausing, disrupting the flow of real-time interactions. In the reflex-intensive game Flappy Bird where pausing is not supported, existing VLM-based methods fail at basic control, and GameSense develops precise control modules through iterative refinement. GameSense exhibits significantly improved real-time performance and adaptation capabilities, setting a new benchmark for game-playing agents. The contributions of this paper are as follows:

- We identify limitations of existing VLM-based game-playing approaches, particularly their inability to handle real-time, high-reactivity tasks.
- We propose a novel paradigm that uses VLMs to develop task-specific execution modules for autonomous real-time interactions.
- We introduce **GameSense**, a framework that enables VLMs to create and refine **Game Sense Modules (GSMs)**.
- Our experiments demonstrate that GameSense outperforms existing methods and is the first to master reflex-intensive games.

## 2 Related Work

### 2.1 Environment for Video Gameplay and RL-based Agents

Researchers have made significant strides in various video game environments, including classic games like Atari games(Bellemare et al., 2013), Minecraft(Fan et al., 2022; Guss et al., 2019), StarCraft II(Ellis et al., 2023). However, these environments rely heavily on open-source code or official APIs, requiring substantial human effort for implementation. This dependency restricts AI accessibility to general games. Recent RL-based approaches have attempted to overcome API dependencies by directly processing game visuals and simulating keyboard-mouse inputs, including DQN-play-sekiro(analoganddigital., 2021). However, these RL methods typically work for specific tasks and exhibit poor generalization, requiring re-training for new scenarios. The challenge of developing agents capable of generalizing across diverse gaming environments without API access remains largely unsolved. This limitation motivates our research toward a more adaptable solution using only visual inputs and key-mouse controls.

## 2.2 LLM/VLM-Driven Gameplay Agent

Current LLM/VLM-driven gameplay agents follow two main approaches. The first relies on game APIs for state observation and control, as seen in Minecraft(Wang et al., 2023a; Liu et al., 2024a) and Starcraft II agents(Ma et al., 2023). While effective, this API dependency limits their application to closed-source commercial games. The second approach uses VLMs to directly process screen information and generate keyboard-mouse controls, as demonstrated by Cradle(Tan et al., 2024). Though eliminating API requirements, this method’s frame-by-frame analysis and decision-making process introduce significant latency. This makes such agents unsuitable for tasks requiring quick reactions or dynamic adaptation. While recent works like SIMA(Raad et al., 2024) and VARP(Chen et al., 2024) attempt to improve performance through behavior cloning, they require extensive human gameplay data for training. The challenge of achieving real-time and adaptive gameplay in VLM-driven agents remains unsolved, motivating our research toward a new paradigm.

## 3 Method

### 3.1 Problem Formulation and Motivation

This work aims to develop a real-time gameplay agent that operates **without** relying on game APIs or pausing the game for action reasoning. The agent solely depends on real-time game screens and outputs key-mouse control commands to interact with the game. This approach aims to create a truly **in-game** agent, mirroring how human players experience and interact with the game environment.

Existing gameplay agents rely on the "pause and plan for each action" paradigm, which exhibits limitations in fast-paced and dynamic game scenarios. In contrast, most human gameplay actions do not stem from deliberate reasoning over each move but from game sense—an intuitive ability to react swiftly based on experience. Motivated by this observation, we propose an agent system capable of developing its form of "game sense," enabling more natural and efficient interaction in gameplay.

### 3.2 Overview of GameSense

GameSense introduces a **paradigm shift** by elevating the VLM from direct controller to developer of task-specific execution modules, termed Game Sense Modules (GSMs). The agent integrates a High-Level VLM Agent and GSMs: the

High-Level VLM Agent is responsible for real-time game screen analysis, historical reflection, and task and action planning. The GSMs, independently developed by the VLM itself, handle tasks requiring rapid response (e.g., combat, shooting, rapid clicks). As shown in Figure 2, the agent operates in a continuous loop: it analyzes real-time game screens, reflects on history, and plans tasks and actions. Depending on the action requirements, the agent either directly generates key-mouse control codes (**VLM-executed actions**) for straightforward actions or invokes GSMs (**GSM actions**) for high-speed processing. This process ensures efficient and natural interaction with the game, mirroring human-like gameplay.

### 3.3 High-Level VLM Agent

The High-Level VLM Agent serves as the brain of the system, responsible for understanding the game environment, reflecting on past experiences, and planning future tasks and actions (both VLM-executed and GSM actions). This module is structured into several core components:

**Game Environment Analysis:** This module leverages VLM’s visual understanding capabilities to generate a textual description of the current game screens. It identifies key elements such as the presence of enemies, bosses, interactable objects, potential threats, and the player character’s status. This textual description is then used for historical reflection and task planning.

**Historical Data Reflection:** This module performs three parallel types of reflection to learn from the past: (1) Previous Task Reflection: evaluate the success of the previous task and suggesting optimizations; (2) Historical Task Summary: summarize the last 10 task executions to extract long-term patterns; and (3) Action Design Reflection: assess **VLM-executed actions**’ effectiveness and generating refinements. This mechanism ensures continuous self-assessment and refinement.

**Memory:** This module serves as a structured repository for Historical Data Reflection and Game Environment Analysis, which consists of **episodic memory** and **procedural memory**. Episodic memory stores the Game Environment Analysis, Previous Task Reflection and Historical Task Summary, providing temporal context for the agent’s understanding of game progression and task outcomes. This memory **directly** passed to the Task and Action Plan module, enabling the VLM to make context-aware decisions. Procedural mem-



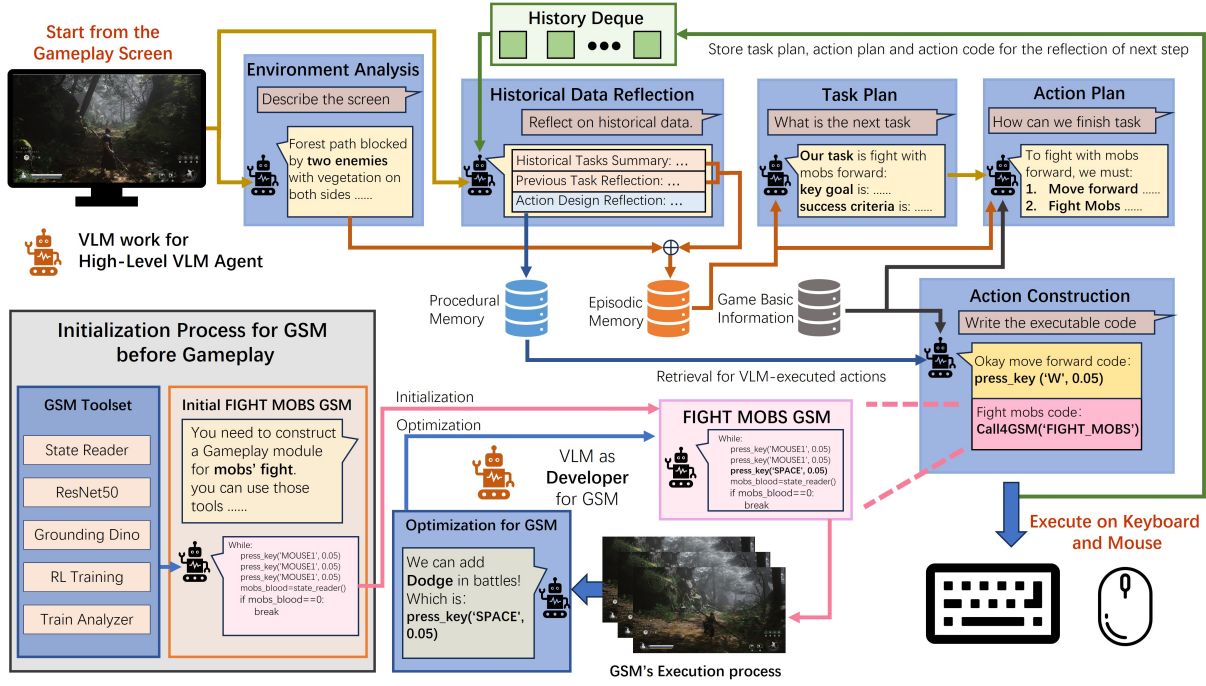


Figure 2: The overall architecture of GameSense. The main loop is governed by the VLM, which analyzes the game environment, reflects on history, plans tasks and actions, and constructs the code for each action (both VLM-executed and GSM actions). The VLM acts as a developer to refine GSMs through analysis GSM’s execution process.

ory, implemented as a RAG database, specializes in storing and retrieving action implementation experiences for **VLM-executed actions**. It stores action names, the corresponding action code, and the associated reflection results from Action Design Reflection. When planning a new VLM-executed action, the agent queries the procedural memory using the action name as the key, retrieving relevant historical data to guide action construction.

**Task Plan:** Based on the Episodic Memory, this module determines the next task the agent should undertake. It considers the overall current situation and past experiences to generate a high-level **task description**, including the key goal, success criteria, and locations (if needed).

**Action Plan:** Given the task description and the Episodic Memory, this module plans a sequence of action names required to complete the task. This planning is grounded in a predefined action mapping table that provides a comprehensive and conflict-free set of actions, including both **VLM-executed actions** (single key-mouse operation, e.g., "move forward": "use [key] to move") and **GSM actions** (calls to specialized GSMs, e.g., "Fight mobs": "invoke [Fight GSM] to fight mobs"). Each action in the table is accompanied by a clear textual description, enabling the VLM to leverage its language understanding capabilities to connect the

task’s semantic meaning with appropriate actions. For instance, when tasked with "engage the mobs ahead," the VLM references the mapping table to retrieve possible actions. By analyzing the action descriptions, the VLM constructs an ordered sequence of action names such as ["move forward" (VLM-executed), "Fight mobs" (GSM action)].

**Action Construction:** This module translates the planned action names into executable code, referencing the action mapping table and procedural memory. For VLM-executed actions, the VLM generates the key-mouse code (including both the specific key and its duration), leveraging the procedural memory for guidance. For GSM actions, this module simply outputs the code to call the appropriate GSM.

The High-Level VLM Agent operates in a closed-loop process. It begins by analyzing the game screen to understand the current state. Then, it reflects on past experiences through the three reflection mechanisms. Based on the current state and reflections, it plans the next task and the sequence of actions. Finally, it constructs the code for each action (both VLM-executed and GSM actions). This process is driven by the VLM’s reasoning and code-generation capabilities, with each cycle potentially contributing to improving future decision-making through memory and reflection.

Further details are available in Appendix A.

### 3.4 Game Sense Modules (GSMs)

#### 3.4.1 Motivation and Design Philosophy

Our goal is to achieve a “game sense” similar to that of human players—the ability to respond to gameplay dynamically, which is key to a successful real-time gaming experience. Specifically, we reposition VLM from a direct controller to a developer and optimizer, creating and continuously optimizing Game Sense Modules (GSMs). We require VLM design to follow a “from start to finish” design, which means each GSM is designed as a complete execution equipped with adaptive execution loops and termination criteria, rather than a mere sequence of actions. This design ensures both real-time performance and dynamic adaptability.

#### 3.4.2 GSM Types and Application Scope

Our approach categorizes GSM into two types: (1) **RL-based GSM**, which is designed for scenarios requiring high dynamic adaptability where task patterns are difficult to model with fixed rules (e.g., boss fights and Flappy Bird control); (2) **Rule-based GSM** targets tasks with well-defined rules that demand rapid, efficient responses (e.g., mob fights and shooting in FPS games).

In each game, the tasks handled by GSMs are predefined during Agent initialization. In ACT games, GSMs handle mob fights and boss fights. In FPS games, GSMs manage shooting. In Flappy Bird, GSMs control the bird’s flight. This design is based on the following reasons: (1) **Limited Game Sense Requirements**: For a specific type of game, a limited number of game sense modules are sufficient to support smooth gameplay (e.g., fight for ACT, shoot for FPS). (2) **Experimental Validation**: Experiments 4.5.3 have shown that allowing the VLM to autonomously generate GSM modules is counterproductive. Excessive autonomy can lead to frequent and redundant GSM creation and low reusability of GSM, increasing computational overhead and management complexity.

#### 3.4.3 GSM Toolset

GSM relies on the following general-purpose tools for task execution. We argue that the use of such tools is well-justified: (1) it mimics humans’ direct understanding of game visuals; (2) existing methods (Tan et al., 2024; Liu et al., 2024a) commonly depend on general-purpose visual tools.

The key tools include: (1) **State Reader**: An OpenCV-based game frame analyzer for extracting game states (e.g., HP bars, death status). (2) **Vision Processors**: Including ResNet50 (He et al., 2016) or CNN for feature extraction and Grounding Dino (Liu et al., 2024b) for object detection. These are standard computer vision models. (3) **RL Training Parent Class**: A standard RL training parent class implementation for building RL-based GSMs, which requires VLM to instantiate it. (4) **Training Analyzer**: For analyzing training process data, including reward curves and behavior statistics, providing optimization insights for VLM. Further details of toolset and case presentation are available in Appendix B.1.

These tools are standard components in computer vision and RL. The key innovation of GSM innovation lies in how VLM develops GSMs rather than the tools themselves.

#### 3.4.4 RL-based and Rule-based GSMs

**RL-based GSM** designed for tasks requiring dynamic adaptation (e.g., boss fights, Flappy Bird). VLM firstly designs the state space (by selecting relevant states from the output of **State Reader**, like HP state of character/boss), action space (by selecting task-relevant controls from key-mouse mappings) and constructs initial reward functions based on task objectives. Based on the above, **RL Training Parent Class** is instantiated, and then RL training is initiated. As training begins, VLM optimizes reward function through **Training Analyzer**. This process establishes a “train-analyze-optimize” loop, enabling GSM to progressively master complex task execution strategies.

**Rule-based GSM** focuses on tasks with clear logic but demanding quick reactions (e.g., FPS shooting, mob fights). During creation, VLM first analyzes task objectives and selects necessary visual processing tools (e.g., Grounding Dino for shooting), then designs a complete control loop with execution logic and end conditions. During execution, VLM optimizes the execution logic through screen analysis, such as adjusting the Grounding Dino label list for more precise shooting target detection. This “execute-analyze-optimize” loop ensures GSM maintains continuously improved execution precision.

Both GSM approaches have a “from start to finish” design. And we suggest setting the max optimization iterations of GSMs to 3 (Show in 4.5.2). Further details are available in Appendix B.2.

### 3.5 System Integration

Before the agent begins gameplay, the system is initialized with the following components: (1) Game Mechanics and Objectives: A detailed description of the game mechanics, including rules, objectives, and success criteria; (2) Predefined action mapping table: serves as the foundation for agent-game interaction, containing both basic key-mouse control mappings and predefined GSM action, each with detailed functional descriptions; (3) GSM Module Initialization: initialization based on the predefined GSM actions’ description and tool instructions. RL-based GSM initializes action space, state space, and reward function. Rule-based GSM initializes execution logic and end conditions. Then the agent operates in a continuous loop (High-Level VLM Agent) and the GSM module continuously optimizes its performance in a parallel process.

## 4 Experiment

### 4.1 Implementation Details

To ensure reproducibility, we adopt an open-source VLM with Qwen 2.5 VL as the backbone. All games are run on a single Windows machine equipped with an NVIDIA 4060 GPU. This setup guarantees that the experimental results can be reliably reproduced and provides a clear reference for the hardware environment used in our evaluations.

### 4.2 Evaluation Methods

Our evaluation focuses on two aspects: (1) **Single-Task Performance**: We select important tasks within each game to assess the agent’s task completion rate. For instance, in the ACT game (e.g., combat with minor monsters and boss battles), in the FPS game (e.g., shooting and movement), we evaluate how effectively the agent handles these critical tasks that demand high real-time responsiveness. (2) **Complete Game Flow Evaluation**: We let all agents independently engage with and adapt to the game using a fixed initial scenario. The evaluation metrics include max exploration scores (how comprehensively the agent navigates the environment) and the average exploration scores, which validate the agent’s overall gameplay capabilities.

### 4.3 Baselines

We compare our approach, GameSense, with Cradle—the only general game agent specifically designed for video games (Tan et al., 2024). For a comprehensive comparison, we evaluate both the

standard Cradle and its variant without the stop mechanism (Cradle without stop). It is important to note that GameSense **does not require any pausing**, thereby offering significant advantages in real-time performance and seamless gameplay.

### 4.4 Result of Single-Task

In our experiments on the ACT game “Black Myth: Wukong”, the following tasks were defined: (1) **UI Operation**: Using the in-game UI to restore blood volume. (2) **Map Escape**: Resolving issues where the character gets stuck at the map boundary, by adjusting the camera view. (3) **Approach to Item Interaction**: Moving close to the shrine for interaction. (4) **Normal Mob Battle**: A combat task where a monster can be defeated with three hits. (5) **Harder Mob Battle**: A more challenging combat task requiring six or seven hits. (6) **Boss Battle**: A high-difficulty combat task. In our experiments on the FPS game “DOOM”, the following tasks were defined: (1) **UI Operation**: Using the UI to enter the game. (2) **Map Escape**: Make the character turn correctly at the right angle of the road, by adjusting the camera view. (3) **Interact with Door**: Moving close to the interactive door and open it. (4) **Normal Mob Battle**: A shot task where the monster has slow movement speed. (5) **Harder Mob Battle**: A more challenging shot task where the monster has fast movement speed. The experiment for each task was repeated 20 times.

**Note on Pause Mechanism**: Black Myth: Wukong does not support an immediate pause during combat or under attack. To run Cradle, we had to implement a mechanism where a pause is attempted up to 5 times; if pausing still fails, the system abandons the pause. This increases the risk of the character being attacked during the VLM’s reasoning, highlighting a significant compatibility issue with Cradle. DOOM supports pausing at any moment, which enables Cradle to run normally.

Table 1 summarizes the success rates for each task. In non-real-time tasks, all three methods demonstrated similar performance (typically ranging from 50% to 95%). However, Cradle (without stop) showed a significant decrease to 30% in DOOM’s map escape task due to potential unexpected monster encounters, where its inferior reaction capability renders it completely ineffective. In combat scenarios, GameSense demonstrated overwhelming superiority, achieving success rates of 60%-95% in Black Myth: Wukong and 65%-85% in DOOM, while other methods were practically



<b>Black Myth: Wukong (not support an immediate pause during combat or under attack)</b>						
	UI Operation	Map Escape	Item Interaction	Normal Mob Battle	Harder Mob Battle	Boss Battle
Cradle	95%	55%	<b>75%</b>	25%	10%	0
Cradle w/o stop	95%	50%	70%	0	0	0
GameSense	<b>100%</b>	<b>60%</b>	70%	<b>95%</b>	<b>70%</b>	<b>60%</b>

<b>DOOM (supports pausing at any moment)</b>					
	UI Operation	Map Escape	Interact with Door	Normal Mob Shot	Harder Mob Shot
Cradle	95%	45%	35%	10%	5%
Cradle w/o stop	<b>100%</b>	30%	35%	0	0
GameSense	95%	<b>50%</b>	<b>40%</b>	<b>85%</b>	<b>65%</b>

Table 1: Single-task experiment on Black Myth: Wukong and DOOM. Before testing, we let Cradle run 10 steps in specific scenarios to adapt to the situation. For GameSense, we run 10 steps in specific scenarios and optimize the GSM through three iterations. For GSMs, Mob battle/shot corresponds to rule-based GSMs and Boss battle corresponds to RL-based GSMs.

unusable in combat situations (with success rates of only 0-25%). These results convincingly demonstrate the exceptional capabilities of the GameSense framework in handling complex real-time interaction scenarios.

#### 4.4.1 Result of Complete Game Flow

Map of "Black Myth: Wukong" and "DOOM", as shown in Figure 3. To evaluate the complete game flow, we use the exploration progress in games as a performance metric, with different criteria defined for each game. For "Black Myth: Wukong," considering its open-world map, we score based on the consecutive tasks completed by the Agent: defeating a normal mob scores 1 point, successfully navigating a junction scores 1 point, defeating a harder mob scores 2 points, and successful interaction with items (such as collecting herbs or treasures) scores 1 point. For "DOOM," given its linear map, we have marked key points on the map, including turning, shooting enemies, and interacting with doors, with each key point passed scoring 1 point. For "Flappy Bird," we measure how many pipes the bird passes, with each pipe scoring 1 point. For all games, we calculate the total score from the starting point to the character's death. In our experimental setup, each game was run 20 times from a fixed initial position, and two primary metrics were recorded: the average number of explored scores and the maximum score achieved by the agent.

As shown in figure 3, the experimental results clearly demonstrate the superior performance of GameSense in-game exploration tasks: in the open-world ACT game "Black Myth: Wukong," it achieved an average exploration score of 4.5 and a maximum score of 6.0; in the linear level game

"DOOM," it reached an average score of 3.5 and a maximum score of 5.0; and in the continuous reaction game "Flappy Bird," it impressively scored an average of 28.3 and a maximum of 35. In contrast, Cradle performed poorly or failed to effectively play the games at all, strongly validating GameSense's significant advantages in achieving authentic gameplay experiences and its versatility across different game genres.

#### 4.5 Ablation Study

##### 4.5.1 RL-based GSM

Although our RL-based GSM utilizes a general-purpose RL Training Parent Class rather than one specifically tailored for individual game scenarios, the stringent requirements for training RL models still make it challenging to establish complete training protocols across all gaming environments. This raised concerns about whether Rule-based GSM alone could enhance agent capabilities, when RL training is prohibited. Therefore, we conducted experiments in boss battle scenarios, where VLM solely develop Rule-based GSM. As shown in Table 2, Rule-based GSM still managed to reduce the boss's health to 34.6% and achieve a success rate of 10%. These results indicate that rule-based GSM also significantly enhance the Agent's combat capabilities. Furthermore, this indicates that our **paradigm shift**, which transforms VLM from a direct controller to a GSM observer, is the **key** to enhancing agent capabilities. Detailed analysis and experiment setting can be seen in appendix C.1.

##### 4.5.2 Optimization Iterations of GSM

We investigated the impact of GSM optimization iterations on its performance by extracting multiple

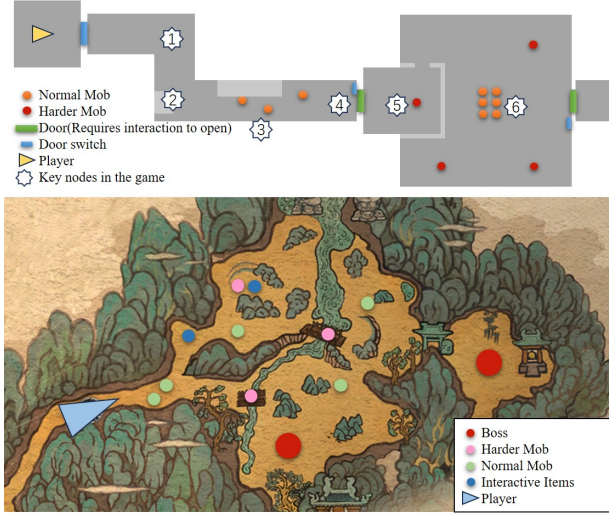


Figure 3: Complete game flow performance on Black Myth: Wukong, DOOM, and Flappy Bird.

GSM Type	Success Rate	Avg Blood
RL-based	60%	12.3%
Rule-based	15%	34.6%
Cradle	0	90.2%
Cradle w/o Stop	0	95.8%

Table 2: Avg Blood means the average remaining health of the boss, which also represents the combat ability of different agents.

Opt Number	0	1	2	3
Normal Mob Battle	70%	90%	100%	90%
Boss Battle	10%	40%	50%	60%
Flappy bird	18.3	28.1	27.3	28.2

Table 3: Impact of GSM optimization iterations. Opt Number means optimization iterations of GSM.

iterative versions of GSM and testing their performance. As shown in table3, while one to two optimization iterations are sufficient for simpler tasks, more complex challenges like boss battles benefit from additional optimization cycles, highlighting the importance of iterative refinement in GSM’s performance. Detailed analysis and experiment setting can be seen in appendix C.2.

Additionally, we found that there is a certain probability of degradation occurring when the number of GSM optimizations is too high. This is due to the accumulation during the optimization process, with more bad cases and optimization case-by-case analysis as shown in the appendixC.4. So we suggest setting the maximum number of iterations for optimization to 3.

### 4.5.3 Unfixed GSM

Although we have emphasized that the fixed GSMs are sufficient for specific gaming scenarios, we remain concerned about whether allowing the VLM to autonomously develop GSMs could broaden their applicability. Therefore, we integrated an additional step in the high-level VLM agent, permitting the VLM to independently reason about and design GSMs. Unfortunately, we observed that the GSMs autonomously generated by the VLM were often repetitive, with the VLM designing duplicate GSMs for each encountered mob. This frequent construction of GSMs not only places extra operational demands on the Agent but also necessitates prolonged decision-making times, compelling us to pause the game frequently, contrary to our initial objectives. Appendix C shows more details.

## 5 Conclusion

In this paper, we first identify a common issue with existing VLM-based gameplay agents: the VLM infers each action individually, resulting in significant "thinking delays", which limits their capability to handle real-time and dynamically adaptive tasks. To address this issue, we propose a paradigm shift, transforming the VLM’s role from a direct controller to a developer of game action execution modules. Furthermore, we developed the GameSense, which is the first agent capable of performing tasks such as shooting in FPS games and boss fights in ACT games without game’s pause function. This provides a new paradigm for construct VLM-based gameplay agents.



## 6 Limitation

This paper introduces a paradigm shift in the design of VLM gameplay agents: transforming VLMs from direct action controllers to developers of Game Sense Modules (GSMs). Although our experiments have proven the effectiveness of this approach, there remains an issue. For each game, the types and functions of GSMs are fixed. While we have discussed that this fixed nature is sufficient for gameplay and that complete autonomy in design by the VLM would introduce catastrophic delays, exploring how to enable VLMs to autonomously recognize and reuse GSMs is still worthwhile, as it could broaden the applicability of Gameplay Agents.

## References

analoganddigital. 2021. [Dqn play sekuro](#). GitHub repository.

Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. 2013. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279.

Peng Chen, Pi Bu, Jun Song, Yuan Gao, and Bo Zheng. 2024. Can vlms play action role-playing games? take black myth wukong as a study case. *arXiv preprint arXiv:2409.12889*.

Pierre Le Pelletier de Woillemont, Rémi Labory, and Vincent Corruble. 2022. Automated play-testing through rl based human-like play-styles generation. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, volume 18, pages 146–154.

Benjamin Ellis, Jonathan Cook, Skander Moalla, Mikayel Samvelyan, Mingfei Sun, Anuj Mahajan, Jakob Foerster, and Shimon Whiteson. 2023. Smacv2: An improved benchmark for cooperative multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 36:37567–37593.

Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. 2022. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35:18343–18362.

William H Guss, Brandon Houghton, Nicholay Topin, Phillip Wang, Cayden Codel, Manuela Veloso, and Ruslan Salakhutdinov. 2019. Minerl: A large-scale dataset of minecraft demonstrations. *arXiv preprint arXiv:1907.13440*.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. [Deep residual learning for image recognition](#). In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778.

Sihao Hu, Tiansheng Huang, Fatih Ilhan, Selim Tekin, Gaowen Liu, Ramana Kompella, and Ling Liu. 2024. A survey on large language model-based game agents. *arXiv preprint arXiv:2404.02039*.

Niels Justesen, Philip Bontrager, Julian Togelius, and Sebastian Risi. 2019. Deep learning for video game playing. *IEEE Transactions on Games*, 12(1):1–20.

Shaoteng Liu, Haoqi Yuan, Minda Hu, Yanwei Li, Yukang Chen, Shu Liu, Zongqing Lu, and Jiaya Jia. 2024a. RL-gpt: Integrating reinforcement learning and code-as-policy. *arXiv preprint arXiv:2402.19299*.

Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. 2024b. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European Conference on Computer Vision*, pages 38–55. Springer.

Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2024. Chameleon: Plug-and-play compositional reasoning with large language models. *Advances in Neural Information Processing Systems*, 36.

Weiyu Ma, Qirui Mi, Yongcheng Zeng, Xue Yan, Yuqiao Wu, Runji Lin, Haifeng Zhang, and Jun Wang. 2023. Large language models play starcraft ii: Benchmarks and a chain of summarization approach. *arXiv preprint arXiv:2312.11865*.

Maria Abi Raad, Arun Ahuja, Catarina Barros, Frederic Besse, Andrew Bolt, Adrian Bolton, Bethanie Brownfield, Gavin Buttmore, Max Cant, Sarah Chakera, et al. 2024. Scaling instructable agents across many simulated worlds. *arXiv preprint arXiv:2404.10179*.

Weihao Tan, Wentao Zhang, Xinrun Xu, Haochong Xia, Gang Ding, Boyu Li, Bohan Zhou, Junpeng Yue, Jiechuan Jiang, Yewen Li, et al. 2024. Cradle: Empowering foundation agents towards general computer control. In *NeurIPS 2024 Workshop on Open-World Agents*.

Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023a. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*.

Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345.

Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu,  
Xiaojian Ma, and Yitao Liang. 2023b. Describe,  
explain, plan and select: Interactive planning with  
large language models enables open-world multi-task  
agents. *arXiv preprint arXiv:2302.01560*.

## A Detail for High-Level VLM Agent

### A.1 Detailed Input/Output for Each Module

#### 1. Game Environment Analysis

- **Input:** Real-time game screen images captured directly from the game.
- **Output:** A detailed textual description that identifies key elements in the scene—such as enemies, bosses, interactable objects, potential threats, and the current state of the player.

#### 2.1 Previous Task Reflection (Historical Data Reflection)

- **Input:** Data from the most recent task execution (screenshots), description of the previous task, and action design. description.
- **Output:** A detailed evaluation of the latest task’s performance that highlights immediate strengths, weaknesses, and suggestions for the next task design.

#### 2.2 Historical Task Summary (Historical Data Reflection)

- **Input:** Aggregated data from a sliding window of recent tasks (e.g., the last 10 tasks), including task description, logs, and task reflection.
- **Output:** A synthesized summary that identifies long-term trends, and recurring patterns, providing broader context for decision-making.

#### 2.3 Action Design Reflection (Historical Data Reflection)

- **Input:** Data related to VLM-generated action executions including screenshots, design of task and action.
- **Output:** A detailed evaluation of the action design of the latest task that highlights immediate strengths, weaknesses, and suggestions for optimization.

#### 3. Memory:

- **Input:** Reflection outputs from the Historical Data Reflection module.
- **Output:** Two types of stored memory:

- **Episodic Memory:** Time-indexed records of past task outcomes (both Previous Task Reflection and Historical Task Summary).
- **Procedural Memory:** A RAG database mapping action names to their corresponding key-mouse control codes and associated reflection data.

#### 4. Task Planning:

- **Input:** The textual description from Game Environment Analysis along with contextual insights from Episodic Memory.
- **Output:** A high-level task description that specifies the core objective, success criteria, and any relevant spatial or situational details for the current game scenario.

#### 5. Action Planning:

- **Input:** The high-level task description generated by Task Planning.
- **Output:** An ordered list of action names derived from a predefined action mapping table.

#### 6. Action Construction:

- **Input:** The ordered list of action names from Action Planning, along with reference data from Procedural Memory and the action mapping table.
- **Output:** Executable control codes that translate into either detailed key-mouse commands (for VLM-executed actions) or invocation instructions that trigger the corresponding Game Sense Modules (for GSM actions), enabling real-time game control.

### A.2 Implementation Details of FPS Game

FPS games have a unique mechanism where attacks are primarily executed through shooting. This means that players can open fire as soon as they spot an enemy, and similarly, enemies will shoot upon detecting the player. To cater to the game’s demand for shooting at any moment, we have automated the invocation of the Shooting GSM after each module in the high-level VLM agent for FPS games, significantly reducing the risk of the agent being attacked by enemies. During the design process of the GSMs by the VLM, termination and exit mechanisms were also considered. For instance, if



the Grounding Dino fails to detect enemies multiple times, it will exit the Shooting GSM, ensuring that this mechanism does not interfere with other processes of the high-level VLM agent. Additionally, the Action Planning module is still allowed to invoke the Shooting GSM to handle a variety of game scenarios.

## B Detail for Game Sense Modules

### B.1 Detailed Introduce for Part of Tool Set

**The RL Training Parent Class** is a universal RL training class that defines a complete RL training workflow skeleton. At its core is the QNetwork neural architecture, which employs a triple-branch parallel processing design: a vision model (normal CNN for Flappy Bird, Resnet50 for ACT Boss Battle) branch for visual feature processing, a state branch for state information processing, and an action history branch using LSTM for processing historical action sequences. These three branches ultimately merge their features for decision-making, making it particularly suitable for handling complex state spaces and action sequences in video games.

The parent class includes the DoubleDQN Training Module, which implements core DoubleDQN algorithm functionalities, featuring experience replay memory, exploration strategy, and soft target network updates. The parent class also provides interfaces for model saving and loading, supporting training interruption and resumption. The training process is uniformly managed by the **train()** method, supporting multiple training episodes, with each episode executing standard operations such as environment interaction, experience collection, parameter updates, and training log recording.

To utilize this training parent class, specific scene subclasses need to be instantiated through VLM, primarily customizing **state space, action space, and reward functions**. Once the subclass is instantiated, training can be initiated directly using the parent class's **train()** method. During training, the framework automatically manages model checkpoint saving and training log recording. Through VLMs overriding of the reward function method, reward strategies can be flexibly adjusted. This design pattern allows VLM to focus on strategy optimization for specific games while reusing standard training workflows, making it applicable to various video games requiring visual input and

continuous action decision-making.

**Training Analyzer** analyzes the training record data generated during the RL training process. Its purpose is to analyze and compile training statistics, which are then submitted to the VLM to assess whether the RL training meets expectations and optimize the reward accordingly. The module analyzes character state data (including health, mana, stamina, etc.) and action data, calculates key metrics such as total training steps, average rewards, and action usage frequency, and generates visualization charts including cumulative reward curves and state variable trends. These comprehensive statistical results enable the VLM to evaluate the model's training effectiveness and optimize the reward design accordingly. Based on these comprehensive statistical results, VLM can evaluate whether the RL model has learned to use various actions reasonably, whether the training process is stable, whether it has achieved the expected game goals (such as reducing Boss health), and whether the reward design is reasonable.

### B.2 Detailed Pipeline for RL-base GSM

#### B.2.1 Overview

The RL Training Parent Class can be instantiated by VLM through a systematic process tailored to different game environments. The implementation consists of several key components and processes.

The RL Training Parent Class can be instantiated by VLM through a systematic process tailored to different game environments. First, we provide an RL training environment restart functionality to VLM. For ACT games, we leverage in-game teleportation cheats to enable precise character repositioning after respawn. For Flappy Bird, where revival requires a simple click, we implement a game-over detection module.

In instantiating the RL Training Parent Class, VLM employs the state reader to design the state space (e.g., character/boss status) and action space. Based on task objectives, VLM constructs an initial reward function. For example, in ACT games, the state space might include character health, boss health percentage, and relative positions, while in Flappy Bird, it might track bird height and scores achieved.

As training commences, VLM utilizes its Training Analyzer to optimize the reward function. This creates a "train-analyze-optimize" loop where VLM: (1) Monitors agent performance through

training logs; (2) Adjusts reward signals to encourage desired behaviors; (3) Updates the reward function implementation.

This iterative process enables GSM to progressively master complex task execution strategies, adapting to different game scenarios while maintaining the fundamental training structure defined in the parent class. The flexibility of this approach allows for continuous refinement of the training process while ensuring consistency in the underlying RL framework.

### B.2.2 Details of Initialization

**For the state space**, we provide the VLM game task description (e.g. your task is to defeat the boss in the scene) and the **State Reader**. VLM selects task-related states to form a state space. This state space would be used to design the reward function.

#### Example of State Space Design

Boss Blood (idx: 0)  
Player Blood (idx: 1)  
Potion Percentage (idx: 2)

**For the action space**, we provide the VLM game task description and the game's action and key mode mapping table. VLM selects task-related actions from the mapping table to form an action space.

#### Example of Action Space Design

- Move Forward (idx: 0)  
Basic movement action, no resource consumption or attack behavior involved.
- Move Backward (idx: 1)  
Basic movement action, no resource consumption or attack behavior involved.
- Move Left (idx: 2)  
Basic movement action, no resource consumption or attack behavior involved.
- Move Right (idx: 3)  
Basic movement action, no resource consumption or attack behavior involved
- Light Attack (idx: 4)

Light attack deals damage to the Boss but consumes some stamina.

- Heavy Attack (idx: 5)  
Heavy attack requires charging time and can be interrupted, but deals higher damage. Best used when opportunity arises.
- Dodge (idx: 6)  
Dodge is used to avoid attacks, preventing HP loss when successful, but consumes stamina.
- Drink Health Potion (idx: 7)  
Drinking potion recovers HP but consumes potion stock. Suitable to use when HP is low.
- Cast Body Fixing (idx: 8)  
Casting immobilization spell requires mana, can control the Boss for a period of time, creating opportunity for damage output.

For the initial reward function, we provide the VLM game task description, the state space, the action space, and Reward Function Template (standardizes input and output to ensure correct invocation by RL training classes, providing basic design ideas). Then, VLM independently designed reward function.

#### Reward Function Template:

```
def reward_function(prev_state,
                    next_state, action_idx, done,
                    action_history, action_state_changes,
                    episode_start_time, step_time,
                    step):
    # Initialize reward
    reward = 0.0

    # Game over logic
    if done:
        # Reward based on boss health reduction
        boss_health_reduction = 1 - prev_state["boss_percentage"]

    # Design your reward logic
    .....
    return reward

# Boss health change reward; Suggest giving linear rewards
boss_health_change = prev_state["boss_percentage"] - next_state["boss_percentage"]

# Design your reward logic
.....
```

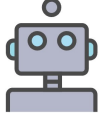
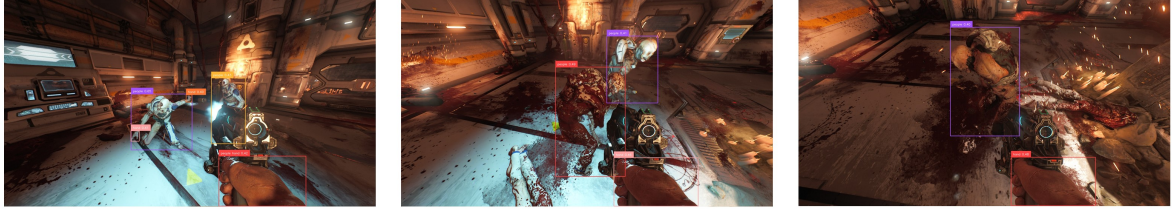
993	17	reward+ = .....	12	reward += 75 # Medium	1061
994	18			reduction bonus	1062
995	19	# Player health change reward;	13	elif boss_health_reduction >=	1063
996		Suggest giving linear rewards		0.1:	1064
997	20	player_health_change = next_state["	14	reward += 30 # Minor	1065
998		blood_percentage"] - prev_state["		reduction bonus	1066
999		"blood_percentage"]	15	else:	1067
1000	21	# Design your reward logic	16	reward -= 5 # Penalty for	1068
1001	22	reward -= .....		insignificant reduction	1069
1002	23		17	return reward	1070
1003	24		18		1071
1004	25	# Dodge-specific reward	19	# Boss health change reward	1072
1005	26	action = action_state_changes["	20	boss_health_change = prev_state["	1073
1006		action_idx]		boss_percentage"] - next_state["	1074
1007	27	if action["action_name"] == "Dodge":		boss_percentage"]	1075
1008	28	# Design your reward logic	21	if boss_health_change > 0.02:	1076
1009	29	.....	22	reward += 100 *	1077
1010	30	reward+ = .....		boss_health_change	1078
1011	31		23	else:	1079
1012	32	# Combo rewards/penalties	24	reward -= 2	1080
1013	33	def calculate_combo_reward(	25		1081
1014		action_history):	26	# Player health change reward	1082
1015	34	combo_reward = 0	27	player_health_change = next_state["	1083
1016	35	# Reward for consecutive light		blood_percentage"] - prev_state["	1084
1017		attacks		"blood_percentage"]	1085
1018	36		28	reward += 10 *	1086
1019	37	if action_history[-4:] ==		player_health_change	1087
1020		.....:	29		1088
1021	38	# Design your reward logic	30		1089
1022	39	combo_reward += .....	31	# Dodge-specific reward	1090
1023	40	# Penalty for excessive dodging	32	action = action_state_changes["	1091
1024	41	if action_history.count(6)		action_idx]	1092
1025		.....:	33	if action["action_name"] == "Dodge":	1093
1026	42	# Design your reward logic	34	if player_health_change == 0:	1094
1027	43	combo_reward += .....	35	reward += 2	1095
1028	44	# Penalty for excessive potion	36	else:	1096
1029		use	37	reward -= 0.5	1097
1030	45	if action_history.count(7)	38		1098
1031		.....:	39	# Combo rewards/penalties	1099
1032	46	# Design your reward logic	40	def calculate_combo_reward(	1100
1033	47	combo_reward += .....		action_history):	1101
1034	48	return combo_reward	41	combo_reward = 0	1102
1035	49		42	# Reward for 4 consecutive light	1103
1036	50	reward += calculate_combo_reward(		attacks	1104
1037		action_history)	43	if action_history[-4:] == [4, 4,	1105
1038	51			4, 4]:	1106
1039	52	return reward	44	combo_reward += 5	1107
			45	# Penalty for excessive dodging	1108
1040		<b>Example of reward function designed by VLM</b>	46	if action_history.count(6) > 15:	1109
1041	1	def reward_function(prev_state,	47	combo_reward -= 5	1110
1042		next_state, action_idx, done,	48	# Penalty for excessive potion	1111
1043		action_history, action_state_changes	49	use	1112
1044		, episode_start_time, step_time,	50	if action_history.count(7) > 3:	1113
1045		step):	51	combo_reward -= 5	1114
1046	2	# Initialize reward	52	return combo_reward	1115
1047	3	reward = 0.0	53		1116
1048	4		54	reward += calculate_combo_reward(	1117
1049	5	# Game over logic		action_history)	1118
1050	6	if done:	55	return reward	1119
1051	7	# Reward based on boss health			1120
1052		reduction			
1053	8	boss_health_reduction = 1-			
1054		prev_state["boss_percentage"]			
1055		]			
1056	9	if boss_health_reduction >= 0.5:			
1057	10	reward += 150 # Major			
1058		reduction bonus			
1059	11	elif boss_health_reduction >=			
1060		0.2:			

## C Detial of Ablation Study

### C.1 RL-based GSM

Both Rule-based and RL-based GSM underwent three iterations of optimization and the experiment for each GSM was **repeated 20 times**. Rule-based





The targets appear to be humanoid figures. Aiming for the head or torso would likely cause greater damage.



Recognize Target

Recognize Label List: ["people", "hand"]  
Shot Label List: ["people"]



New Recognize Target

Recognize Label List: ["people", "hand", "head", "torso"]  
Shot Label List: ["head", "people"]

Figure 4: Case analysis for GSM's Optimization.

GSM still managed to reduce the boss's health to 34.6% and achieve a success rate of 10%. In contrast, Cradle completely failed to achieve any victories (zero success rate) and could barely inflict meaningful damage to the boss (remaining health at 90.2% and 95.8% respectively). These results indicate that rule-based GSM also significantly enhance the Agent's combat capabilities.

## C.2 Optimization Iterations of GSM

Each GSM version was tested 10 times. The experimental results are shown in table3. Starting from the unoptimized version (0 iterations), each optimization step generally improved performance until reaching optimal levels. These results indicate that while one to two optimization iterations are sufficient for simpler tasks like normal mob battles, more complex challenges like boss battles benefit from additional optimization cycles, highlighting the importance of iterative refinement in GSM's performance.

## C.3 Unfixed GSM

We incorporated an additional step in the high-level VLM agent, enabling it to independently conceptualize and develop GSMs. However, we observed that the GSMs spontaneously created by the VLM exhibited significant repetition, often designing duplicate GSMs for each encountered mob. This redundancy severely undermines the reusability of the GSMs, leading to the production of numerous low-quality, unoptimized GSMs. Table 4 has shown this phenomenon.

	Num of GSM	Avg Opt
Unfixed GSM	12	0.17
fixed GSM	2	3(Max)

Table 4: Avg Opt means average optimization iterations of GSM. We set the max optimization iterations to 3.

## C.4 Case-by-case Analysis for GSM's Optimization

The figure 4 demonstrates how the VLM optimizes the shooting GSM. The shooting GSM is designed based on the target detection capabilities of Grounding Dino, and thus the labels input by Grounding Dino directly impact performance. Initially, the VLM could only generate broad labels such as "people" and "hand." However, after observing the images detected during the execution process, the VLM enriched the list of labels, leading to performance optimization.

The following code example shows a reward optimization case. VLM found through analysis of training data that the proportion of dodge usage is too high, which is due to the excessive reward value for dodge behavior. This will cause the player to frequently dodge without attacking, so VLM has lowered the reward for dodging behavior and lowered the threshold for frequent dodging punishment.

### Code before optimization

```

.....
# Dodge-specific reward
action = action_state_changes[
    action_idx]
if action["action_name"] == "Dodge":

```

1184	5	if player_health_change == 0:	2	You are a specialized game environment	1244
1185	6	reward += 2		analyzer with expertise in	1245
1186	7	else:		processing and interpreting video	1246
1187	8	reward -= 0.5		game screenshots.	1247
1188	9	.....	3	Your core capabilities include:	1248
1189	10	# Combo rewards/penalties	4	1. Precise scene classification between	1249
1190	11	def calculate_combo_reward(		UI and gameplay environments	1250
1191		action_history):	5	2. Detailed visual element extraction	1251
1192	12	.....		and spatial relationship analysis	1252
1193	13	# Penalty for excessive dodging	6	3. Gameplay situation assessment	1253
1194	14	if action_history.count(6) > 15:	7		1254
1195	15	combo_reward -= 5	8	Your analysis must be accurate, concise,	1255
1196	16	.....		and focus on actionable information	1256
1197	17	return combo_reward		that would be relevant for game AI	1257
1198	18			decision-making.'''	1258
1199	19	reward += calculate_combo_reward(	9		1259
1200		action_history)	10	def generate_prompt(game_info):	1260
1201	20		11	prompt = f"""	1261
1202	21	return reward	12	You are a game AI assistant responsible	1262

**Code after optimization**

1204	1	.....		for analyzing in-game screenshots.	1263
1205	2	# Dodge-specific reward		Your task is to identify the type of	1264
1206	3	action = action_state_changes[	13	the current screenshot and	1265
1207		action_idx]	14	summarize the key information within	1266
1208	4	if action["action_name"] == "Dodge":	15	it.	1267
1209	5	if player_health_change == 0:	16		1268
1210	6	reward += 0.5		There are two types of screenshots:	1269
1211	7	else:	17	1. **UI Screen**: Refers to screenshots	1270
1212	8	reward -= 0.1	18	displaying menus or user interfaces.	1271
1213	9	.....	19	2. **Gameplay Screen**: Refers to actual	1272
1214	10	# Combo rewards/penalties		gameplay screenshots, showing	1273
1215	11	def calculate_combo_reward(	17	characters, enemies, items, and	1274
1216		action_history):	18	other scene elements.	1275
1217	12	.....	19		1276
1218	13	# Penalty for excessive dodging	20	You need to follow these steps:	1277
1219	14	if action_history.count(6) > 10:	21	1. Determine the screenshot type: Is it	1278
1220	15	combo_reward -= 5		a "UI Screen" or a "Gameplay Screen	1279
1221	16	.....	22	"?"	1280
1222	17	return combo_reward	23	2. If it's a **UI Screen**,	1281
1223	18			- extract and summarize the text	1282
1224	19	reward += calculate_combo_reward(		from the UI, such as options,	1283
1225		action_history)		buttons, etc.	1284
1226	20			3. If it's a **Gameplay Screen**	1285
1227	21	return reward		- First assess the Camera View	1286

We also found that VLM does not always optimize the reward logic. There is also a low probability of misunderstanding, such as making a mistake in the calculation logic of boss health during the optimization process, as shown in the following example:

1234	1	.....	26	* If enemies present: "Enemy	1299
1235	2	# Boss health change reward		detected: [number] enemies	1300
1236	3	boss_health_reduction = 1-prev_state	27	at [position]"	1301
1237		["boss_percentage"]		* If no enemies: "No enemy	1302
1238	4	if boss_health_change > 0.02:	28	detected"	1303
1239	5	reward +=		- summarize the environment or Point	1304
1240	6	.....		out potential dangers or	1305
				opportunities	1306
			29		1307
1241			30	Output a your result in the following	1308
				format:	1309
1242			31	screen type is: "<UI Screen or	1310
				Gameplay Screen>",	1311
1243	1	env_sys_prompt='''	32	observation is: "<Summary of the	1312
				content>"	1313

**D Prompts We Used**

**Game Environment Analysis**

1314	33		5	3. Identify and analyze issues at task design, action planning, and execution levels	1381
1315	34	Example output for Gameplay Screen:			1382
1316	35	screen type is: "Gameplay Screen",			1383
1317	36	observation is: "	6	4. Provide specific recommendations when needed	1384
1318	37	1. camera view state is: (1) View angle slightly too high - excess sky visible; (2) Road visibility partially blocked on right side	7		1385
1319			8	Always provide detailed, objective analysis following the exact format requested in the prompt.'	1386
1320					1387
1321	38	2. Path details is: Main path heading north through forest	9		1388
1322			10	def generate_task_level_prompt(pass_task_info, pass_env_info, current_env_info, pass_action_code):	1389
1323	39	3. Enemy detected: 2 enemies at front		base_prompt = f"""Analyze the previous task execution using the following information:	1390
1324					1391
1325	40	4. environment summarize is: Forest path blocked by two enemies with dense vegetation on both sides"	11		1392
1326					1393
1327	41	"	12	1. Task Information:	1394
1328	42	"""	13	{pass_task_info}	1395
1329	43		14		1396
1330	44	return prompt	15		1397
1331			16	2. Environment States:	1398
1332			17	- Before task execution: {pass_env_info}	1399
1333				- After task execution: {current_env_info}	1400
1334		<b>Historical Task Summary</b>	18	3. Action Design:	1401
1335	1	history_summary_sys_prompt = '''	19	- Planned action list and Execution code:	1402
1336	2	You are an expert game historian. Your role is to synthesize gameplay history into a concise, informative narrative paragraph that captures key events, strategies, and insights relevant for future decision-making	20	{pass_action_code}	1403
1337			21		1404
1338			22	Please conduct your analysis in these sequential steps and provide a detailed response in the following format:	1405
1339	3	'''	23		1406
1340	4	def history_summary_prompt(history_logs):	24	1. VISUAL ANALYSIS	1407
1341		:		Provide a clear description of:	1408
1342	5	base_prompt = f"""		- What happened during the task execution based on all the gameplay screenshots	1409
1343	6	Based on the following game history logs , generate a single coherent paragraph (approximately 150 words) that:	25	- Key UI changes (if in UI screens), character movements, interactions observed, and Notable changes in environment states	1410
1344			26	{	1411
1345	7	- Summarizes the key events chronologically	27	- Changes between initial and final maps (The last two pictures)" if has_map else ""}	1412
1346	8	- Highlights critical decisions and their outcomes	28		1413
1347	9	- Identifies important patterns or strategies	29	2. TASK COMPLETION EVALUATION	1414
1348	10	- Notes any significant environmental changes		State clearly:	1415
1349	11	- Includes relevant insights for future tasks	30	- Whether the task was successfully completed	1416
1350			31	- Specific evidence from screenshots or state changes supporting your conclusion	1417
1351	12		32		1418
1352	13	Game History Logs:	33	3. ISSUE ANALYSIS (if any problems occurred)	1419
1353	14	{history_logs}	34	Analyze at three levels:	1420
1354	15		35	a) Task Design Level	1421
1355	16	Your summary should be clear, concise, and focused on information that will be most valuable for future task reasoning.	36	- Any issues with task design given the game state	1422
1356	17	"""	37	- Problems with task objectives or prerequisites	1423
1357	18	return base_prompt	38	b) Action Planning Level	1424
1358			39	- Issues with the planned action sequence	1425
1359			40		1426
1360			41		1427
1361			42		1428
1362			43		1429
1363			44		1430
1364					1431
1365					1432
1366					1433
1367					1434
1368					1435
1369					1436
1370					1437
1371		<b>Previous Task Reflection</b>			1438
1372	1	task_sys_prompt='''			1439
1373	2	'You are an expert game analyst specializing in task reflection and evaluation. Your role is to:			1440
1374					1441
1375	3	1. Analyze all gameplay screenshots and state changes to understand what happened during task execution			1442
1376					1443
1377	4	2. Evaluate task completion status with concrete evidence			1444
1378					1445
1379					1446
1380					1447



1451	45	- Problems with action strategy or logic	<b>Action Design Reflection</b>	1520
1452				
1453	46		1 action_sys_prompt=''	1521
1454	47	c) Action Execution Level	2 You are an expert game action analyst	1522
1455	48	- Problems with specific control inputs	specializing in analyzing and	1523
1456			improving game control	1524
1457	49	- **Issues with duration of actions**	implementations. Your role is to:	1525
1458			3 1. Analyze gameplay screenshots to	1526
1459	50		understand the execution effects of	1527
1460	51	4. NEXT STEP RECOMMENDATION	each action	1528
1461	52	If task failed:	4 2. Evaluate action code design and	1529
1462	53	- Specific suggestions to complete the task in the **CURRENT** state	implementation quality	1530
1463			5 3. Provide reusable insights for similar	1531
1464			actions in the future	1532
1465	54		6 4. Suggest specific improvements for	1533
1466	55	If task succeeded:	action code design	1534
1467	56	- Simply state that the task was completed successfully and no modifications are needed	7	1535
1468			8 Always provide detailed, objective	1536
1469			analysis following the exact format	1537
1470	57		requested in the prompt.	1538
1471	58	Please provide your analysis in the following format:	9 ''	1539
1472			10	1540
1473	59	VISUAL ANALYSIS:	11 def generate_action_level_prompt(	1541
1474	60	<Describe the sequence of events observed in gameplay screenshots, including UI changes (if in UI screens), character actions, and any significant state changes>	pass_task_info, pass_action_code):	1542
1475			base_prompt = f"""	1543
1476			Analyze the	1544
1477			previous action execution using	1545
1478			the following information:	1546
1479	61	{} <Describe any relevant changes observed between initial and final maps>" if has_map else ""}	13	1547
1480			14 1. Screenshot Sequence Rules:	1548
1481			- For WASD movement actions lasting	1549
1482	62		over 2 seconds:	1550
1483	63	TASK COMPLETION EVALUATION:	* Screenshots are captured every	1551
1484	64	Status: <SUCCESS/FAILURE>	2 seconds during the	1552
1485	65	Evidence: <List specific observations from screenshots or state changes that support your status determination>	movement	1553
1486			- For all other key/mouse actions:	1554
1487			* Only two screenshots are	1555
1488			captured: one before and one	1556
1489	66		after the action	1557
1490	67	ISSUE ANALYSIS:	This helps track continuous	1558
1491	68	Task Design Level:	movements and precise action	1559
1492	69	<Evaluate if there are any issues with how the task was designed and specified. If no issues, explicitly state that>	effects.	1560
1493			20	1561
1494			21 2. Task Context:	1562
1495			{pass_task_info}	1563
1496	70		22	1564
1497	71	Action Planning Level:	23 3. Action plan and code list:	1565
1498	72	<Analyze if the planned sequence of actions was appropriate and complete . Identify any logical gaps or problems>	{pass_action_code}	1566
1499			24	1567
1500			25 Please conduct your analysis in	1568
1501			these sequential steps and	1569
1502	73		provide a detailed response in	1570
1503	74	Action Execution Level:	the following format:	1571
1504	75	<Assess if there were any issues with the specific implementation of actions, such as timing or input problems>	28	1572
1505			29 1. ACTION EXECUTION ANALYSIS	1573
1506			30 For each action in the sequence,	1574
1507			analyze:	1575
1508	76		31 - Initial state and final state from	1576
1509	77	NEXT STEP RECOMMENDATION:	screenshots	1577
1510	78	<If task failed: Provide specific suggestions for task completion given the current state>	- Whether the action achieved its	1578
1511			intended effect	1579
1512			- Timing and smoothness of execution	1580
1513	79	<If task succeeded: Simply state that the task was completed successfully and no modifications are needed>	- Any unexpected behaviors or side	1581
1514			effects	1582
1515			35 2. ACTION CODE EVALUATION	1583
1516	80		36 For each action implementation,	1584
1517	81	"""	37 evaluate:	1585
1518	82		38 - Appropriateness of key/mouse	1586
1519	83	return base_prompt	mapping choices	1587
			39 - Timing duration settings	1588

1589	40	- Action sequence coordination	analysis of what	1659
1590	41	- Code efficiency and reliability	worked/didn't work	1660
1591	42		>",	1661
1592	43	3. SUCCESS/FAILURE ANALYSIS	"reusability": {{	1662
1593	44	For each action, determine:	"	1663
1594	45	- Whether it succeeded or failed	applicable_scenarios	1664
1595	46	- Root causes of any failures:	": "<list of	1665
1596	47	a) Input mapping issues	potential reuse	1666
1597	48	b) Timing problems	cases>",	1667
1598	49	c) Sequence coordination issues	"prerequisites": "<	1668
1599	50	d) Environmental factors	required	1669
1600	51		conditions>",	1670
1601	52	4. REUSABILITY ANALYSIS	"limitations": "<	1671
1602	53	Analyze each action's potential for	known	1672
1603		reuse:	constraints>"	1673
1604	54	- Common scenarios where this action	}},	1674
1605		pattern could apply	"improvements": "<	1675
1606	55	- Required prerequisites and	specific suggestions	1676
1607		conditions	for implementation	1677
1608	56	- Potential adaptations needed for	improvements>"	1678
1609		different contexts	}}	1679
1610	57	- Limitations and constraints	}},	1680
1611	58		# ... repeat for each action	1681
1612	59	5. IMPROVEMENT RECOMMENDATIONS	]	1682
1613	60	Provide specific suggestions for:	```	1683
1614	61	- Better key/mouse mapping choices		1684
1615	62	- Optimal timing parameters	Ensure your response ends with this	1685
1616	63	- Enhanced sequence coordination	structured list for easy parsing	1686
1617	64	- More robust implementation	. Format it exactly as shown	1687
1618		patterns	above.	1688
1619	65		"""	1689
1620	66	Note that:	return base_prompt	1690
1621	67	1. output will be directly evaluated		
1622		using Python eval(), so it must be a		
1623		valid Python list of dicts	<b>Task Planning</b>	1691
1624	68	2. No additional text or explanation	1 task_planner_sys='''You are an	1692
1625		should be added between or after	intelligent game AI assistant	1693
1626		these sections	specializing in strategic task	1694
1627	69	After completing your analysis,	planning and execution.	1695
1628		output a list of dictionaries in		1696
1629		the following format:	2 Key Responsibilities:	1697
1630	70		3 1. Analyze game situations	1698
1631	71	```python	comprehensively considering:	1699
1632	72	[	4 - Current state and environment	1700
1633	73	{{	5 - Historical context and past	1701
1634	74	"action_name_description":	6 experiences	1702
1635		"<original action	7 - Game objectives and constraints	1703
1636		description from	8	1704
1637		action_name_description	9 2. For ALL tasks (not just movement),	1705
1638		>",	provide:	1706
1639	75	"action_code": "<	10 - Clear, specific, and actionable	1707
1640		corresponding action	objectives	1708
1641		code tuple from	11 - Precise success criteria	1709
1642		action_code>",	12 - Required resources or conditions	1710
1643	76	"reflection": {{	13 - Risk assessment and mitigation	1711
1644	77	"execution_analysis": "<	strategies	1712
1645		summary of execution	14	1713
1646		analysis>",	15 3. For movement-related tasks, MUST	1714
1647	78	"code_evaluation": {{	provide precise location	1715
1648	79	"status": "<SUCCESS/	descriptions using:	1716
1649		PARTIAL SUCCESS/	16 - Relative position to character (	1717
1650		FAILURE>",	using character height as scale)	1718
1651	80	"quality_analysis":	17 - Directional instructions (up/down/	1719
1652		"<implementation	left/right or compass directions)	1720
1653		quality summary	18 - Safe path recommendations	1721
1654		>"	considering terrain	1722
1655	81	}},	19	1723
1656	82	"	20 4. Special Considerations:	1724
1657		success_failure_analysis	21 - Prioritize agent safety and	1725
1658		": "<detailed	objective completion	1726

1727	22	- Balance exploration with risk management	68	Consider:	1797
1728			69	1. Previous task outcomes and lessons learned	1798
1729	23	- Adapt strategy based on previous task outcomes	70	2. Current environmental constraints	1799
1730			71	3. Progress toward game objectives	1800
1731	24	- Consider resource management and efficiency	72	4. Safety and risk management"""	1801
1732			73		1802
1733	25		74	output_format = ""	1803
1734	26	Your task is to make informed decisions that progress game objectives while maintaining agent safety and efficiency.	75	Based on your analysis, provide your response in the following format:	1804
1735			76		1805
1736			77	reasoning process:	1806
1737			78	1. Current State Analysis: "<analyze current environment and immediate situation>"	1807
1738	27	'''	79	2. Historical Context: "<analyze relevant history and reflections>"	1808
1739	28		80	3. Strategic Evaluation: "<evaluate opportunities, risks, and priorities>"	1809
1740	29	def construct_task_prompt(current_frame, pass_task_history_summary, pass_task_reflection, env_info, game_info, step):	81		1810
1741			82	task details:	1811
1742			83	goal: "<specific, actionable objective>"	1812
1743			84		1813
1744	30		85	location details:	1814
1745	31	base_prompt = f"""	86	- screen_position: "<describe target position. Example: '3 meters to the right'>"	1815
1746	32	Analyze the current situation and plan the most appropriate next task considering:	87	key_requirements: "<essential conditions or resources needed>"	1816
1747			88	success_criteria: "<main condition that must be met>"	1817
1748			89		1818
1749	33		90	Note:	1819
1750	34	1. Game Objectives:	91	- The direction of camera adjustment ** MUST** be consistent, and there should be no angle that switches left and then right, or up and then down	1820
1751	35	{game_info.get('Global_task')}	92	- For movement-related tasks, always specify both screen-relative positions (using character height as scale). For non-movement tasks, mark position fields as 'N/A' if not relevant. """	1821
1752	36	2. Additional Task Context: {game_info.get('additional_task_info4_task_plan')}	93		1822
1753			94		1823
1754			95	return analysis_prompt + output_format	1824
1755	37				1825
1756	38	Your design task should be broken down into the following specific Available Controls:			1826
1757					1827
1758					1828
1759	39	{game_info.get('control_info')}			1829
1760	40				1830
1761	41	Required Analysis Steps:			1831
1762	42	1. Evaluate current environment and state			1832
1763					1833
1764	43	2. Consider historical context and lessons learned			1834
1765					1835
1766	44	3. Assess risks and opportunities			1836
1767	45	4. Determine priority actions"""			1837
1768	46				1838
1769	47				1839
1770	48	current_state = f"""			1840
1771	49	Current Environment Status:			1841
1772	50	{env_info}"""			1842
1773	51				1843
1774	52	if step == 1:			1844
1775	53	analysis_prompt = f"""{base_prompt}			1845
1776					1846
1777	54	{current_state}			1847
1778	55				
1779	56	This is the initial step. Focus on understanding the current situation and establishing a safe starting point. """			1848
1780					1849
1781					1850
1782					1851
1783	57				1852
1784	58	else:			1853
1785	59	history_context = f"""			1854
1786	60	Historical Context:			1855
1787	61	Task History Summary: {pass_task_history_summary}			1856
1788					1857
1789	62				1858
1790	63	Previous Task Reflection: {pass_task_reflection}"""			1859
1791					1860
1792	64	analysis_prompt = f"""{base_prompt}			1861
1793					1862
1794	65	{current_state}			1863
1795	66	{history_context}			1864
1796	67				1865
				<b>Action Planning</b>	
			1	action_prompt = f"""Based on the task you just planned, break it down into specific executable actions.	
			2	Please list the specific actions needed to complete this task.	
			3		
			4	Available Controls:	
			5	{game_info.get('control_info')}	
			6		
			7	Note that:	
			8	1. output will be directly evaluated using Python eval(), so it must be a valid Python list	
			9	2. No additional text or explanation should be added between or after these sections	
			10	3. Ignore actions such as ' wait 'and'	

1866	observe' that cannot be associated	36	3. Ensure precise timing for each	1935
1867	with available controls		control input	1936
1868	11 4. Action list is *no longer than 5!!*.	37	4. Consider safety in all actions	1937
1869	12	38		1938
1870	13 Output Format MUST be exactly as follows	39	Output Format MUST be exactly as follows	1939
1871	14 :	40	:	1940
1872	14 ["Action1: <action name> - <detailed	40	[	1941
1873	description including precise	41	{	1942
1874	measurements and requirements>","	42	"action_name_description": "<	1943
1875	Action2: <action name> - <detailed		original action description	1944
1876	description including precise		>","	1945
1877	measurements and requirements>","	43	"action_code": [("<key>", <	1946
1878	...]		duration>), ...]	1947
1879	15	44	}},	1948
1880	16 """	45	...	1949
1881	<b>Action Construction</b>	46	]	1950
1882	1 action_sys_prompt = '''	47		1951
1883	2 You are an expert game AI action planner	48	Example Output:	1952
1884	specializing in converting high-	49	[	1953
1885	level tasks into precise, executable	50	{	1954
1886	action sequences.	51	"action_name_description": "Move	1955
1887	3		Forward - Move 3 meters	1956
1888	4 Key Responsibilities:	52	forward",	1957
1889	5 1. Convert task descriptions into	53	"action_code": [("W", 3.0)]	1958
1890	specific control sequences	54	}},	1959
1891	6 2. Ensure accurate timing and duration	55	{	1960
1892	for each action		"action_name_description": "Jump	1961
1893	7 3. Maintain action safety and efficiency		and Interact - Jump over	1962
1894	8 4. Generate properly formatted action	56	obstacle and press button",	1963
1895	code that can be directly executed		"action_code": [("SPACE", 0.1),	1964
1896	9	57	("E", 0.1)]	1965
1897	10 Important Guidelines:	58	}}	1966
1898	11 1. All outputs must be in valid Python	59	]	1967
1899	dictionary list format	60	Note:	1968
1900	12 2. Each action must include both	61	1. Output will be evaluated using Python	1969
1901	description and corresponding		ast.literal_eval()	1970
1902	control code	62	2. Use only valid control keys: {list(	1971
1903	13 3. Control codes must use only valid		game_info.get('Mapping_info', {}).keys())}	1972
1904	game controls		3. All durations must be positive	1973
1905	14 4. All durations must be reasonable and	63	numbers	1974
1906	safe	64	4. Maintain exact format with no	1975
1907	15 '''		additional text	1976
1908	16	65	"""	1977
1909	17 def generate_action_prompt(game_info,	66	return prompt	1978
1910	reason_task, action_plan):			1979
1911	prompt = f"""			1980
1912	19 You are an expert game AI action planner			
1913	specializing in converting high-			
1914	level action into precise,			
1915	executable action sequences.			
1916	20			
1917	21 Your current task:			
1918	{reason_task}			
1919	22			
1920	23			
1921	24 The action plan for task:			
1922	{action_plan}			
1923	25			
1924	26			
1925	27 Available Controls:			
1926	{game_info.get('control_info')}			
1927	28			
1928	29			
1929	30 Additional Action Information:			
1930	{game_info.get('additional_action_info')}			
1931	}			
1932	31			
1933	32			
1934	33 Requirements:			
	34 1. Convert each action into specific			
	control sequences			
	35 2. Provide both action description and			
	control code			