

# The Dark Side of VLM Rewards: Understanding and Mitigating False Positive Noise

Anonymous Author(s)

Affiliation

Address

email

1           **Abstract:** Vision-Language Models (VLMs) have shown promise as reward mod-  
2       els for reinforcement learning (RL) agents, potentially providing rich reward sig-  
3       nals based on how well an agent’s trajectory aligns with expert instructions. How-  
4       ever, our observations reveal that RL agents trained with VLM rewards often un-  
5       derperform compared to those employing only intrinsic (exploration-driven) re-  
6       wards, contradicting expectations set by recent work. To understand this unex-  
7       pected outcome, we analyze the noisy reward signals from VLMs, particularly  
8       focusing on the prevalence and consequences of false positives and false nega-  
9       tives. Our analysis revealed that false positive rewards of VLMs were a major  
10      factor contributing to the poor learning efficacy. To address this, we developed  
11      a novel reward function, BiMI, that prioritizes minimizing false positive rewards  
12      while accepting a slight increase in false negatives as a necessary trade-off. Experi-  
13      ments demonstrate that BiMI substantially improves agent performance across  
14      challenging environments, highlighting the critical importance of addressing re-  
15      ward noise when applying VLM-based rewards to complex real-world scenarios.

16           **Keywords:** RL, VLM-based Reward Model, Sparse Reward Environments

## 17      1 Introduction

18      Natural language instructions are increasingly recognized as a valuable source of reward signals for  
19      reinforcement learning (RL) agents to learn complex tasks. In particular, a growing trend in agent  
20      learning involves using vision-language models (VLMs) for reward modeling. This approach mea-  
21      sures the semantic similarity – often quantified by cosine similarity – between the embedding rep-  
22      resentations of an agent’s behaviors (i.e., past trajectories) and the provided instructions, all within  
23      the same embedding space [1, 2, 3, 4].

24      However, we observed that RL agents trained with VLM reward models, while effective in simpli-  
25      fied settings, often struggled with tasks involving complex dynamics and longer action horizons.  
26      This is evident in several recent works – for instance, Goyal et al. [2] reported the effective use  
27      of VLM rewards in Montezuma’s Revenge, a notoriously challenging Atari game. However, we  
28      observed that this success was confined to individual sub-tasks and the agent struggled when at-  
29      tempting to scale up to the full game. Similarly, Du et al. [4] demonstrated impressive performance  
30      of VLM rewards in guiding agents within a 2D survival game. However, their study was conducted  
31      in a modified environment with a reduced observation and action space using internal game state  
32      information and manually defined macro actions. Consequently, when tested in the original, un-  
33      modified environment, their agent’s performance did not exceed that of agents using only *intrinsic*  
34      (exploration-driven) rewards.

35      The consistent underperformance of VLM-based models in complex environments, **particularly**  
36      **their unexpected failure to outperform even simple exploration-driven rewards**, raised sig-  
37      nificant concerns about their reliability in real-world applications. This discrepancy, where VLM

38 reward models underperformed contrary to their perceived potential, prompted us to investigate the  
 39 underlying causes of this performance gap. Our findings indicate that **noisy reward estimates** in  
 40 VLMs are a key factor contributing to poor learning efficacy. We attribute this noise primarily to  
 41 the approximation errors inherent in the *cosine similarity metric* commonly used in these models.  
 42 Specifically, our analysis focuses on two types of noise: false positives, which involve rewarding  
 43 unintended behaviors, and false negatives, which occur when correct behaviors are not rewarded.

44 We posit that false positive reward noise (see Figure 1)  
 45 are not only more prevalent but potentially more detrimental to the learning process than false negatives, a  
 46 hypothesis supported by our empirical and theoretical  
 47 findings. To this end, we propose a novel reward function,  
 48 **BIMI** (**B**inary **M**utual **I**nformation **R**eward). It  
 49 leverages binary reward signals to directly reduce the  
 50 occurrence of false positives and incorporates a mutual  
 51 information term to prevent overfitting to potentially  
 52 noisy signal sources. Our experiments demonstrate  
 53 that BIMI significantly improves the learning efficacy  
 54 of agents trained by VLM rewards across various chal-  
 55 lenging environments.

## 57 2 Related Work

58 **Use of VLMs as Reward Models.** VLMs have been  
 59 pivotal in robotics studies, serving as reward models  
 60 that guide agents to follow instructions [5, 6, 7]. While  
 61 most research primarily focuses on leveraging VLMs  
 62 to overcome the challenge of manual reward design  
 63 for complex tasks [8], the impact of reward noise and  
 64 its implications for policy convergence rates are often  
 65 overlooked. As mentioned in Section 1, some work  
 66 sidesteps the noisy reward problem by accessing inter-  
 67 internal state information from the game engine as well as  
 68 providing predefined action macros [4, 9], thereby pre-  
 69 venting the accumulation of noise over longer horizon.  
 70 However, understanding the impact of reward noise  
 71 from VLMs is crucial for developing reliable language  
 72 interface for RL agents, as it directly affects their ability to learn effectively in real-world tasks.

73 **Mitigating Reward Model Failures.** Research by Ghosal et al. [10] and Fu et al. [11] has intro-  
 74 duced methods to counteract unreliable reward signals from learned VLMs. These strategies involve  
 75 employing a parallel exploration-based policy alongside the reward-maximizing policy, thereby re-  
 76 ducing reliance on potentially misspecified VLM rewards. Our work contributes to this growing  
 77 body of research by proposing a novel reward function that directly addresses the issue of noisy  
 78 rewards in VLM-based models, complementing approaches that use exploration policies to escape  
 79 local optima. Furthermore, we have tested the synergy of combining our reward function with explo-  
 80 ration strategies, demonstrating how these approaches can be integrated for improved performance.

## 81 3 Formal Problem Statement

82 In the context of using instructions as auxiliary guide for RL agents, we frame our task as a MDP  
 83 defined by a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, s_0, r^e, \gamma \rangle$ , where  $\mathcal{S}$  represents a set of states  $s \in \mathcal{S}$ ,  $\mathcal{A}$  represents a set  
 84 of actions  $a \in \mathcal{A}$ , and  $\mathcal{T}(s' | s, a)$  describes the dynamics of the environment.  $s_0 \in \mathcal{S}$  is the initial  
 85 state, and  $\gamma \in [0, 1]$  is the reward discount factor.  $r^e(s, a, s')$  is the environmental reward function.

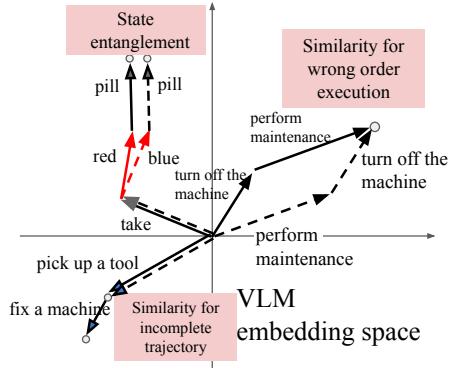


Figure 1: Schematic diagram of false positives in a VLM embedding space. The unintended agent’s trajectory (dashed line) may exhibit high cosine similarity to the instruction (solid line) in the embedding space, as indicated by the proximity of their endpoints in the embedding space. Despite this apparent similarity, the unintended trajectories fundamentally fail to fulfill the intended instruction. Therefore, the rewards guides the agent towards incorrect behaviors. The figure shows three distinct cases of false positives. Refer to Section 4 for more details.

86 In this work, we focus on a sparse reward setting, where the agent receives a non-zero reward only  
 87 when reaching goal states  $S_G \subset \mathcal{S}$ , and 0 otherwise, with  $|S_G| \ll |\mathcal{S}|$ . An agent's trajectory is a  
 88 sequence of states and actions  $\tau = \langle s_t, a_t, \dots, s_{t+T} \rangle$ , where  $t$  is the initial time step and  $T$  is the  
 89 length of the trajectory. The objective is to learn a policy  $\pi$  that maximizes the expected cumulative  
 90 reward of the trajectories  $\tau^\pi = \{\tau_1, \dots, \tau_n \mid \tau_i \text{ induced by } \pi\}$ .

91 A sparse reward function only *realizes* a set of acceptable policies, rather than distinguishing be-  
 92 tween them more finely (as derived from Abel et al. [12]; see Appendix A.1 for details). Building  
 93 on this, we define the following:

94 (1) *Policy Universe*,  $\Pi$ : the set of all possible policies for the MDP.

95 (2) *Acceptable Policies*,  $\Pi_G$ : the policies in  $\Pi_G$  are those policies whose start-state values  $V^\pi(s_0)$   
 96 are within a small  $\epsilon$  of the optimal value, while all other policies have strictly lower values.

97 (3) *Unacceptable Policies*,  $\Pi_B$ : formally,  $\Pi_B = \Pi \setminus \Pi_G$ , which refers to the set of policies that do  
 98 not meet the near-optimal criteria and are therefore considered not acceptable.

99 A task instruction  $L$  is a sequence of natural language sentences  $L = \{l_1, l_2, \dots, l_n\}$ , where each  $l_k$   
 100 induces a set of desired sub-trajectories  $\tau^{l_k}$ . A complete task execution is represented by a trajectory  
 101 that concatenates instructions' corresponding sub-trajectories in order, traversing from the initial  
 102 state  $s_0$  to a goal state in  $S_G$ .

103 An auxiliary VLM-based reward model provides more  
 104 frequent rewards by evaluating the semantic similarity  
 105 between the agent's trajectory and the specified in-  
 106 struction sentence. This can be formally represented as  
 107  $r^v(\tau, l_k) = p(l_k|\tau)$  (see Figure 2), where  $r^v$  denotes the  
 108 VLM-based reward function. The use of *non-Markovian*  
 109 reward functions in MDP environments has been well-  
 110 established in the literature, particularly through the  
 111 work on reward machines [13, 14]. Such a VLM-based  
 112 reward function consists of two components: an embed-  
 113 ding encoder  $\phi(\cdot)$  and a similarity function. The reward  
 114 is calculated as the cosine similarity between the embed-  
 115 dings, which can be expressed as  $\frac{\phi(\tau) \cdot \phi(l_k)}{\|\phi(\tau)\| \|\phi(l_k)\|}$ .

116 **Reduction of Convergence Rate.** To understand the  
 117 effects of reward noise on policy learning, it's crucial to  
 118 analyze its impact on the convergence rate. In sparse re-  
 119 ward settings, the convergence rate – defined as the number of training iterations needed to learn an  
 120  $\epsilon$ -optimal policy – can be significantly affected by false positive rewards that are incorrectly given  
 121 for policies that do not lead to the goal states.

122 To illustrate this, we analyze the convergence rate by measuring the probability of learning accept-  
 123 able policies ( $P(\pi \in \Pi_G)$ ) over each update iteration. The following theorem demonstrates how  
 124 false positive rewards impact the convergence rate:

125 **Theorem 3.1** (Convergence rate reduction). *In Actor-Critic algorithm, gradient ascent on*  
 126  $Q(s, a)\pi_{\theta_i}(a|s)$  *pushes the next updated policy*  $\pi_{\theta_{i+1}}$  *in the direction provided by the Q value func-*  
 127 *tion. In the direction of maximizing*  $P(\pi \in \Pi_G)$ , *the gradient of*  $Q(s, a)\pi(a|s)$  *can be expressed as*  
 128 *follows:*

$$\begin{aligned} \nabla_\theta Q_\phi(s, a)\pi_\theta(a|s) &= \text{const} \cdot \mathbb{E}[G^{\pi \in \Pi_B}] \nabla_\theta \pi_\theta(a|s) \\ &\quad + (\mathbb{E}[G^{\pi \in \Pi_G}] - \mathbb{E}[G^{\pi \in \Pi_B}]) P(\pi \in \Pi_G) \nabla_\theta \pi_\theta(a|s) \end{aligned}$$

129 where  $\Pi_G$  is the set of acceptable policies;  $\Pi_B$  is the set of unacceptable policies that follow in-  
 130 structions but fail to reach the goal;  $\mathbb{E}[G^\pi]$  is expectation of the cumulative rewards by executing  
 131 policy  $\pi$  in one episode.

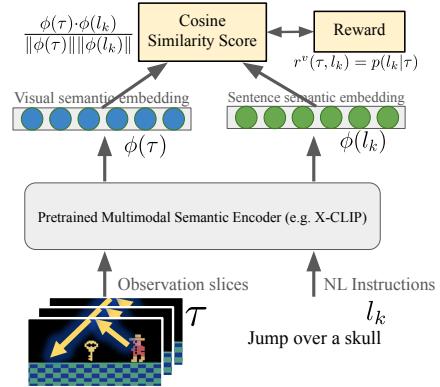


Figure 2: Illustration of reward signals from a vision-language reward model

132 The proof is in Appendix A.2. In this context, false positive rewards correspond to rewards assigned  
133 to policies within the set  $\Pi_B$ , which represents undesirable policies. Therefore,  $\mathbb{E}[G^{\pi \in \Pi_B}]$  can  
134 be interpreted as the expected cumulative false positive rewards. This theorem highlights how the  
135 presence of  $\mathbb{E}[G^{\pi \in \Pi_B}]$  reduce the convergence rate by: (1) introducing deviation in the policy update  
136 directions (the first term), and (2) reducing the gradient of the desired policy update (the second  
137 term). Both effects can significantly hinder the convergence process.

138 The issue of noisy rewards is not merely a theoretical concern; it has practical implications for  
139 the effectiveness of instruction-guided RL in real-world applications. In the next section, we will  
140 investigate how cosine similarity metric used in VLM-based reward models inevitably introduce  
141 false positive rewards. Additionally, we will empirically evaluate the impact of both false positives  
142 and false negatives on the learning efficacy of agents in sparse reward environments.

## 143 4 Approximation Error of Cosine Similarity

144 In this section, we will discuss two fundamental issues associated with cosine similarity scores in RL  
145 contexts: *state entanglement* and *composition insensitivity*. The former issue, state entanglement,  
146 refers to the metric’s inability to recognize trajectories  $\{\tau_1, \dots, \tau_n \mid \tau_i \text{ induced by } \pi \in \Pi_B\}$  that,  
147 while being cosine similar to the target instruction in the embedding space, fail to reach the goal  
148 states in  $S_G$ . The latter issue refers to the metric’s tendency to reward trajectories that execute  
149 the sub-tasks (corresponding to individual instruction sentences  $l_k$ ) in an incorrect order. This is  
150 problematic in RL tasks where the specific sequence of executing sub-tasks, as defined by the order  
151 of instructions in  $L = \{l_1, l_2, \dots, l_n\}$ , is crucial for reaching the goal states.

152 **The Issue of State Entanglement** State entanglement refers to the issue where the cosine similarity  
153 metric erroneously pays more attention to lexical-level similarity while lacking comprehension  
154 of the underlying state transitions. Consequently, rewards are given to trajectory-instruction pairs  
155 that are cosine similar in embedding space but in fact result in distinct state transitions. For in-  
156 stance, consider the significant contrast between “take the *red* pill” and “take the *blue* pill”. Despite  
157 their lexical similarity, they lead to vastly different states. However, the cosine similarity metric  
158 may represent them as similar due to the shared action “take” and shared object “pill”, disre-  
159 garding the critical difference in state outcomes. Understanding state transitions is crucial in sequential  
160 decision-making scenarios. Otherwise, rewards may be given to trajectories that lead to unintended  
161 states, resulting in the learning of unacceptable policies.

162 **The Issue of Composition Insensitivity** Composition insensitivity in cosine similarity metrics  
163 gives rise to two issues: (1) *rewarding incomplete task execution* – cosine similarity may incorrectly  
164 reward partial task completion, as even incomplete trajectories can receive high similarity score in  
165 the embedding space. For instance, in a task to “pick up a tool, then fix a machine,” the model  
166 might prematurely reward the agent for merely picking up the tool, neglecting the crucial repair  
167 action. This can lead to suboptimal learning, with the agent potentially getting stuck in easier,  
168 incomplete subtasks rather than fully accomplishing the task. (2) *insensitivity to the ordering of*  
169 *execution* – cosine similarity often fails to adequately penalize incorrect execution sequences. In a  
170 safety protocol requiring an agent to “turn off the machinery, then perform maintenance,” the metric  
171 might assign high rewards based merely on the presence of relevant actions, disregarding their order.  
172 While some large language models have become sensitive to the order of tokens [15], the compact  
173 visual and sentence embeddings from multimodal VLMs remains largely insensitive to sequential  
174 information [16]. This limitation can lead to potentially dangerous outcomes in scenarios where  
175 execution sequence is critical, such as in safety protocols or complex manufacturing processes.

176 Figure 1 illustrates various scenarios where high similarity scores are erroneously assigned to false  
177 positive cases. To empirically demonstrate the issue, Section 4.1 presents experiments on false  
178 positive rewards and their impact on agent learning in sparse reward environments.

## 179 4.1 Experiments on Reward Noise Impact

180 Our experiments test the following hypothesis: **(H1)** The two issues of *state entanglement* and *composition insensitivity* exist; **(H2)** *false positive* rewards are prevalent during training; **(H3)** similarity-based reward models lacking noise handling mechanisms underperform against strong intrinsic reward models in sparse reward environments; **(H4)** *false negatives* may not be as harmful as *false positives*.

185 **Setup.** We evaluate these hypotheses through various challenging sparse-reward environments: 186 (1) *Crafter*, an open-ended 2D Minecraft environment [17]; (2) *Montezuma’s Revenge*, a classic 187 hard adventure game in Atari [18]; and (3) *Minigrid ‘Go To Seq’*, a hard task involving long-horizon 188 navigation and object interactions [19]. Two cosine similarity-based reward models were tested: 189 (1) *Pix2R* by Goyal et al. [3], which uses only the current video frame to determine if the goal 190 state specified in the instruction has been reached; and (2) *ELLM-*, a variant of *ELLM* by Du et al. 191 [4]. Unlike *ELLM*, which queries instructions from LLMs in real-time, *ELLM-* directly uses pre-set 192 expert instructions and compares them with the transition differences of the agent’s trajectory. 193 The VLM backbones used are: (1) *CLIP* [20], pretrained by image-text pairs; and (2) *X-CLIP* [21], 194 pretrained by video-text pairs. To ensure high-quality finetuning data, we used internal information 195 from the game engine to annotate expert trajectories from expert agents. To demonstrate how noisy 196 reward signals hinder learning, we selected a strong intrinsic reward model *DEIR* [22] for com- 197 parison. It provides auxiliary rewards based on observation novelty to encourage exploration. See 198 Appendix A.3 for detailed implementation of the experiments.

199 **Evaluation Metric.** We used the *score* metric adapted from the Crafter benchmark [17] for perfor- 200 mance evaluation, as it better captures an agent’s consistent performance across multiple subtasks. 201 This metric follows the intuitive ‘higher is better’ principle. Unlike *maximum total rewards* metric, 202 which fail to capture consistent performance in sparse reward settings, the score metric provides a 203 more reliable measure of learning progress.

204 **Reward Noise Issue.** To investigate **H1**, we evaluated the models’ sensitivity by examining how 205 cosine similarity scores change for manipulated trajectory-instruction pairs. The state entanglement 206 test involved reversing trajectories and negating instructions (i.e., “do not do  $l_k$ ”). The composition 207 insensitivity test examined concatenated pairs of trajectory-instruction data. For example, given 208  $(\tau_1, l_1)$  and  $(\tau_2, l_2)$ , we create a concatenated pair  $(\tau_1 + \tau_2, l_1 + l_2)$ . We then test two types of 209 manipulations – (1) swapping the order within one modality: e.g.,  $(\tau_2 + \tau_1, l_1 + l_2)$ ; and (2) truncating 210 one modality: e.g.,  $(\tau_1, l_1 + l_2)$ . Overall, we categorized the trajectory-instruction pairs for evalua- 211 tion into three types: (1) matched pairs – these are positive examples where the trajectory correctly 212 corresponds to the given instruction. (2) not-matched pairs – these are negative examples where the 213 trajectory does not correspond to the given instruction. (3) manipulated pairs: derived from matched 214 pairs by altering either the trajectory or the instruction. Ideally, manipulated pairs should receive low

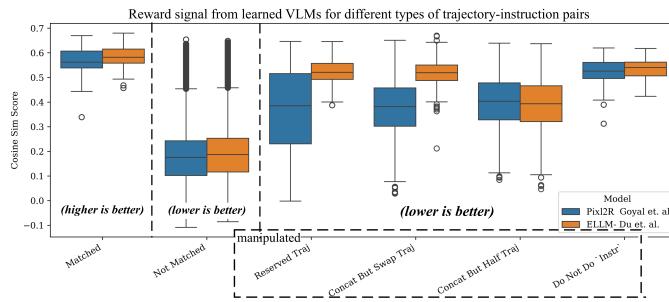


Figure 3: Learned VLM models differentiate between matched and not-matched pairs, but struggle with O.O.D. cases. They incorrectly assign high scores to manipulated pairs, which should be low as the trajectories in the manipulated pairs fail the instruction

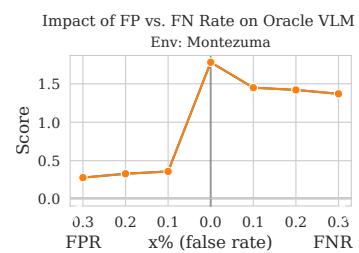


Figure 4: The false positive vs. false negative oracle model. The false positive model get a more severe drop in the final training score.

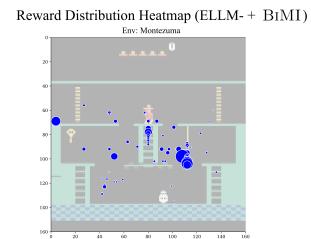
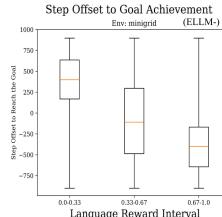
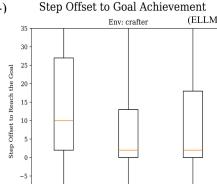
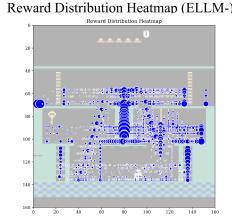


Figure 5: The heatmap shows the cumulative rewards received at various locations, with larger circle sizes indicating higher rewards. The later figures shows the offsets between the state where rewards are given and the actual goal-reaching state. Agents are getting both issues of false positives and false negatives during training

Figure 6: The ratio of false positive rewards is significantly reduced after applying BiMI

similarity scores as essentially the trajectory are not fulfilling the instruction. However, our results reveal that the reward model assigns high scores to these manipulated pairs, particularly in the “do  $l_k$ ” case. This finding highlights the noise issue in cosine similarity-based reward models (see Figure 3). It’s worth noting that the poor performance in the negation case aligns with broader challenges in natural language processing. Recent studies [23, 24] have highlighted that negation is central to language understanding but is not properly captured by modern language models. This limitation extends to VLMs and significantly impacts their ability to provide accurate rewards in complex scenarios and instructions.

**Prevalence of False Positives.** To address **H2**, we analyzed reward distribution heatmap from VLM-based reward models during training. The heatmap revealed a concerning trend: RL agents engage in reward hacking, receiving rewards across vast areas of the environment rather than just at goal states. For instance, in *Montezuma* environment where the goal is to grab the key and escape the room, we observed that agents received rewards even for falling off cliffs, which undoubtedly contribute to the detrimental  $\mathbb{E}[G^{\pi \in \Pi_B}]$  specified in Theorem 3.1. For environments without fixed camera views, we calculated the step offset between the current rewarded state and the actual goal state. A positive offset indicates that a reward was given before reaching the goal state, signifying a false positive reward, while a negative offset suggests that the agent reached the goal state but the reward model failed to recognize it, indicating a false negative reward (see Figure 5). Interestingly, we observed a large amount of negative step offsets in Minigrid environments. We attribute this to Minigrid’s abstract shape-based visual representations, which fall outside the VLM’s pretraining distribution.

**Impact on Learning.** We trained agents using learned VLM reward models and compared their learning efficacy against intrinsic reward models. As shown in Table 7, our results confirmed **H3**: *instruction-guided RL agents using learned VLM reward models without noise handling consistently underperform compared to DEIR, the intrinsic reward-based RL agent*. To investigate the impact of false negatives versus false positives (**H4**), we designed an oracle Pixl2R model with two variants: a false negative model and a false positive model. The false negative model only rewards the agent for reaching subgoal states described in the instruction, with a probability of  $x\%$  that some rewarding states in the map are removed. In contrast, the false positive model rewards the agent for reaching every subgoal, but also introduces a 0.1 one-off reward for certain locations, covering  $x\%$  of the map. The results indicate that false negatives may be less detrimental to agent performance than false positives (see Figure 4). This performance difference can be explained through our theoretical framework:

Figure 7: Score metric across environments (equivalent to total rewards, higher is better). \* denotes baseline intrinsic reward model. VLM reward models without noise handling underperformed. All models are based on PPO.

Models	Type	Monte.	Minigrid	Crafter	% vs. DEIR
PPO	Pure	0.151	24.993	16.863	-28%
DEIR *	Intrinsic	0.174	55.556	19.758	-
Pixl2R	VLM	0.142	12.422	9.409	-49%
ELLM-	VLM	0.150	19.406	10.826	-41%
Pixl2R + DEIR	VLM + intr.	0.176	17.372	10.440	-38%
ELLM- + DEIR	VLM + intr.	0.178	30.985	11.857	-27%

254 (1) False negative models maintain the coefficient  $\mathbb{E}[G^{\pi \in \Pi_G}] - \mathbb{E}[G^{\pi \in \Pi_B}]$ , ensuring steady gra-  
 255 dient ascent towards maximizing  $P(\pi \in \Pi_G)$ ;  
 256 (2) By decreasing  $\mathbb{E}[G^{\pi \in \Pi_B}]$ , false negative models minimize deviations from the target direction,  
 257 leading to more stable learning.  
 258 In contrast, scenarios with high  $\mathbb{E}[G^{\pi \in \Pi_B}]$  (i.e., the false positive case) significantly reduce the  
 259 gradient ascent rate and also introduce deviations from the target direction. The findings thus chal-  
 260 lenged the prevailing assumption that false negatives are the primary issue in the degeneracy of  
 261 instruction-guided RL. See Appendix A.4 for more experimental details.

## 262 5 Solution to the Reward Noise Issue

### 263 5.1 Binary Signal and Conformal Prediction Thresholding

264 Our experiments have shown that the issue of false positives may be more detrimental to learn-  
 265 ing than false negatives. Based on these findings, we propose a reward function that provides  
 266 one-time binary rewards only when the similarity between the agent’s current trajectory and the  
 267 instruction exceeds a high confidence threshold. This method effectively reduces the likelihood of  
 268 rewarding unacceptable trajectories. To achieve this, we introduce a thresholding mechanism using  
 269 a calibration set of true positive trajectory-instruction pairs. This threshold, denoted as  $\hat{q}$ , is set to  
 270 the empirical quantile of cosine similarity scores at the significance level  $1 - \alpha$ . Pairs whose sim-  
 271 ilarity scores fall below this threshold  $\hat{q}$  receive no reward. Conversely, pairs exceeding  $\hat{q}$  receive  
 272 a one-time +1 reward, i.e.,  $r_{\text{BI}}^v(\tau, l_k) = \mathbf{1}_{\{p(l_k | \tau) \geq \hat{q}\}}$ . This approach statistically guarantees a high  
 273 probability (at least  $1 - \alpha$ ) that true positive pairs are recognized while minimizing frequency of  
 274 false positives errors [25]. See Appendix A.5 for detailed threshold calculation.

### 275 5.2 Mutual Information Maximization

276 Intuitively, when we observe rewards coming from a particular signal source too frequently, we  
 277 tend to downplay the significance of that signal to avoid over-reliance. This intuition is effectively  
 278 captured by incorporating a *mutual information maximization* term into the reward function. Specif-  
 279 ically, the updated reward function  $r_{\text{MI}}^v(\tau, l_k)$  measures the mutual information between the agent’s  
 280 trajectory and the instruction. Mathematically, it can be expressed as:

$$\begin{aligned}
 r_{\text{MI}}^v(\tau, l_k) &= I(l_k; \tau) = D_{KL}(p(l_k, \tau) || p(l_k)p(\tau)) \\
 &= \mathbb{E}_{\tau \sim \pi_\theta, l_k \sim L} [\log p(l_k | \tau) - \log p(l_k)]
 \end{aligned} \tag{1}$$

281 where  $p(l_k | \tau)$  comes from the similarity score provided by VLMs, and  $p(l_k)$  is the frequency  
 282 of the instruction  $l_k$  being fulfilled by the agent’s trajectory. Therefore, the second term in the  
 283 equation serves as a regularization term that downplays the significance of the reward signal when  
 284 it is too frequent. For instance, if a VLM frequently detects that the agent’s actions are fulfilling  
 285 the “climbing the ladder” instruction, even when the agent is performing unrelated tasks, the reward  
 286 signal for this instruction will be downplayed.  $p(l_k)$  is calculated as follows:

$$p(l_k) = \mathbb{E}_{\tau \sim \pi_{\theta-1}} \left[ \frac{\sum_{t=1}^{T_\tau} \mathbf{1}_{\{p(l_k | \tau_t) \geq \hat{q}\}}}{T_\tau} \right] \tag{2}$$

287 Here,  $\tau_t$  is the agent’s trajectory up to time  $t$ , and if the VLM deems the trajectory as fulfilling  
 288 the instruction (i.e.,  $p(l_k | \tau_t) \geq \hat{q}$ ), we increment the count. Dividing the count by the total  
 289 trajectory length  $T_\tau$  gives the empirical frequency of the instruction being fulfilled. The subscript  
 290  $\theta-1$  in  $\pi_{\theta-1}$  indicates that the trajectories are sourced from rollouts in the previous policy iteration,  
 291 acknowledging the impracticality of real-time computation of  $p(l_k)$  during an ongoing episode.

292 To enhance the stability of the training process, we adopt a linearized version of the mutual infor-  
 293 mation maximization approach, as proposed by Li et al. [26]. Overall, BiMI, the proposed reward

294 function that enhances the noise resilience of VLM-based reward models, can be expressed as fol-  
 295 lows:

$$r_{\text{BiMI}}^v(\tau, l_k) = \max(\mathbf{1}_{\{p(l_k|\tau) \geq \hat{q}\}} - p(l_k), 0) \quad (3)$$

296 It's important to note that the BiMI approach primarily mitigates false positives (FP) rather than  
 297 false negatives (FN). By implementing a high confidence threshold and a binary reward signal, we  
 298 effectively reduce the likelihood of rewarding incorrect trajectories, thus addressing the FP issue.  
 299 While this conservative approach may potentially increase FNs, as some correct trajectories might  
 300 fall below the threshold, our experiments indicate that this trade-off is beneficial. We found that false  
 301 positives are more detrimental to learning in instruction-guided RL tasks, justifying our preference  
 302 for a method that prioritizes FP reduction.

## 303 6 Experiments

304 We evaluate the BiMI reward function using baseline agents (*Pixl2R* and *ELLM-*) and their BiMI-  
 305 enhanced counterparts, while also exploring potential synergies with the intrinsic reward model  
 306 DEIR. We follow the same experimental setup as in Section 4.1, with additional details in Ap-  
 307 pendix A.6.

Table 1: Model score across various environments.  $\star$  is the baseline agents with a learned VLM-based reward model to compare with. BiMI significantly improves performance in *Montezuma* and *Minigrid*, while showing mixed results in *Crafter* due to task-specific characteristics

Methods	Montezuma	% vs. $\star$	Minigrid	% vs. $\star$	Crafter	% vs. $\star$
<i>Pixl2R</i> *	$0.142 \pm 0.003$	—	$12.422 \pm 2.439$	—	$9.409 \pm 0.022$	—
<i>Pixl2R</i> + Bi	$0.137 \pm 0.009$	-3.5%	$31.236 \pm 2.040$	151%	$10.735 \pm 0.784$	14%
<i>Pixl2R</i> + BiMI	$0.162 \pm 0.022$	14%	$37.507 \pm 7.832$	199%	$7.951 \pm 0.351$	-15%
<i>Pixl2R</i> + DEIR	$0.176 \pm 0.009$	23%	$17.372 \pm 0.514$	39%	$10.440 \pm 1.015$	10%
<i>Pixl2R</i> + BiMI + DEIR	$0.267 \pm 0.016$	88%	$57.759 \pm 2.157$	364%	$11.014 \pm 0.190$	17%
<i>ELLM</i> *	$0.150 \pm 0.004$	—	$19.406 \pm 10.067$	—	$10.826 \pm 1.017$	—
<i>ELLM</i> + Bi	$0.151 \pm 0.016$	0.6%	$29.788 \pm 1.290$	53%	$11.175 \pm 0.601$	3.2%
<i>ELLM</i> + BiMI	$0.156 \pm 0.014$	4.0%	$33.683 \pm 3.990$	74%	$9.425 \pm 0.267$	-12%
<i>ELLM</i> + DEIR	$0.178 \pm 0.029$	20%	$30.985 \pm 3.507$	59%	$11.857 \pm 1.152$	9.5%
<i>ELLM</i> + BiMI + DEIR	$0.279 \pm 0.078$	86%	$56.281 \pm 6.193$	190%	$13.170 \pm 0.393$	22%

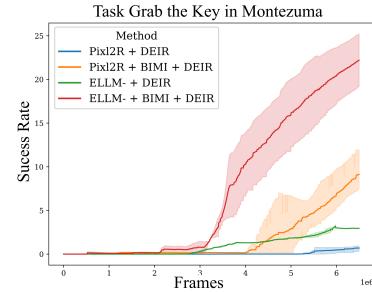


Figure 8: BiMI reward showed faster and higher success rates on difficult tasks in Montezuma

### 308 6.1 Overall Performance and Ablation Study

309 The overall improvements were substantial. As shown in Table 1, BiMI led to a 67% improve-  
 310 ment for *Pixl2R* and a 22% improvement for *ELLM-*. These results are also illustrated in Figure 8  
 311 Figure 13 in Appendix A.7. Our ablation study highlights the distinct contributions of the binary  
 312 reward (Bi) and Mutual Information (MI) components within the BiMI framework. The binary  
 313 reward mechanism alone accounted for a substantial 36.5% improvement in performance. When  
 314 excluding the results from *Crafter*, MI component further contributes a 23% improvement over the  
 315 binary reward alone. Please see Appendix A.7 for result details for each test environment.

## 316 7 Conclusion

317 Our research reveals two key findings in agent learning using VLM-based reward models: (1) false  
 318 positive rewards, rather than false negatives, are more detrimental to policy learning; and (2) our  
 319 proposed BiMI reward function significantly mitigates the slowdown in learning caused by false  
 320 positives, thereby paving the way for more reliable and effective “language as an interface” RL  
 321 systems in practical settings.

322 **Limitation.** Our study primarily focused on linear sequences of language instructions, excluding  
 323 more complex cases. Future research should investigate conditional and ambiguous instructions,  
 324 which likely introduce additional challenges for VLM-based reward models. We also did not explore  
 325 finetuning the VLM during agent training, a useful strategy as discussed by Fu et al. [11].

326 **References**

- 327 [1] R. Kaplan et al. Beating atari with natural language guided reinforcement learning. *ArXiv*,  
328 abs/1704.05539, 2017. URL <https://api.semanticscholar.org/CorpusID:6022828>.
- 329 [2] P. Goyal et al. Using natural language for reward shaping in reinforcement learning.  
330 In *International Joint Conference on Artificial Intelligence*, 2019. URL <https://api.semanticscholar.org/CorpusID:70350059>.
- 332 [3] P. Goyal et al. Pixl2r: Guiding reinforcement learning using natural language by mapping  
333 pixels to rewards. In *Language in Reinforcement Learning Workshop at ICML 2020*, 2020.
- 334 [4] Y. Du, O. Watkins, Z. Wang, C. Colas, T. Darrell, P. Abbeel, A. Gupta, and J. Andreas. Guid-  
335 ing pretraining in reinforcement learning with large language models. In *International Confer-  
336 ence on Machine Learning*, 2023. URL <https://api.semanticscholar.org/CorpusID:256846700>.
- 338 [5] X. E. Wang, Q. Huang, A. Celikyilmaz, J. Gao, D. Shen, Y. fang Wang, W. Y. Wang,  
339 and L. Zhang. Reinforced cross-modal matching and self-supervised imitation learning for  
340 vision-language navigation. *2019 IEEE/CVF Conference on Computer Vision and Pattern  
341 Recognition (CVPR)*, pages 6622–6631, 2018. URL <https://api.semanticscholar.org/CorpusID:53735892>.
- 343 [6] M. Shridhar et al. Cliprot: What and where pathways for robotic manipulation. In *Conference  
344 on robot learning*, pages 894–906. PMLR, 2022.
- 345 [7] P. Mahmoudieh et al. Zero-shot reward specification via grounded natural language. In *Inter-  
346 national Conference on Machine Learning*, 2022. URL <https://api.semanticscholar.org/CorpusID:250340669>.
- 348 [8] J. Clark. Faulty reward functions in the wild. [https://openai.com/index/  
349 faulty-reward-functions/](https://openai.com/index/faulty-reward-functions/), 2016. Accessed: 2024-08-06.
- 350 [9] Z. Wang, S. Cai, G. Chen, A. Liu, X. Ma, and Y. Liang. Describe, explain, plan and se-  
351 lect: Interactive planning with llms enables open-world multi-task agents. In *Neural Infor-  
352 mation Processing Systems*, 2023. URL <https://api.semanticscholar.org/CorpusID:268042457>.
- 354 [10] G. R. Ghosal, M. Zurek, D. S. Brown, and A. D. Dragan. The effect of modeling human ratio-  
355 nality level on learning rewards from multiple feedback types. In *AAAI Conference on Artificial  
356 Intelligence*, 2022. URL <https://api.semanticscholar.org/CorpusID:251740992>.
- 357 [11] Y. Fu, H. Zhang, D. Wu, W. Xu, and B. Boulet. Furl: Visual-language models as fuzzy re-  
358 wards for reinforcement learning. In *Forty-first International Conference on Machine Learn-  
359 ing*, 2024.
- 360 [12] D. Abel, W. Dabney, A. Harutyunyan, M. K. Ho, M. Littman, D. Precup, and S. Singh. On  
361 the expressivity of markov reward. *Advances in Neural Information Processing Systems*, 34:  
362 7799–7812, 2021.
- 363 [13] R. T. Icarte, T. Q. Klassen, R. A. Valenzano, and S. A. McIlraith. Using reward machines  
364 for high-level task specification and decomposition in reinforcement learning. In *Inter-  
365 national Conference on Machine Learning*, 2018. URL <https://api.semanticscholar.org/CorpusID:51868784>.
- 367 [14] J. Corazza et al. Reinforcement learning with stochastic reward machines. In *AAAI Conference  
368 on Artificial Intelligence*, 2022. URL <https://api.semanticscholar.org/CorpusID:250297195>.

- 370 [15] J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu. Roformer: Enhanced transformer with  
371 rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- 372 [16] T. Pham, T. Bui, L. Mai, and A. Nguyen. Out of order: How important is the sequential order  
373 of words in a sentence in natural language understanding tasks? In *Findings of the Association  
374 for Computational Linguistics: ACL-IJCNLP 2021*, pages 1145–1160, 2021.
- 375 [17] D. Hafner. Benchmarking the spectrum of agent capabilities. In *Deep RL Workshop NeurIPS  
376 2021*, 2021.
- 377 [18] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment:  
378 An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:  
379 253–279, 2013.
- 380 [19] M. Chevalier-Boisvert, D. Bahdanau, S. Lahlou, L. Willems, C. Saharia, T. H. Nguyen, and  
381 Y. Bengio. Babyai: A platform to study the sample efficiency of grounded language learning.  
382 In *International Conference on Learning Representations*, 2018.
- 383 [20] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,  
384 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from  
385 natural language supervision. In *International Conference on Machine Learning*, 2021. URL  
386 <https://api.semanticscholar.org/CorpusID:231591445>.
- 387 [21] Y. Ma, G. Xu, X. Sun, M. Yan, J. Zhang, and R. Ji. X-clip: End-to-end multi-grained  
388 contrastive learning for video-text retrieval. *Proceedings of the 30th ACM International  
389 Conference on Multimedia*, 2022. URL <https://api.semanticscholar.org/CorpusID:250607505>.
- 390 [22] T. Zhang, H. Xu, X. Wang, Y. Wu, K. Keutzer, J. E. Gonzalez, and Y. Tian. Noveld: A simple  
391 yet effective exploration criterion. In *Neural Information Processing Systems*, 2021. URL  
392 <https://api.semanticscholar.org/CorpusID:245877021>.
- 393 [23] M. M. Hossain et al. An analysis of negation in natural language understanding corpora. *arXiv  
394 preprint arXiv:2203.08929*, 2022.
- 395 [24] T. H. Truong, T. Baldwin, K. Verspoor, and T. Cohn. Language models are not naysayers: an  
396 analysis of language models on negation benchmarks. *arXiv preprint arXiv:2306.08189*, 2023.
- 397 [25] M. Sadinle et al. Least ambiguous set-valued classifiers with bounded error levels. *Journal of  
398 the American Statistical Association*, 114(525):223–234, 2019.
- 399 [26] M. Li, X. Zhao, J. H. Lee, C. Weber, and S. Wermter. Internally rewarded reinforcement  
400 learning. In *International Conference on Machine Learning*, pages 20556–20574. PMLR,  
401 2023.
- 402 [27] A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan. On the theory of policy gradient  
403 methods: Optimality, approximation, and distribution shift. *J. Mach. Learn. Res.*, 22(98):  
404 1–76, 2021.
- 405 [28] K. Hornik et al. Multilayer feedforward networks are universal approximators. *Neural net-  
406 works*, 2(5):359–366, 1989.
- 407 [29] S. Moon, J. Yeom, B. Park, and H. O. Song. Discovering hierarchical achievements in rein-  
408 forcement learning via contrastive learning. In *Neural Information Processing Systems*, 2023.

410 **A Technical Appendix**

411 Continued from the main text of *Overcoming Reward Model Noise in Instruction-Guided Reinforce-*  
412 *ment Learning*, the technical appendix consist of the following:

- 413 • §A.1 **Property of Sparse Reward Function**, which discuss why sparse reward function  
414 results in the categorization of policies into a set of acceptable policies and a set of unac-  
415 ceptable policies but not others.
- 416 • §A.2 **Proof of the Reduction of Convergence Rate**, which is referred in Section 3.
- 417 • §A.3 **Implementation Details of the Experiments**, which specify the implementation de-  
418 tails for both the first and second stage experiments in Section 4.1 and Section 6 respec-  
419 tively.
- 420 • §A.4 **Additional Details of the Experiments of False Positive Rewards**, which is referred  
421 in Section 4.1.
- 422 • §A.5 **Pseudo-code for Empirical Quantile Calculation for Binary Signal Threshold**,  
423 which is referred in Section 5.1.
- 424 • §A.6 **Additional Implementation Details of the Experiments of BiMI Reward Func-  
425 tion**, which is referred in Section 6.
- 426 • §A.7 **Detailed Experiment Results of BiMI Reward Function**, which provides detailed  
427 experiment results of BiMI reward function, which is referred in Section 6.

428 **A.1 Sparse Reward Function and Range-SOAP**

429 Formally, Abel et al. [12] defined that:

430 **Theorem A.1** (Abel et al. [12]). *A reward function realizes a Range Set of Acceptable Policies*  
431 *(Range-SOAP)  $\Pi_G$  when there exists an  $\epsilon \geq 0$  such that every  $\pi_g \in \Pi_G$  is  $\epsilon$ -optimal in start-state*  
432 *value,  $V^*(s_0) - V^{\pi_g}(s_0) \leq \epsilon$ , while all other policies are worse.*

433 When the reward signal is sparse, the agent only receives a reward upon reaching some goal states,  
434 and the reward function does not provide any feedback during the intermediate steps. We argue  
435 that the sparse reward function realizes a Range-SOAP, leading to the categorization of policies into  
436 acceptable and unacceptable policies.

437 **Proposition A.2.** *Sparse reward function “realizes” Range-SOAP (i.e., Range Set of Acceptable  
438 Policies).*

439 **Justification:**

440 A sparse reward function, where the agent only receives a reward upon completing the task, cannot  
441 prefer policies that lead to shorter task completion times. This is because either the agent completes  
442 the goal very quickly or slowly, they will receive nearly the same amount of cumulative rewards,  
443 and the reward function will not show a strong preference.

444 Since the sparse reward function does not induce a strict partial ordering on the policies, we say  
445 this reward function cannot realize a Partial Ordering (PO) task. Specifically, a PO on policies is  
446 a generalization of a Set of Acceptable Policies (SOAP) task. In a PO, the agent specifies a partial  
447 ordering on the policy space, where some policies are identified as “great”, some as “good”, and  
448 some as “bad” to strictly avoid, while remaining indifferent to the rest.

449 Therefore, the sparse reward function can realize a Set of Acceptable Policies (SOAP), where there  
450 is a set of policies that are all considered “good” or near-optimal, while all other policies are worse.

451 Furthermore, the sparse reward function will lead to a Range-SOAP, rather than an Equal-SOAP.  
452 Specifically, Equal-SOAP is a SOAP where all the acceptable policies are equally optimal in start-  
453 state value. This is because the good policies in the SOAP may differ slightly in their start-state  
454 values, as some may reach multiple goal states in the environment and thereby receiving different

455 cumulative rewards. Therefore, the sparse reward function will realize a Range-SOAP, where there  
456 is a range of acceptable policies that are all near-optimal in start-state value.

457 **A.2 Proof of the Reduction of Convergence Rate**

458 Specifically, the update rule of Actor-Critic algorithm is:

459 • **Critic:**

$$\phi \leftarrow \phi - \alpha_\phi \nabla_\phi (\delta)^2 \quad (4)$$

460 where

461  $\delta = \mathbb{E}_{\pi_\theta}[G^{\pi_\theta} - Q_\phi(s, a)]$  is the Monte Carlo (MC) estimation error,  $Q_\phi(s, a)$  is the  $Q$ -  
462 value function which measures the expected discounted cumulative reward given the  
463 state  $s$  and action  $a$ , and  $\pi_\theta$  is the policy.

464  $G^{\pi_\theta}$  is the rollout cumulative rewards from the trajectory  $\tau^{\pi_\theta}$  generated from  $\pi_\theta$ .

465 • **Actor:**

$$\theta \leftarrow \theta + \alpha_\theta \frac{Q_\phi(s, a) \nabla_\theta \pi_\theta(a|s)}{\pi_\theta(a|s)} \quad (5)$$

466 We need to make the following assumptions to simplify the theoretical analysis:

- 467 1. We use  $\Pi_G$  to represent the set of acceptable policies and use  $\Pi_B$  as an arbitrary set of “not  
468 acceptable” policies that are learned to follow instruction but fail to reach subgoals, i.e.,  
469  $\Pi_B \cap \Pi_G = \emptyset$ .
- 470 2. Assume the policy class parameterized by  $\theta$  should be expressive enough to capture optimal  
471 or near-optimal policies, and the policy is initialized randomly from uniform distribution,  
472 i.e.,  $\pi_{\theta_0} \sim \mathcal{U}(\Pi)$ . Meanwhile, the Q-value is initialized as  $Q_{\phi_0}(s, a) = 0, \forall s \in S \forall a \in A$ .
- 473 3. We assume that  $|\Pi_G|$  and  $|\Pi_B|$  is a fixed number predefined by the task environment while  
474  $\mathbb{E}[G^\pi | \pi \in \Pi_B]$  is non-zero as false positive rewards are unavoidable in real-world VLMs.

475 Since the update rule for  $Q$ -value is a gradient descent on  $\|\mathbb{E}[G^{\pi_\theta} - Q_\phi(s, a)]\|^2$  and also we have  
476 that  $\{\Pi_B, \Pi_G, \Pi \setminus (\Pi_B \cup \Pi_G)\}$  is a countable partition of the policy universe  $\Pi$ , the updated  $Q$ -value  
477 will approach as follows:

$$\begin{aligned} Q_\phi(s, a) &\rightarrow \mathbb{E}[G^{\pi_\theta}] \\ &= P(\pi_\theta \in \Pi_B) \cdot \mathbb{E}[G^{\pi_\theta} | \pi_\theta \in \Pi_B] \\ &\quad + P(\pi_\theta \in \Pi_G) \cdot \mathbb{E}[G^{\pi_\theta} | \pi_\theta \in \Pi_G] \\ &\quad + P(\pi_\theta \in \Pi \setminus (\Pi_B \cup \Pi_G)) \\ &\quad \cdot \mathbb{E}[G^{\pi_\theta} | \pi_\theta \in \Pi \setminus (\Pi_B \cup \Pi_G)] \\ &= P(\pi_\theta \in \Pi_B) \cdot \mathbb{E}[G^{\pi_\theta} | \pi_\theta \in \Pi_B] \\ &\quad + P(\pi_\theta \in \Pi_G) \cdot \mathbb{E}[G^{\pi_\theta} | \pi_\theta \in \Pi_G] \end{aligned} \quad (6)$$

478 Given that the update rule for policy  $\pi$  is the gradient ascent on  $Q_\phi(s, a)\pi_\theta(a|s)$ , we have the  
 479 following:

$$\begin{aligned}
 & \nabla_\theta Q_\phi(s, a)\pi_\theta(a|s) \\
 &= (\nabla_\theta P(\pi_{\theta_{old}} \in \Pi_B) \cdot \mathbb{E}[G^{\pi_{\theta_{old}}} | \pi_{\theta_{old}} \in \Pi_B] \pi_\theta(a|s)) \\
 &\quad + (\nabla_\theta P(\pi_{\theta_{old}} \in \Pi_G) \\
 &\quad \cdot \mathbb{E}[G^{\pi_{\theta_{old}}} | \pi_{\theta_{old}} \in \Pi_G] \pi_\theta(a|s)) \\
 &= P(\Pi_B) \cdot \mathbb{E}[G^{\Pi_B}] \nabla_\theta \pi_\theta(a|s) \\
 &\quad + P(\Pi_G) \cdot \mathbb{E}[G^{\Pi_G}] \nabla_\theta \pi_\theta(a|s) \\
 &= (1 - P(\Pi_G) - P(\Pi_B)) \\
 &\quad \cdot \mathbb{E}[G^{\Pi_B}] \nabla_\theta \pi_\theta(a|s) + P(\Pi_G) \cdot \mathbb{E}[G^{\Pi_G}] \nabla_\theta \pi_\theta(a|s) \\
 &= const \cdot \mathbb{E}[G^{\Pi_B}] \nabla_\theta \pi_\theta(a|s) \\
 &\quad + (\mathbb{E}[G^{\Pi_G}] - \mathbb{E}[G^{\Pi_B}]) P(\Pi_G) \nabla_\theta \pi_\theta(a|s)
 \end{aligned} \tag{7}$$

480 Justification of the assumptions.

- 481 Regarding the non-zero probability of recovering the optimal policy at initialization, it is  
 482 standard in theoretical analyses to assume a uniform distribution of a random variable at  
 483 initialization (see Agarwal et al. [27]). This assumption does not contradict the conclusion  
 484 about the convergence rate deterioration in Theorem 3.1.
- 485 In reference to the realizability condition implied by Assumption 2, the expressiveness of  
 486 the policy class parameterized by  $\theta$  is an underlying assumption for deep learning models,  
 487 supported by the The Universal Approximation Theorem [28].

488 **Observation.** Since the goal of the learning agent is to maximize  $P(\pi \in \Pi_G)$  (i.e., to converge  
 489 to an acceptable policy), we can see that the second term provides the target direction with rate  
 490  $(\mathbb{E}[G^{\pi \in \Pi_G}] - \mathbb{E}[G^{\pi \in \Pi_B}])$ . Therefore, the ascent rate will decrease when the cumulative rewards  
 491 from trajectories that cannot reach the goal state (i.e., the false positive rewards) gets higher. More-  
 492 over, the first term  $const \cdot \mathbb{E}[G^{\Pi_B}] \nabla_\theta \pi_\theta(a|s)$  can be regarded as the deviation of target direction. It  
 493 shows that the level of deviation is also positively proportional to the magnitude of rewards coming  
 494 from failed trajectories. This theorem follows the intuition that the presence of false positive rewards  
 495 can slow down the convergence rate of the learning agent.

### 496 A.3 Implementation Details of the Experiments

#### 497 A.3.1 Environment Details

498 We describe each testing environment used in our experiments. More details introduction can be  
 499 found in on the official project homepage of each benchmark [19, 18, 17].

- 500 • **Crafter** features randomly generated 2D worlds where the player needs to forage for food  
 501 and water, find shelter to sleep, defend against monsters, collect materials, and build tools.  
 502 The original Crafter environment does not have a clear goal trajectory or instructions;  
 503 agents are aimed at surviving as long as possible and exploring the environment to un-  
 504 lock new crafting recipes. We modified the environment to include a preset linear sequence  
 505 of instructions to guide the agent to mine diamond. However, this instruction was found  
 506 to hinder the agent’s performance. The nature of the task requires dynamic strategies and  
 507 real-time decision-making, but the fixed instructions limited the agent. For example, the  
 508 instruction did not account for what to do when the agent is attacked by zombies.
- 509 • **Montezuma’s Revenge** is a classic adventure platform game where the player must nav-  
 510 igitate through a series of rooms to collect treasures and keys. The game is known for its  
 511 sparse rewards and challenging exploration requirements. We manually annotate 97 in-  
 512 structions for the agent to follow, guiding it to conquer the game. The instructions were

513 designed to guide the agent through the game’s key challenges, such as avoiding enemies,  
 514 collecting keys, and unlocking doors.

- 515 • **Minigrid ‘Go to seq’ Task:** We use the ‘Go to seq’ task in the Minigrid environment,  
 516 where the agent must navigate through a sequence of rooms and touch target objects in the  
 517 correct order. This is a sparse reward task where the agent receives a reward of 1 only upon  
 518 completing the entire sequence correctly. During the training phase, we randomly generate  
 519 50 different tasks, each with a room size of 5, 3 rows, and 3 columns. Each task features  
 520 a unique room layout and target object sequence. The instruction complexity is set to 3,  
 521 meaning there are at least 3 target objects to interact with in a specific order.

### Montezuma and Instructions



Figure 9: Illustration of the Montezuma’s Revenge task. The agent must navigate through a series of rooms to collect treasures and keys.

### Minigrid and Instructions

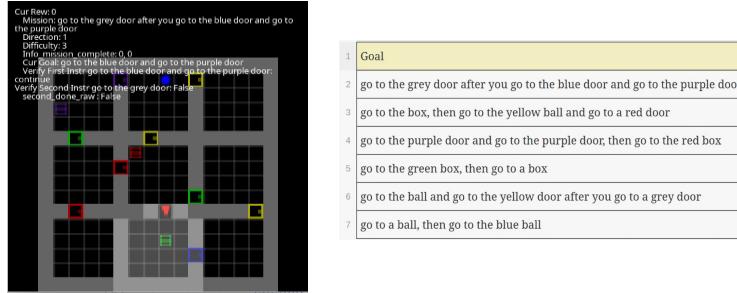


Figure 10: Illustration of the Minigrid ‘Go to seq’ task. The agent must navigate through a sequence of rooms and touch target objects in the correct order.

### Crafter and Instructions



Figure 11: Illustration of the Crafter task. The agent must survive as long as possible and explore for new crafting recipes.

522 **A.3.2 Instruction-Guide Procedure Details**

523 The VLM-based reward model will have a pointer to the sequence of the instruction sentence, start-  
 524 ing at the first sentence. For original models *Pixel2R* and *ELLM*-, we follow the setting in there  
 525 original work where for each instruction sentence, the reward model will have a maximum cap of  
 526 rewards (2.0) it can assign to the agent in one episode. When the cap is reached, the reward model  
 527 will move its pointer to the next instruction sentence. For the BiMI reward model, the reward model  
 528 will move its pointer to the next instruction sentence when the binary signal is triggered.

529 **A.3.3 Finetuning VLM-based Reward Models**

530 In contrast to previous work on Instruction-guided RL where they rely on hand-crafted oracle multi-  
 531 modal reward models, we use actual pretrained VLMs to generate reward signals. 2 VLM backbone  
 532 models are used in our experiments: 1) *CLIP* [20], pretrained by image-text pairs; and (2) *X-CLIP*  
 533 [21], pretrained by video-text pairs. In particular, *Pixel2R* uses *CLIP* because it only uses the single  
 534 latest frame as input. In contrast, *ELLM*- takes a slice of trajectory (i.e., multiple frames) as input,  
 535 and thus uses either *X-CLIP* or *CLIP* with additional RNN encoder as the reward model.

536 Due to the cartoonish and abstract visuals of the testing environments, we need to further fine-  
 537 tune the VLMs to adapt to this new visual domain. We use well-trained expert agents based on  
 538 Moon et al. [29] to generate expert trajectories for the Crafter environments and annotate them  
 539 with instructions using internal information from the game engine. For Minigrid environments, we  
 540 use classical search-based planning robots to generate expert trajectories and annotate them with  
 541 the corresponding task instructions. For Montezuma’s Revenge, we manually annotate the expert  
 542 trajectories.

543 For Minigrid and Crafter, we have 80,000 training pairs, while for Montezuma’s Revenge, we have  
 544 around 300 training pairs. These training data are of high quality, as we have made every effort to  
 545 avoid false positive rewards due to poor training data quality. **To enhance our models’ robustness,**  
 546 **we also employed contrastive learning techniques during VLM training, utilizing similar ma-**  
 547 **nipulated data as hard negatives.** However, despite the fine-tuning process, false positive rewards  
 548 remain unavoidable.

```

1  data_id,instruction,trajectory_chunk_file,trajectory_local_idx
2  0,climb down the middle ladder,montezuma/expert_traj_chunk_0.pkl,0
3  1,walk to the right side of the conveyor belt,montezuma/expert_traj_chunk_0
4  2,jump right to the yellow rope,montezuma/expert_traj_chunk_0.pkl,2
5  3,jump right to the right platform,montezuma/expert_traj_chunk_0.pkl,3
6  4,climb down the right ladder,montezuma/expert_traj_chunk_0.pkl,4
7  5,jump over the skull,montezuma/expert_traj_chunk_0.pkl,5
8  6,climb up the left ladder,montezuma/expert_traj_chunk_0.pkl,6
9  7,jump to grab the key,montezuma/expert_traj_chunk_0.pkl,7
10 8,jump left to the left roof ,montezuma/expert_traj_chunk_0.pkl,8
11 9,use key to open the left door,montezuma/expert_traj_chunk_0.pkl,9
12 10,walk left when the laser gate disappears,montezuma/expert_traj_chunk_0.p
13 11,walk to the middle when the laser gate disappears,montezuma/expert_traj_
14 12,wait until the laser gate disappears,montezuma/expert_traj_chunk_0.pkl,1
15 13,approach to the gem,montezuma/expert_traj_chunk_0.pkl,13
16 14,jump to grab the gem,montezuma/expert_traj_chunk_0.pkl,14
17 15,walk to the middle when the laser gate disappears,montezuma/expert_traj_
18 16,climb down the middle ladder,montezuma/expert_traj_chunk_0.pkl,16
19 17,wait until the spider goes away,montezuma/expert_traj_chunk_0.pkl,17

```

Figure 12: Example of training data for the Montezuma environment.

549 We used the threshold  $\hat{q}$  introduced in Section 5.1 to make binary classification on the testing pairs  
 550 to evaluate the performance of the fine-tuned VLM-based reward models. We found that VLM  
 551 models had difficulty achieving high accuracy on Minigrid environment, which is likely due to the

552 too abstract and cartoonish nature of the environment, causing the VLMs to struggle to learn the  
 553 visual-textual correspondence. We also found that X-CLIP did not perform better than CLIP in  
 554 our experiments. We hypothesize that the cartoonish nature of the testing environments may have  
 555 caused the X-CLIP model to struggle to learn the visual-textual correspondence. Thus, we used  
 556 CLIP as the backbone model throughout our following experiments. The performance of the fine-  
 557 tuned VLM-based reward models is shown in Table 2. **Even when the precision score reaches 0.98,**  
 558 **indicating that only 2% of the rewards are false positives in the validation set, the agent can**  
 559 **still significantly underperform in the testing environments. The core issue is that in out-of-**  
 560 **distribution (O.O.D.) testing environments, false positive rewards are prevalent and inevitable.**  
 561 **Therefore, it is crucial to design a reward function that is robust to reward noise.**

Table 2: Performance of fine-tuned VLM reward model on the testing dataset using the 90th percentile empirical quantile as threshold

Environment	Precision	Accuracy	F1 Score	Recall	Model
Crafter	0.9847	0.9466	0.8538	0.9702	CLIP ELLM-
Crafter	0.9799	0.9028	0.7618	0.9842	CLIP Pixl2R
Crafter	0.2095	0.2514	0.2868	0.9657	XCLIP ELLM-
Minigrid	0.7260	0.9200	0.7849	0.9763	CLIP ELLM-
Minigrid	0.6992	0.9086	0.7592	0.9616	CLIP Pixl2R
Minigrid	0.1716	0.2310	0.2642	0.9704	XCLIP ELLM-
Montezuma	0.8838	0.9638	0.8825	0.9478	CLIP ELLM-
Montezuma	0.8343	0.9108	0.7652	0.9842	CLIP Pixl2R
Montezuma	0.8044	0.9259	0.8045	0.9657	XCLIP ELLM-

### 562 A.3.4 Hyperparameters for Instruction-Guided RL Agents

563 In the experiments, all methods are implemented based on PPO with same model architecture. The  
 564 Minigrid and Crafter environments use the same training hyperparameters as the Achievement Dis-  
 565 tillation paper [29]. For Montezuma’s Revenge, we found that the performance of the agent was  
 566 sensitive to the gamma and GAE lambda parameters. To improve the performance of agents in  
 567 Montezuma’s Revenge, we took two additional steps: (1) normalizing the observation inputs when  
 568 computing the rewards, and (2) not normalizing the advantage during the GAE calculation. The  
 569 hyperparameters are shown in the following tables.

Table 3: Model Parameters

Parameter	Value
model_cls	“PPORNNModel”
hidsize	1024
gru_layers	1
impala_kwarg	
- chans	[64, 128, 128]
- outsize	256
- nblock	2
- post_pool_groups	1
- init_norm_kwarg	
- batch_norm	false
- group_norm_groups	1
dense_init_norm_kwarg	
- layer_norm	true

Table 4: Crafter and Minigrid RL Parameters

Parameter	Value
gamma	0.95
gae_lambda	0.65
algorithm_cls	“PPOAlgorithm”
algorithm_kwargs	
- ppo_nepoch	3
- ppo_nbatches	8
- clip_param	0.2
- vf_loss_coeff	0.5
- ent_coeff	0.01
- lr	3.0e-4
- max_grad_norm	0.5
- aux_freq	8
- aux_nepoch	6
- pi_dist_coeff	1.0
- vf_dist_coeff	1.0

Table 5: Montezuma RL Training Parameters

Parameter	Value
gamma	0.99
gae_lambda	0.95
int_rew_type	“rnd”
pre_obs_norm_steps	50
algorithm_cls	“PPOAlgorithm”
algorithm_kwargs	
- update_proportion	0.25
- ppo_nepoch	3
- ppo_batch_size	256
- clip_param	0.1
- vf_loss_coeff	0.5
- ent_coeff	0.001
- lr	1.0e-4

#### 570 A.4 Additional Details of the Experiments on the Impact of Noisy Rewards

571 **Evaluation Metric Details** In our experiments, we used a score metric adapted from the Crafter  
 572 paper to evaluate agent performance across different environments. This score metric aggregates  
 573 success rates for individual subtasks using a geometric mean. This metric was chosen over the  
 574 *maximum total rewards* metric for several reasons:

- 575 **Consistency in Sparse Reward Settings:** Sparse reward environments often pose significant  
 576 challenges for reinforcement learning agents. An agent might occasionally achieve  
 577 high rewards by chance in one rollout but fail to replicate this success consistently in sub-  
 578 sequent rollouts. This variability can lead to misleading evaluations if only the maximum  
 579 total rewards are considered. The Score metric, by measuring the success rate of achieving  
 580 each subgoal, provides a more stable and consistent measure of an agent’s performance.
- 581 **Capturing Learning Stability:** The Score metric evaluates the agent’s ability to consist-  
 582 ently reproduce successful behaviors across multiple episodes. This is crucial in sparse  
 583 reward settings, where the agent’s performance can fluctuate significantly. By focusing on  
 584 the success rates of individual subtasks, the Score metric offers a more granular and reliable  
 585 assessment of the agent’s learning progress and stability.
- 586 **Crafter Benchmark Standard:** The Crafter benchmark, which introduces the Score met-  
 587 ric, is a well-regarded standard.

588 Crafter codebase provides *score* metric calculation by default. For Minigrid and Montezuma environments,  
 589 we use the internal information from the game engine to detect whether the subtasks are  
 590 completed, thus facilitating the calculation of the *score* metric.

#### 591 A.4.1 Details of Manipulated Trajectory-Instruction Pairs to Evaluate Robustness

592 We evaluated the models' sensitivity by examining how cosine similarity scores change for manipulated  
 593 trajectory-instruction pairs. These manipulations were designed to test the robustness of the  
 594 models against various types of noise. Here's a detailed breakdown of the manipulations:

- 595     1. **Trajectory Reversal:** We inverted the sequence of frames within each trajectory (i.e.,  
       596        `frames = frame[::-1]`) to assess the model's ability to detect reversed state transitions.  
       597        This manipulation tests whether the model can distinguish between forward and backward  
       598        progression in the state transition.
- 599     2. **Instruction Negation:** We modified the original instructions by adding negation (e.g.,  
       600        changing "do  $l_k$ " to "do not do  $l_k$ " or "avoid  $l_k$ "). This tests the model's sensitivity to  
       601        semantic changes in the instruction that fundamentally alter the goal.
- 602     3. **Instruction Rephrasing:** We rephrase the original instructions while maintaining their  
       603        core meaning. This evaluates the model's robustness to linguistic variations and its ability  
       604        to capture the essential semantic content of instructions.
- 605     4. **Concatenation and Order Swapping:** Given two trajectory-instruction pairs  $(\tau_1, l_1)$  and  
       606         $(\tau_2, l_2)$ , we created concatenated pairs and then swapped the order in one modality. For  
       607        example:

- 608        • Original concatenation:  $(\tau_1 + \tau_2, l_1 + l_2)$
- 609        • Swapped trajectory:  $(\tau_2 + \tau_1, l_1 + l_2)$
- 610        • Swapped instruction:  $(\tau_1 + \tau_2, l_2 + l_1)$

611 This tests the model's sensitivity to the order of components in multi-step tasks.

- 612     5. **Concatenation with Partial Content:** We concatenated pairs but truncated one modality.  
       613        For instance:

- 614        • Truncated trajectory:  $(\tau_1, l_1 + l_2)$
- 615        • Truncated instruction:  $(\tau_1 + \tau_2, l_1)$

616 This assesses the model's ability to detect partial mismatches in longer sequences.

#### 617 A.5 Pseudo-code for Empirical Quantile Calculation for Binary Signal Threshold

618 Using empirical quantile as threshold guarantees a high probability (at least  $1 - \alpha$ ) that the true  
 619 positive pairs are recognized while minimizing the average number of mistakes predicting false  
 620 positives [25]:

---

##### Algorithm 1 Calculate Empirical Quantile ( $\hat{q}$ )

---

**Require:** Calibration set  $\{\tau, l\}_n$ , where  $l$  is the instruction sentence,  $\tau$  is the corresponding trajectory, and  $n$  is the number of samples;

Significance level  $\alpha$ ;

VLM model reward model  $v$

- 1:  $\triangleright$  Obtain the similarity-based score
  - 2:  $\{r\}_n \leftarrow \{v(\tau, l)\}_n$
  - 3:  $\triangleright$  Compute the quantile level
  - 4:  $q_{level} \leftarrow \frac{\lceil (n-1) \times (1-\alpha) \rceil}{n}$
  - 5:  $\triangleright$  Compute the empirical quantile
  - 6:  $\hat{q} \leftarrow np.quantile(\{r\}_n, q_{level}, method='lower')$
  - 7: **return**  $\hat{q}$
-

621 **A.6 Implementation Details of the Experiments of BiMI Reward Function**

622 We set confidence level for empirical quantile calculation to be  $1 - \alpha = 0.9$ . We adhered to the  
623 standard requirement of limiting the training budget to 1 million frames [17]. This constraint poses  
624 a significant challenge, particularly in sparse reward settings, as it demands that agents both explore  
625 efficiently and exploit their knowledge effectively within this limited budget.

626 To achieve the 1 million frame budget, we used the following configuration:

- 627 • nproc: 8 (Number of processes used for parallel environments)  
628 • nstep: 512 (Length of the rollout stored in the buffer)  
629 • nepoch: 250 (Number of epochs to train the RL policy)

630 This configuration results in approximately 1 million steps: 250 epochs  $\times$  512 steps  $\times$  8 processes  
631 = 1,024,000 frames.

632 In the case of Montezuma’s Revenge, we found that the 1 million frame limit used in Crafter was  
633 insufficient due to the game’s complexity and sparse reward structure. To address this, we extended  
634 the training budget to 8 million frames. It’s important to note that even with this increased frame  
635 count, agents were still unable to fully solve the task. As Zhang et al. [22] pointed out, about 1  
636 billion frames are required to truly master Montezuma’s Revenge. This vast difference in required  
637 training time (8 million vs 1 billion frames) underscores the exceptional difficulty of Montezuma’s  
638 Revenge as a sparse reward task.

639 The implementation details for the BiMI reward function are consistent with those outlined in the  
640 first stage of experiments.

641 **A.7 Detailed Experiment Results of BiMI Reward Function**

642 **A.7.1 Montezuma’s Revenge**

643 In *Montezuma*, Pixl2R+BiMI demonstrated an 14% performance increase compared to the original  
644 models (see Table 1), which is slightly below our expectations. We attribute this result to BiMI’s in-  
645 tentional strategy of providing less frequent discrete rewards. While this strategy effectively reduces  
646 false positives, it does not substantially mitigate the inherent reward sparsity issue in *Montezuma*.  
647 However, **we discovered a remarkable synergy between BiMI and intrinsic reward models**.  
648 While previous models showed no significant improvements with *DEIR* (the intrinsic reward model)  
649 alone, combining BiMI and *DEIR* led to a 65% performance gain. The gap in collaboration ef-  
650 fectiveness can be attributed to two factors. In the previous setup, the consistent presence of false  
651 positive rewards misled agents towards unacceptable behaviors and hindered further exploration.  
652 Now, BiMI’s less frequent but more meaningful rewards provide anchor points for the agent’s learn-  
653 ing. Meanwhile, *DEIR*’s intrinsic rewards fill the gaps between these anchor points, encouraging the  
654 agent to explore efficiently in the interim.

655 As illustrated in Figure 6, BiMI rewards are now concentrated on key locations. A significant  
656 improvement is the minimal rewards given for falling off cliffs, which was a common source of  
657 false positives in the original model. Figure 8 demonstrates higher success rate grabbing the key  
658 in the first room, one of the most difficult task in *Montezuma*, highlighting the effectiveness of the  
659 proposed reward function and its collaboration with intrinsic reward models in guiding agents to  
660 solve difficult sparse-reward tasks.

661 **A.7.2 Minigrid**

662 ELLM-+BiMI achieved a remarkable 74% improvement in performance compared to the original  
663 models. This substantial gain is particularly noteworthy given the unique challenges presented by  
664 *Minigrid*. The abstract, shape-based visuals in *Minigrid* diverge drastically from the natural images  
665 used in VLMs’ pretraining, preventing the models from effectively utilizing their prior pretraining

knowledge. Consequently, VLMs struggled to accurately assess similarities between Minigrid's abstract visuals and textual instructions, resulting in highly noisy reward signals. The significant improvement demonstrated by BiMI underscores its effectiveness in handling noisy signals, directly addressing our primary research challenge. This capability is crucial for deploying instruction-guided agents in real-world, unfamiliar scenarios, where visual inputs often deviate from the VLMs' training distribution, leading to noisy reward signals.

### 672 A.7.3 Crafter

673 We observed an intriguing pattern of results. The BI component alone led to 14% and 3.2% improvement in performance over the original models. However, contrary to our observations in other environments, the addition of the MI component actually decreased this improvement. This unexpected outcome can be attributed to the unique nature of *Crafter* task, where agents must repeatedly achieve the same subtasks (e.g., drinking water) for survival. The MI component, designed to discourage over-reliance on frequently occurring signals, inadvertently penalized the necessary repetition of survival-critical actions. Furthermore, note that instruction-guided RL agents, regardless of the reward model employed, were unable to outperform pure RL agents in *Crafter*. This discrepancy is due to the open-world nature of *Crafter*, which requires dynamic strategies and real-time decision-making that our testing setups did not fully capture. Despite these challenges, it is noteworthy that BI alone still managed to improve performance over vanilla VLM-based reward models, suggesting that reducing false positives is still beneficial across all testing environments. The combination of BiMI with *DEIR* (the intrinsic reward model) also showed promising results, indicating a productive balance between exploration (driven by *DEIR*) and exploitation (guided by BiMI instruction reward).

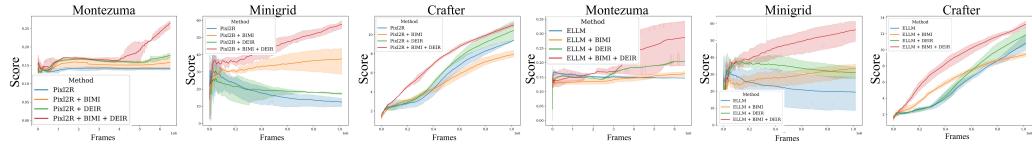


Figure 13: Besides the improvements of the score performance of agents across different environments with the BiMI reward function, it also collaborates well with intrinsic rewards. Combining both can lead to significant performance improvements

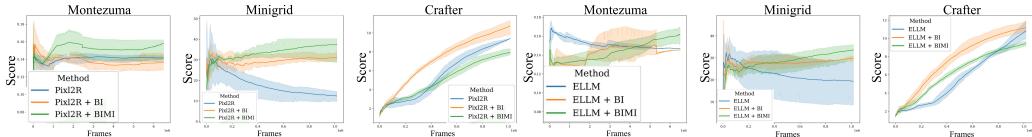


Figure 14: Ablation on the components of BiMI reward function. The binary reward (Bi) alone led to a 36.5% improvement compared to original models. Excluding Crafter, Mutual Information (MI) provided a 23% further improvement over Bi alone