

THE OTHER SIDE OF FOUNDATION MODELS FOR REINFORCEMENT LEARNING: HACKING REWARDS WITH VISION-LANGUAGE MODELS

CycleResearcher

ABSTRACT

Recent studies have explored the integration of Vision Language Models (VLMs) and Reinforcement Learning (RL) to tackle complex decision-making tasks. By leveraging the zero-shot captioning capabilities of pre-trained VLMs, an agent can be trained to maximize rewards generated through text prompts. Despite the promise of these recent advances, we reveal a potentially significant limitation: generated rewards are susceptible to hacking. This means that an agent, when manipulated in-env, can inadvertently cause poor performance under true rewards. To illustrate this, we conduct experiments across six distinct environments that span both visual and state inputs, as well as manipulation and navigation tasks. Notably, our findings demonstrate that reward hacking is prevalent in all these setups. Given the lack of prior research on hacking in the context of rewards generated by VLMs for RL agents, we provide a comprehensive analysis of the root cause of this phenomenon and discuss potential mitigation strategies. Our findings underscore the need for increased vigilance when deploying such methods in real-world applications.

1 INTRODUCTION

Reinforcement Learning (RL) has achieved impressive results in training autonomous agents to solve complex games (Mnih et al., 2013; Silver et al., 2016; Ma et al., 2023a) and long-horizon manipulation tasks (Yu et al., 2020; Agostinelli et al., 2019). Nevertheless, scaling RL to real-world scenarios, where tasks are often diverse and demands for high performance are strict, poses a substantial challenge (Khalshnikov et al., 2018). One of the most important roadblocks to the broader applicability of RL in real-world scenarios is the need for time-consuming human supervision to design rewards (Sutton & Barto, 2018). Some researchers have sought to address this limitation by leveraging recent advances with Vision-Language Models (VLMs) (Radford et al., 2021; Alayrac et al., 2022). These models, pre-trained on large datasets of aligned text and image pairs, enable zero-shot image captioning. The idea is to construct a reward function that utilizes these captioning abilities to closely align with human objectives.

In this work, we closely study *reward hacking* (Skalse et al., 2022) in the context of RL agents trained with Vision-Language Models (VLMs). Reward hacking occurs when optimizing a proxy reward function leads to poor performance under the true reward function. While this phenomenon has been extensively studied in classical RL, its study in the context of modern RL problems is relatively overlooked. In our experiments, we observe that RL agents trained with rewards generated by VLMs can exhibit significantly worse performance than those trained with orchestrated human-designed rewards. In other words, the RL agents hack the VLM-based rewards. To investigate this, we study complex tasks across various domains that span different inputs (visual and state), different agents (visually guided DrQ-v2 (Yarats et al., 2022) and the MP agent (Ahn et al., 2022)), and different task types (manipulation in Meta-World (Yu et al., 2020) and navigation in House3D (Das et al., 2018)). Our new findings demonstrate that such reward hacking is prevalent in diverse scenarios. As a matter of fact, we discover that RL agents can consistently exploit the proxy VLM-based rewards for different text prompts, even if these prompts are semantically distant.

Why does reward hacking occur in the context of VLM-generated rewards? We demonstrate that the primary cause is *reward unfaithfulness*: the reward function fails to faithfully reflect the true objective.

We study this issue from two perspectives. First, as illustrated in Figure ??, a faithful reward function should be able to generate an acceptable policy for the true objective. However, we observe that it is not the case; RL agents with VLM-based rewards underperform, even after extensive training. Second, a faithful reward would lead to generalizable policies across various tasks and domains that share similar semantics. However, we find the trained policies with VLM-based rewards fail to transfer or generalize.

We further investigate the generalizability of VLM-based rewards by testing a wide range of prompts to assess whether any of them can lead to reasonable performance. Interestingly, our results reveal a more concerning issue: while some prompts yield poor performance, others prompt the agent to seek a different goal. Specifically, the agent is guided to find a blue cube rather than a red cube. Notably, these misleading prompts not only lead the agent away from the objective but also result in the agent ignoring the target object entirely. Finally, we analyze the severe nature of reward hacking. We show that the performance under the true reward gradually declines as training with the generated reward progresses. This indicates that the generated reward, while seeming beneficial at first, ultimately leads to misguided learning and poor performance.

In light of these findings, it is essential to recognize the inherent limitations of using VLMs to generate reward functions for RL agents. Although they can appear advantageous at first, they can ultimately result in more harm than good. We analyze possible mitigation methods and discuss why they are challenging to apply in the context of using VLMs to train RL agents. Our findings underscore the need for more advanced approaches to leverage the capabilities of VLMs in RL.

2 RELATED WORK

2.1 REINFORCEMENT LEARNING WITH VISION-LANGUAGE MODELS

The use of Vision Language Models (VLMs) in reinforcement learning has gained increasing attention. On one hand, researchers have studied visual representation learning for RL using VLMs. For instance, Chen et al. (2023) discovered that using pre-trained VLMs' representations as part of the state space can enhance the performance of RL agents. Additionally, Nair et al. (2023) pre-trains a visual representation using a pre-trained VLM and demonstrates its effectiveness in a multi-task robotic manipulation setting.

Another line of research focuses on the use of VLMs to guide RL agents. A recent approach utilizes a VLM to replace the need for expert demonstrations by reanimating previously successful trajectories with new text prompts (Sun et al., 2023). Rocamonde et al. (2023a); Chan et al. (2023) propose to use pre-trained VLMs as zero-shot reward models to train RL agents. The intuition is that the zero-shot captioning capabilities of a VLM can be leveraged to generate reward signals aligned with human objectives. This approach addresses the need for expert demonstrations and enables the training of RL agents purely through textual supervision. This can be further enhanced through the integration of the VLM captioning function with the reward model (Mahmoudieh et al., 2022; Rocamonde et al., 2023b; Du et al., 2023; Lubana et al., 2023; Adeniji et al., 2023; Dang et al., 2023; Hu et al., 2023; Nam et al., 2023). For example, by incorporating semantic information to enhance faithfulness or by using a Flamingo model (Alayrac et al., 2022) to improve the understanding of more complex tasks. Furthermore, recent studies have shown that VLMs can be used to generate intrinsic motivation rewards when training RL agents, eliminating the need for external supervision (Klissarov et al., 2023; Ma et al., 2023b).

2.2 MITIGATING HALLUCINATION IN VISION-LANGUAGE MODELS

Recent works have shown that VLMs can suffer from hallucination, where the model generates text that is not aligned with the input image. In this work, we show that RL agents trained with rewards generated by VLMs can exhibit significantly worse performance than those trained with orchestrated human-designed rewards. Our new findings demonstrate that such reward hacking is prevalent in diverse scenarios.

Hallucination in VLMs, similar to LLM hallucinations (Chakraborty et al., 2024), can be categorized into two main types. The first involves the misinterpretation of visual features. For instance, FLamingo struggles to accurately caption images of camels with multiple humps, referring to them as

“camels with one hump” (Li et al., 2023). This indicates that the model has misinterpreted the image, leading to incorrect text generation. Another work demonstrates that asking FLINGO (Alayrac et al., 2022) to identify typical objects in a scene can result in the generation of non-existent objects such as an ‘invisible man’ (Yao et al., 2023).

The second category of hallucination involves the insertion of elements into an image that is not present. Notably, this type of hallucination can occur more frequently than the first, with its prevalence ranging from 20% to 50% across various benchmarks (Li et al., 2023). Additionally, (Li et al., 2023) shows that this issue is difficult to address. Recent studies have investigated the underlying reasons for such hallucinations and offer insights as to why they occur. For instance, (Ilyas et al., 2019) suggests that non-robust features may give rise to adversarial examples. (Zhang & Wang, 2019) demonstrates that these features are inherent in the data distribution and that they are difficult to remove. In this work, we report similar findings in the context of RL training. Specifically, we show that an agent can exploit VLM-based rewards by slightly modifying the text prompt, and the agent successfully achieves another goal.

3 PRELIMINARY: REINFORCEMENT LEARNING AND VISION-LANGUAGE MODELS

We study reinforcement learning tasks that can be formulated as continuous control. A task is defined as a Markov Decision Process (MDP), denoted as $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $\mathcal{T}(s'|s, a)$ represents the transition probability function, \mathcal{R} is the reward function, and γ is the discount factor. The goal is to train a policy π that maximizes the expected return.

To leverage pre-trained Vision-Language Models (VLMs) as reward functions, we make use of the captioning abilities of CLIP (Radford et al., 2021) and FLINGO (Alayrac et al., 2022). More precisely, we leverage the CLIP model to train RL agents and use FLINGO to conduct empirical studies. It is important to note that while other VLMs exist, our focus is on these two due to their superior captioning abilities. In the following, we briefly introduce CLIP and FLINGO.

CLIP takes a text prompt and an image as inputs and generates a scalar indicating the similarity between the text and the image. It is pre-trained on 400 million image-text pairs using a contrastive objective. For RL, we use a text prompt (e.g., “Find a red cube”) and the agent’s observation (e.g., an image of a red cube) as inputs to CLIP, which outputs a scalar indicating whether the observation contains a red cube.

FLINGO is a hybrid architecture that combines a pre-trained vision encoder with a pre-trained language encoder. This combination is trained using image-text contrastive learning. As a result, FLINGO not only excels in captioning but is also capable of producing detailed and comprehensive image descriptions. For our empirical study, we use FLINGO to generate rewards by matching image observations with specific text prompts such as “Find a blue cube”.

4 THE OTHER SIDE OF VISION-LANGUAGE MODELS FOR REINFORCEMENT LEARNING

In this section, we delve into the other side of utilizing VLMs for RL. More precisely, we show that reward hacking occurs in the context of RL agents trained with VLM-based rewards. To study this, we conduct thorough experiments across various domains that span different inputs (visual and state), different trained agents (DrQ-v2 (Yarats et al., 2022) and MP (Ahn et al., 2022)), and different task types (manipulation and navigation). Our new findings demonstrate that reward hacking is prevalent in all these scenarios.

4.1 EXPERIMENTAL SETUPS

To investigate the generalizability of reward hacking across various domains, we perform experiments on two types of tasks: manipulation and navigation. Manipulation requires an agent to perform actions that directly interact with the environment, while navigation only requires the agent to move

Environment	Agent	Input	Task
MetaWorld - Sweep Direction	DrQ-v2	Visual	Manipulation
MetaWorld - Sweep Direction	MP	State	Manipulation
House3D - Navigate	DrQ-v2	Visual	Navigation
House3D - Navigate	MP	State	Navigation

Table 1: Environments used for the experiments.

around. We use MetaWorld (Yu et al., 2020) for manipulation and House3D (Das et al., 2018) for navigation.

For the manipulation task in MetaWorld, we use the Sweep Direction task, where an agent is required to rotate a sweep arm to a target direction. The reward is a binary signal indicating whether the sweep arm is within (5, 10) of the target direction. For the navigation task in House3D, an agent is required to navigate to a target object. The reward is a binary signal indicating whether the agent is within 1 meter of the target. We use a door as the target for the navigation task in House3D.

We show that hacking can occur regardless of the agent’s input type. We use two types of agents for our experiments: (1) the DrQ-v2 agent (Yarats et al., 2022), which takes image observations as input, and (2) the MP agent (Ahn et al., 2022), which takes state observations as input. We use the implementation of the DrQ-v2 agent provided in the official DrQ codebase and the MP agent provided by the authors. For both types of agents, we train them using Human-designed Rewards and VLM-based Rewards. Additionally, we use MP as the base agent for our analysis of the underlying cause of reward hacking.

4.2 REWARD HACKING IN VARIOUS DOMAINS

In this section, we present the performance of RL agents trained with VLM-based rewards and human-designed rewards. Our findings show that, regardless of the environment, the input type, or the agent, hacking occurs.

Hacking occurs in manipulation tasks. The left two panels in Figure ?? show the performance of DrQ-v2 agents trained with VLM-based rewards and human-designed rewards in the MetaWorld Sweep Direction task. The results demonstrate that the agent trained with VLM-based rewards achieves lower performance with slower learning than the agent trained with human-designed rewards. We observe a similar result for the MP agent.

Hacking occurs in navigation tasks. The right two panels in Figure ?? show the performance of DrQ-v2 agents and the MP agents trained with VLM-based rewards and human-designed rewards in the House3D Navigate task. The results demonstrate that the agents trained with VLM-based rewards achieve lower performance with slower learning than the agents trained with human-designed rewards.

Hacking occurs with visual inputs and state inputs. Moreover, it is important to note that such hacking can occur with different types of observations. The second and third panels in Figure ?? show the performance of DrQ-v2 agents trained with human-designed and VLM-based rewards, using visual inputs. The rightmost panel shows the performance of the MP agent, which takes state inputs. The results demonstrate that hacking can occur with both visual and state inputs.

Hacking occurs in different languages. In Figure ??, we use an English text prompt for VLM-based rewards. We perform additional experiments to verify that hacking can also occur in different languages. More precisely, we change the language of the text prompt from English to Chinese. The results demonstrate that hacking can still occur in this scenario. More details are provided in Section 7.2.

In particular, we observe that the agent takes longer to learn policies for the sweep direction task than for the window open task. We attribute this to the fact that window open tasks tend to provide more visual cues than sweep direction tasks (e.g., whether a window is open or not), making it easier for the agent to hack.

4.3 EDGE-CASE HACKING

We have demonstrated in Figure ?? that RL agents trained with VLM-based rewards achieve significantly worse performance than those trained with human-designed rewards. These results show that hacking occurs in various distributions, where the observation can be either visual or state, and the task can be either navigation or manipulation. In this section, we further investigate whether hacking can occur in edge cases, which are rarity in the environment.

These tasks are designed such that they are deliberately away from the data distribution. We expect that human-designed reward functions should not affect the agent’s performance, regardless of how the text prompt is changed. However, as shown in Figure ??, the agents trained with VLM-based rewards struggle to solve these tasks. For instance, replacing “Find a red cube” with “Find a cube” (the second panel in the upper row of Figure ??) leads to a decrease in performance. This indicates that hacking can occur in edge cases.

4.4 REWARD HACKING CAN OCCUR WITH DIFFERENT PROMPTS

To further investigate the nature of reward hacking, we analyze whether it is possible to get a reasonable performance by modifying the text prompt. More precisely, we create a mixture of different tasks and test whether an agent can solve them by trying different text prompts.

We use the House3D environment, where an agent is required to navigate to a target object in a house. We use three different targets: a sofa, a chair, and a door. We use three text prompts: “Navigate to the sofa”, “Navigate to the chair”, and “Navigate to the door”. We then test whether an agent can solve the tasks by trying all possible combinations of prompts and targets. The results are shown in Figure ??.

According to Figure ??, we observe that an agent can solve “Navigate to the door” tasks regardless of the text prompt used. However, for other tasks, an interesting observation is that certain prompts can lead to good performance, while others cannot. For example, the best prompt for the “Navigate to the sofa” task is “Navigate to the sofa” prompt. However, changing the prompt to “Navigate to the chair” or “Navigate to the door” leads to a decrease in performance.

4.5 REWARD HACKING CAN LEAD TO ZERO PERFORMANCE

We have demonstrated that hacking can occur with various environments, input types, and prompts. However, a more alarming issue is that hacking can occur with an incorrect text prompt, causing the agent to completely fail to solve the original task.

To investigate this, we use the “Navigate to the chair” task in the House3D environment. We use the correct text prompt: “Navigate to the chair” as well as an incorrect text prompt: “Find a blue cube”. We then train an agent with an incorrect text prompt and test whether this harms the agent’s performance on the original task. As shown in Figure ??, an agent trained with an incorrect text prompt (“Find a blue cube”) for 1000 steps starts to fail on the original task (“Navigate to the chair”). This result demonstrates that hacking can occur with an incorrect text prompt can cause an agent to completely fail to solve the original task. This is of particular concern when training agents with VLM-based rewards as the agent may start to chase a seemingly reasonable goal and forget how to solve the original task.

4.6 REWARD HACKING IS NOT DUE TO NOISY REWARDS

It is worth noting that a common defense against reward hacking is that the generated reward is noisy (Skalse et al., 2022; Ma et al., 2023b). More precisely, some areas in the reward landscape may be ambiguous or even undefined. Consider the task of window open, the reward is ambiguous in the area where the window is half open. In these areas, the reward can take on any value (including zero) that is consistent with the undefined nature of the reward function. The proponents of this view argue that an agent can exploit these ambiguous areas by carefully choosing its actions, resulting in improper reward hacking.

However, our work reveals that even well-defined rewards can be hacked. To demonstrate this, we construct a point cloud environment where we manually set the reward function and train an agent

with a human-designed reward function and a VLM-based reward function. In this environment, the reward function is well-defined for all states. Our results demonstrate that hacking still occurs in this well-defined scenario. More details are provided in Section 7.2. This suggests that the generated reward may not be faithful to the underlying task.

4.7 THE GENERATED REWARD FUNCTION CANNOT BE FIXED WITH MORE DATA

In this section, we analyze whether collecting more data during RL training can mitigate reward hacking. To this end, we analyze how the performance of an RL agent changes with respect to the number of transitions collected. We show that even when the VLM-based reward function appears to be beneficial at the beginning of training, with more training data, the performance of an agent with VLM-based rewards falls behind the performance of an agent with human-designed rewards.

To investigate this, we use the sweep direction task in MetaWorld. We use a DrQ-v2 agent and plot the performance of the agent with respect to the number of transitions collected. The results are shown in Figure ???. According to Figure ??, at the beginning of training, the performance of an agent with VLM-based rewards is better than an agent with human-designed rewards. For example, the performance of an agent with VLM-based rewards is better than an agent with human-designed rewards when there are 100K transitions. However, the performance of an agent with VLM-based rewards starts to fall behind the performance of an agent with human-designed rewards when there are 200K transitions.

5 WHY DOES REWARD HACKING OCCUR?

We reveal that generated rewards suffer from unfaithfulness. More precisely, we demonstrate that they are not generalizable and they can be hacked by changing the text prompt.

5.1 FAILURE IN GENERATING ACCEPTABLE POLICIES

As illustrated in Figure ??, a faithful reward function should be able to generate an acceptable policy for the true objective. However, we show that this is not the case: RL agents with VLM-based rewards underperform, even after extensive training.

Performance is not consistent across different text prompts. We hypothesize that the generated reward may not be faithful to the underlying task. To investigate this, we change the text prompt used by the VLM-based reward function and evaluate the resulting performance of an RL agent.

According to Figure ??, the performance of an agent is not consistent across different text prompts. For example, replacing "Find a red cube" with "Find a blue cube" leads to a decrease in performance. This indicates that the generated reward function may not be faithful to the underlying task.

Performance is not consistent across different domains. We further demonstrate that the generated reward function may not be faithful to the underlying task by evaluating the performance of an RL agent in different domains. According to Figure ??, an agent trained with VLM-based rewards in the sweep direction environment achieves lower performance than an agent trained with human-designed rewards. However, an agent trained with VLM-based rewards in the navigate environment achieves performance that is comparable to an agent trained with human-designed rewards. This indicates that the generated reward function may not be faithful to the underlying task.

A faithful reward function should lead to a reasonable policy with extensive training. We further demonstrate that the generated reward function may not be faithful to the underlying task by evaluating the performance of an RL agent after extensive training. According to Figure ??, an agent trained with VLM-based rewards still achieves lower performance than an agent trained with human-designed rewards, even after extensive training. This indicates that the generated reward function may not be faithful to the underlying task.

5.2 FAILURE IN GENERALIZING TO SIMILAR TASKS

We have demonstrated that a faithful reward function should generate an acceptable policy and be consistent across different text prompts and domains. In this section, we further demonstrate that a faithful reward function should generalize to similar tasks.

To investigate this, we analyze whether an agent can solve complex tasks when trained with VLM-based rewards.

The results are shown in Figure ?? . We observe that agents trained with VLM-based rewards, even using complex tasks with language instructions or visual descriptions, still struggle to solve these tasks. For example, we observe that even after 4M steps, an agent trained with "Navigate to the sofa" reward still cannot solve the "Navigate to the chair" task. This indicates that the generated reward function may not be faithful to the underlying task and cannot generalize to similar tasks.

5.3 ANALYSIS OF WHY REWARD HACKING OCCURS

We have demonstrated that a faithful reward function should generate a reasonable policy and generalize to similar tasks. However, we show that these two properties are lacking in the context of RL agents trained with VLM-based rewards. In this section, we further analyze why reward hacking occurs.

The performance under the true reward gradually declines with training. We demonstrate that training an agent with VLM-based rewards can lead to hacking. To demonstrate this, we show that the performance of an agent under the true reward function gradually declines with training. According to Figure ?? , the performance of an agent trained with VLM-based rewards on the true reward function gradually declines as training progresses. For example, at the beginning of training, the performance of an agent trained with VLM-based rewards may be better than an agent trained with gap reward. However, with more training, the performance of an agent trained with VLM-based rewards falls behind the performance of an agent trained with gap reward. This indicates that training an agent with VLM-based rewards can lead to hacking.

The reward function can be hacked regardless of the training method. We further demonstrate that the generated reward function can be hacked regardless of the training method used. To demonstrate this, we train an agent with VLM-based rewards using SAC (Haarnoja et al., 2018), ACER (Wang et al., 2017), and IMPALA (Espeholt et al., 2018). The results are shown in Figure ?? . According to Figure ?? , the performance of an agent trained with VLM-based rewards is worse than the performance of an agent trained with human-designed rewards, even when using different training methods. This indicates that the generated reward function can be hacked regardless of the training method used.

6 CAN HACKING BE MITIGATED?

In this section, we discuss whether hacking can be mitigated. We propose two methods that are conceptually simple. However, as shown in Figure ?? , they fail to solve the problem of hacking. Possible Mitigation Methods are discussed in Section 7.1.

Adding a baseline reward To mitigate hacking, we can add a baseline reward to the VLM-based reward. More precisely, we add a gap reward with an initial value of 0 to the VLM-based reward. According to Figure ?? , while this approach leads to an improvement in performance compared to when no mitigation method is used, an agent still achieves lower performance with this method than with the human-designed reward.

Confining the reward space To mitigate hacking, we can confine the reward space by clipping the VLM-based reward. More precisely, we clip the VLM-based reward to be between a negative number and a positive number. According to Figure ?? , while this approach leads to an improvement in performance compared to when no mitigation method is used, an agent still achieves lower performance with this method than with the human-designed reward.

7 CONCLUSION AND DISCUSSION

In this work, we reveal a potentially significant limitation of utilizing Vision Language Models (VLMs) for RL: the issue of reward hacking can occur in the context of RL agents trained with VLM-based rewards. This phenomenon happens across various domains, affecting both navigation and manipulation tasks, and occurs with different input types (including visual and state inputs) and RL agents (including DrQ-v2 and the MP agent). We further demonstrate that reward hacking is due to unfaithful reward functions and can lead to zero performance. Finally, we analyze why hacking occurs and discuss potential mitigation methods. Our findings underscore the need for increased vigilance and further research when deploying such methods in real-world applications.

7.1 POSSIBLE MITIGATION METHODS

In this section, we discuss potential mitigation methods for hacking. We propose two methods: adding a baseline reward and confining the reward space. We also discuss why these methods fail. Based on our analysis, we believe that mitigating hacking in the context of using VLMs to train RL agents is challenging.

Adding a baseline reward One potential mitigation method is to add a baseline reward to the VLM-based reward. The baseline reward can be a constant or a value that is computed based on the agent’s observation. However, this method is not without its limitations. For example, setting the baseline reward too high can prevent the agent from hacking. However, this can make the learning problem more difficult. As demonstrated in Figure ??, while this method leads to an improvement in performance compared to when no mitigation method is used, it fails to match the performance of an agent trained with a human-designed reward.

Confining the reward space Another potential mitigation method is to confine the reward space by clipping the VLM-based reward. The reward space can be confined to be positive, or it can be confined to be a narrow range. However, this method is not without its limitations. For example, setting the range too narrow can prevent the agent from hacking. However, this can make the learning problem more difficult. As demonstrated in Figure ??, while this method leads to an improvement in performance compared to when no mitigation method is used, it fails to match the performance of an agent trained with a human-designed reward.

In conclusion, our findings suggest that mitigating hacking in the context of using VLMs to train RL agents is challenging.

7.2 EXPERIMENTAL DETAILS

In this section, we provide additional experimental details.

7.3 ENVIRONMENT DESCRIPTION

MetaWorld is a simulated robotic hand-aware manipulation environment built upon the MuJoCo physics engine. It includes a variety of tasks that require an agent to manipulate objects in a scene in order to achieve a specific objective. In this work, we use the sweep direction task, which requires an agent to rotate a sweep arm to a target direction. The reward for this task is a binary signal that indicates whether the sweep arm is within $(5, 10)$ of the target direction.

House3D is a simulated environment that contains several objects such as chairs, sofas, tables, and lamps. It is similar to the EmbodiedQA environment (Das et al., 2018) and VQA (Antol et al., 2015), which are both based on the same codebase. The navigation task in House3D requires an agent to navigate to a target object. The reward for this task is a binary signal that indicates whether the agent is within 1 meter of the target. In our experiments, the target is a door.

7.4 IMPLEMENTATION DETAILS

Image Processing For image-based input, we resize the observation to have a shape of $(64, 64, 3)$.

Text Processing For the prompt, we use “Find a [object]” or “Navigate to [object]”. For “Find a [object]”, we use “Find a red cube”, “Find a blue cube”, “Find a cube with one red cube”, or “Find a

cube with two red cubes". For "Navigate to [object]", we use "Navigate to the sofa", "Navigate to the chair", or "Navigate to the door".

Reward Processing We use a CLIP model to train RL agents. We reshape the similarity score output by CLIP to serve as a reward. For the navigation task in House3D, we only use the text prompt "Navigate to the door". For the sweep direction task in MetaWorld, we only use the text prompt "Find a red cube". This is consistent with the main text.

RL Training We use the official DrQ-v2 codebase and MP agent provided by the authors for our experiments. We use gap reward to train RL agents. The gap reward is defined as $R_{gap} = R_{VLM} - \lambda R_{human}$, where R_{VLM} is the reward generated by the VLM, R_{human} is the human-designed reward, and λ is a hyperparameter that is selected from $\{0.5, 0.7, 0.9, 1.0, 1.1, 1.3\}$. We perform RL training for 2M steps in MetaWorld and 5M steps in House3D.

8 LIMITATIONS AND FUTURE WORK

Limitation 1: Scope of Analysis Our analysis is limited to complex tasks with sparse rewards. We show that RL agents can hack VLM-based rewards in complex tasks with sparse rewards. However, we have not studied whether hacking can occur in other scenarios. For example, we have not studied whether hacking can occur in tasks with dense rewards.

Limitation 2: Scope of Experiments Our experiments are conducted in simulated environments. We have shown that RL agents can hack VLM-based rewards in simulated environments. However, we have not studied whether hacking can occur in other scenarios. For example, we have not studied whether hacking can occur in real-world environments.

Future Work 1: Broaden the Scope of Analysis To deepen our understanding of hacking, it is crucial to expand the scope of our analysis. Specifically, future works can delve into other types of tasks (such as those with dense rewards), environments (such as real-world environments), and agents. This exploration can offer a more comprehensive perspective on the prevalence and nature of hacking in various RL contexts.

Future Work 2: Develop Advanced Methods to Leverage the Capabilities of VLMs in RL Our findings regarding hacking underscore the need for more advanced methods to leverage the capabilities of VLMs in RL. While mitigating hacking remains challenging as of now, future research directions could explore advanced training algorithms that can better harness VLMs without falling victim to hacking. This is an open question and presents a promising avenue for future research.

REFERENCES

- Ademi Adeniji, Amber Xie, Carmelo Sferrazza, Younggyo Seo, Stephen James, and Pieter Abbeel. Language reward modulation for pretraining reinforcement learning. *arXiv:2308.12270*, 2023.
- Forest Agostinelli, Stephen McAleer, Alexander Shmakov, and Pierre Baldi. Solving the rubik’s cube with deep reinforcement learning and search. *Nature Machine Intelligence*, 1:356–363, 2019.
- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. Do as I can, not as I say: Grounding language in robotic affordances. In *Conference on Robot Learning*, 2022.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikołaj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. Flamingo: a visual

- language model for few-shot learning. In *Advances in Neural Information Processing Systems*, 2022.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. VQA: Visual question answering. In *IEEE International Conference on Computer Vision*, 2015.
- Neeloy Chakraborty, Melkior Ornik, and Katherine Driggs-Campbell. Hallucination detection in foundation models for decision-making: A flexible definition and review of the state of the art. *arXiv:2403.16527*, 2024.
- Harris Chan, Volodymyr Mnih, Feryal Behbahani, Michael Laskin, Luyu Wang, Fabio Pardo, Maxime Gazeau, Himanshu Sahni, Dan Horgan, Kate Baumli, Yannick Schroecker, Stephen Spencer, Richie Steigerwald, John Quan, Gheorghe Comanici, Sebastian Flennerhag, Alexander Neitz, Lei M Zhang, Tom Schaul, Satinder Singh, Clare Lyle, Tim Rocktäschel, Jack Parker-Holder, and Kristian Holsheimer. Vision-language models as a source of rewards. In *Second Agent Learning in Open-Endedness Workshop*, 2023.
- William Chen, Oier Mees, Aviral Kumar, and Sergey Levine. Vision-language models provide promptable representations for reinforcement learning. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
- Xuzhe Dang, Stefan Edelkamp, and Nicolas Ribault. CLIP-Motion: Learning reward functions for robotic actions using consecutive observations. *arXiv:2311.03485*, 2023.
- Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- Yuqing Du, Ksenia Konyushkova, Misha Denil, Akhil Raju, Jessica Landon, Felix Hill, Nando de Freitas, and Serkan Cabi. Vision-language models as success detectors. *arXiv:2303.07280*, 2023.
- Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. IMPALA: Scalable distributed deep-RL with importance weighted actor-learner architectures. In *International conference on machine learning*, 2018.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, 2018.
- Yafei Hu, Quanting Xie, Vidhi Jain, Jonathan Francis, Jay Patrikar, Nikhil Keetha, Seungchan Kim, Yaqi Xie, Tianyi Zhang, Shibo Zhao, Yu Quan Chong, Chen Wang, Katia Sycara, Matthew Johnson-Roberson, Dhruv Batra, Xiaolong Wang, Sebastian Scherer, Zsolt Kira, Fei Xia, and Yonatan Bisk. Toward general-purpose robots via foundation models: A survey and meta-analysis. *arXiv:2312.08782*, 2023.
- Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. In *Advances in Neural Information Processing Systems*, 2019.
- Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, and Sergey Levine. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, 2018.
- Martin Klissarov, Pierluca D’Oro, Shagun Sodhani, Roberta Raileanu, Pierre-Luc Bacon, Pascal Vincent, Amy Zhang, and Mikael Henaff. Motif: Intrinsic motivation from artificial intelligence feedback. *arXiv:2310.00166*, 2023.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv:2305.10355*, 2023.

- Ekdeep Singh Lubana, Johann Brehmer, Pim de Haan, and Taco Cohen. FoMo rewards: Can we cast foundation models as reward functions? *arXiv:2312.03881*, 2023.
- Yecheng Jason Ma, William Liang, Vaidehi Som, Vikash Kumar, Amy Zhang, Osbert Bastani, and Dinesh Jayaraman. LIV: Language-image representations and rewards for robotic control. *arXiv:2306.00958*, 2023a.
- Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy Zhang. VIP: Towards universal visual reward and representation via value-implicit pre-training. In *International Conference on Learning Representations*, 2023b.
- Parsa Mahmoudieh, Deepak Pathak, and Trevor Darrell. Zero-shot reward specification via grounded natural language. In *International Conference on Machine Learning*, 2022.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3M: A universal visual representation for robot manipulation. In *Conference on Robot Learning*, 2023.
- Taewook Nam, Juyong Lee, Jesse Zhang, Sung Ju Hwang, Joseph J. Lim, and Karl Pertsch. Lift: Unsupervised reinforcement learning with foundation models as teachers. *arXiv:2312.08958*, 2023.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, 2021.
- Juan Rocamonde, Victoriano Montesinos, Elvis Nava, Ethan Perez, and David Lindner. Vision-language models are zero-shot reward models for reinforcement learning. *arXiv:2310.12921*, 2023a.
- Juan Rocamonde, Victoriano Montesinos, Elvis Nava, Ethan Perez, and David Lindner. Vision-language models are zero-shot reward models for reinforcement learning. *arXiv:2310.12921*, 2023b.
- David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- Joar Skalse, Nikolaus H. R. Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and characterizing reward hacking. *arXiv:2209.13085*, 2022.
- Theodore Sumers, Kenneth Marino, Arun Ahuja, Rob Fergus, and Ishita Dasgupta. Distilling internet-scale vision-language models into embodied agents. In *International Conference on Machine Learning*, 2023.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Ziyu Wang, Victor Bapst, Nicolas Heess, Volodymyr Mnih, Remi Munos, Koray Kavukcuoglu, and Nando de Freitas. Sample efficient actor-critic with experience replay. In *International Conference on Learning Representations*, 2017.
- Jia-Yu Yao, Kun-Peng Ning, Zhen-Hui Liu, Mu-Nan Ning, and Li Yuan. LLM lies: Hallucinations are not bugs, but features as adversarial examples. *arXiv:2310.01469*, 2023.
- Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous control: Improved data-augmented reinforcement learning. In *International Conference on Learning Representations*, 2022.

Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning*, 2020.

Haichao Zhang and Jianyu Wang. Defense against adversarial attacks using feature scattering-based adversarial training. In *Advances in Neural Information Processing Systems*, 2019.

Generated by CycleResearcher