

OCALM: Object-Centric Assessment with Language Models

Anonymous authors

Paper under double-blind review

Abstract

Properly defining a reward signal to efficiently train a reinforcement learning (RL) agent is a challenging task. Designing balanced objective functions from which a desired behavior can emerge requires expert knowledge, especially for complex environments. Learning rewards from human feedback or using large language models (LLMs) to directly provide rewards are promising alternatives, allowing non-experts to specify goals for the agent. However, black-box reward models make it difficult to debug the reward. In this work, we propose Object-Centric Assessment with Language Models (OCALM) to derive inherently interpretable reward functions for RL agents from natural language task descriptions. OCALM uses the extensive world-knowledge of LLMs while leveraging the object-centric nature common to many environments to derive reward functions focused on relational concepts, providing RL agents with the ability to derive policies from task descriptions.

1 Introduction

Defining reward functions for reinforcement learning (RL) agents is a notoriously challenging task (Amodei et al., 2016; Knox et al., 2023; Delfosse et al., 2024). Consequently, reward functions are often unavailable or sub-optimal, suffering from issues such as reward sparsity (Andrychowicz et al., 2017) or difficult credit assignment (Raposo et al., 2021; Wu et al., 2023). While standard RL benchmark are equipped with predefined reward functions, real-world tasks typically lack explicit reward signals. Existing approaches, such as reinforcement learning from human feedback (RLHF) (Christiano et al., 2017; Ouyang et al., 2022; Kaufmann et al., 2023), circumvent the reward specification problem by learning a reward model from human feedback. However, it generally requires to learn reward models from scratch, which can be slow and inefficient. Further, the *black box* nature of the reward model complicates the understanding and adjustment of the given reward signal.

In contrast to RL agents, humans can learn to solve tasks without clear external rewards, deriving their own objectives from task context (Deci & Ryan, 2013) (*cf.* Figure 1). Given such context, humans formulate their own goals and generate a corresponding reward signal autonomously (Spence, 1947; Oudeyer & Kaplan, 2008). This capability stems from our rich understanding of the world, enabling us to derive specific goals from potentially vague task descriptions. Conversely, RL agents typically lack common sense and are trained tabula rasa, devoid of any world knowledge. In this paper, we demonstrate that large language models (LLMs) are capable of a similar feat, using their acquired world knowledge to derive goals from task descriptions that can be used by RL agents.

While previous works have demonstrated that LLMs can provide RL agents with a reward signal derived from context (Ma et al., 2024; Xie et al., 2024), these approaches do not capitalize on the object-centric and relational nature prevalent in environments that incorporate relational reasoning challenges. Assuming object-centricity offers a powerful inductive bias, enabling agents to reason about the world in terms of objects and their interactions rather than through raw pixels or other low-level features (Delfosse et al., 2023b; Luo et al., 2024). We demonstrate that by directing the LLM to concentrate on the relationships between objects we can significantly enhance the effectiveness of the generated reward functions and, consequently, improve the final agent’s performance.

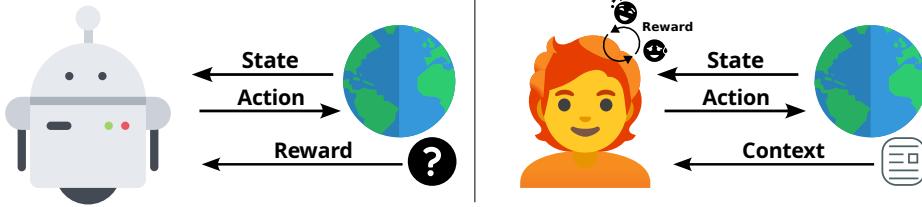


Figure 1: **Contrary to RL agents, humans infer objectives from context.** The RL setting assumes the existence of an external reward function, whereas humans are able to infer rewards from information about the environment and task context.

We introduce Object-Centric Assessment with Language Models (OCALM, *cf.* Figure 2) as an approach to derive inherently interpretable reward functions for RL agents from the natural-language context of tasks. OCALM leverages both the extensive world-knowledge of LLMs and the object-centric nature of many environments to equip RL agents with a rich understanding of the world and the ability to derive goals from task descriptions. We leverage the powerful inductive bias of object-centric reasoning, directing the LLM to focus on the relationships between objects in the environment using a multi-turn interaction. OCALM comprises two main components: (1) a language model that generates a symbolic reward function from text-based task context, and (2) an RL agent that trains based on this derived reward function.

In our evaluations on the iconic Atari Learning Environment (ALE) (Mnih et al., 2013), we provide experimental evidence of OCALM’s performance, particularly its capability to learn policies comparable to those of agents trained with ground-truth reward functions. We demonstrate the benefits of object-centric reasoning through the relational inductive bias, which significantly enhances the quality of the learned reward functions. Additionally, we highlight the interpretability of the learned reward functions and OCALM’s applicability to environments lacking ground-truth rewards.

In summary, our specific contributions are:

- (i) We introduce OCALM, an approach for inferring relational (*object-centric*) reward functions from text-based task descriptions for RL agents.
- (ii) We show that OCALM produces learnable reward functions, that lead to RL agents performing on par with agents trained on the original reward.
- (iii) We empirically demonstrate the importance of object-centric reasoning for enhancing the performance of OCALM.
- (iv) We establish that OCALM provides inherently interpretable reward functions.

In the remainder of the paper, we provide a detailed description of OCALM and its components (Section 2), followed by experimental evaluations and analysis (Section 3). We address related work (Section 5) before concluding (Section 6).

2 Object-Centric Assessment with Language Models

OCALM provides RL agents with inherently interpretable reward functions derived from text-based task descriptions. We follow a multistep approach, as depicted in Figure 2, to achieve this goal.

(1) Context definition. We start by gathering a *natural-language task description* and extracting an *object-centric state abstraction* from the raw input state. The task descriptions (listed in Appendix A.4) are based on the short descriptions of each atari environment (Towers et al., 2023), slightly modified to add missing information. The object-centric state abstractions include the properties of each object, such as their class, position, size, and color. It is given by the classes provided by the OCAtari framework (Delfosse et al., 2023a), *i.e.*, the parent game object class and the game-specific objects (examples are provided in Appendix A.3.2). Game objects related to the score were

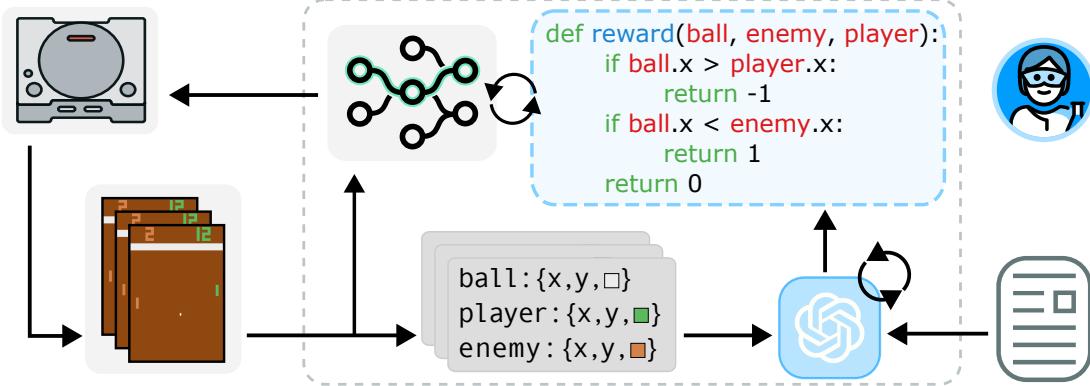


Figure 2: Object-Centric Assessment with Language Models. OCALM extracts a *neurosymbolic abstraction* from the raw state, provided to a language model together with the game’s context, to generate a *symbolic reward function* (in python). The language model first generates relational utility functions, that are then used in the reward function. This transparent reward can be inspected and used to train the policy.

omitted, since we assume a reward-free environment. The task description and the object-centric state abstraction form the task context, which is provided to the language model.

(2) LLM-driven reward generation. The large language model (LLM) processes the task context to generate a symbolic reward function in the form of Python source code. We use a guided multi-turn approach to direct the LLM to focus on the relationships between objects in the environment, similar to chain-of-thought reasoning (Wei et al., 2022).

(2.1) Relational concept extraction. The LLM is tasked with generating relational functions that describe the relationships between objects in the environment (*cf.* Listing 4, Appendix A.3), which are important to understand the game states.

(2.2) Reward generation. Given the task context and the created utility functions, the LLM generates a symbolic reward function in the form of Python source code (prompt given in Listing 5 in Appendix A.3).

(2.3) Reward scaling. As a last step, we prompt the LLM (Listing 6, Appendix A.3) to adjust the created reward function in such a way that the rewards are on a scale from -1 to 1 .

The resulting reward function takes the form of a Python function that maps the object-centric state abstraction to a scalar reward, that includes semantic descriptions. The code is inherently interpretable and allows experts to inspect and verify the reward function before proceeding with the next step. We also present an ablated version of OCALM, calling it OCALM (no relations), where the LLM directly generates the reward function without the relational and reward scaling steps. While the no relations version may still use relational concepts, we do not prompt it to do so. We use a modified prompt (Listing 3 in Appendix A.3) in that case. We provide a shortened example of a generated reward function in Listing 10 and the full reward functions in Appendix A.4.

(3) Policy training. The derived reward function is used to train an RL agent, which learns a policy that maximizes the reward. The agent can be trained using any conventional RL algorithm, Proximal Policy Optimization (PPO, Schulman et al., 2017) in our experiments.

3 Experimental Evaluation

We evaluate the OCALM approach to answer the following research questions:

(Q1) Does OCALM generate rewards that correspond to learnable tasks?

(Q2) Can OCALM agents master Atari environments without access to the true game score?

Listing 1: An example reward function generated by OCALM (full). Implementation of relational utility function elided and unused utilities removed. The full version is in Listing 10.

```

1 def detect_collision(chicken, car): [...]
2 def has_reached_top(chicken, screen_height): [...]
3 def progress_made(chicken, screen_height): [...]
4
5 def reward_function(game_objects) -> float:
6     # Initialize reward
7     reward = 0.0
8
9     # Constants
10    SCREEN_HEIGHT = 160
11    COLLISION_PENALTY = -1.0 # Scaled down to fit within [-1, 1]
12    PROGRESS_REWARD = 0.1 # Scaled down to incrementally increase reward
13    SUCCESS_REWARD = 1.0 # Maximum reward for reaching the top
14
15    # Filter out chickens and cars from game_objects
16    chickens = [obj for obj in game_objects if isinstance(obj, Chicken)]
17    cars = [obj for obj in game_objects if isinstance(obj, Car)]
18
19    # Assume control of the leftmost chicken (player's chicken)
20    if chickens:
21        player_chicken = min(chickens, key=lambda c: c.x)
22
23        # Check if the chicken has reached the top
24        if has_reached_top(player_chicken, SCREEN_HEIGHT):
25            reward += SUCCESS_REWARD
26
27        # Reward based on progress towards the top
28        reward += progress_made(player_chicken, SCREEN_HEIGHT) * PROGRESS_REWARD
29
30        # Check for collisions with any car
31        for car in cars:
32            if detect_collision(player_chicken, car):
33                reward += COLLISION_PENALTY
34                break # Only penalize once per time step
35
36    # Ensure reward stays within the range [-1, 1]
37    reward = max(min(reward, 1.0), -1.0)
38
39    return reward

```

(Q3) How does relation-focused reward derivation influence performance and interpretability?

(Q4) How interpretable are the reward functions generated by OCALM agents?

Experimental setup: We evaluate OCALM on four Atari games (Pong, Freeway, Skiing, and Seaquest) from the Atari Learning Environments (ALE) (Bellemare et al., 2013), the most used RL framework (Delfosse et al., 2023a). All results are averaged over three seeds for each agent configuration, with standard deviation indicated. We use exponential moving average (EMA) smoothing for our plots, with an effective window size of 50, resulting in a smoothing factor $\alpha = 2/(1+50) \approx 0.039$. Note that rewards are not always reported at the same timestep, in which case we ignore missing values when computing the average, relying on the EMA smoothing to provide a avoid excessive fluctuations and noise. We use Proximal Policy Optimization (PPO, Schulman et al., 2017) as the base architecture due to its success in Atari games. The input representation is a stack of four gray-scaled 84×84 images, introduced by Mnih et al. (2015). All agents are trained using 10M frames with the implementation by Huang et al. (2022) and default hyperparameters (*cf.* Appendix A.1).

We compare our OCALM agents trained with the ‘true’ reward functions given by the ALE environment, typically based on game score. Both types of agents are evaluated against the true game score. All evaluations use the latest *v5* version of the ALE environments, following best-practices

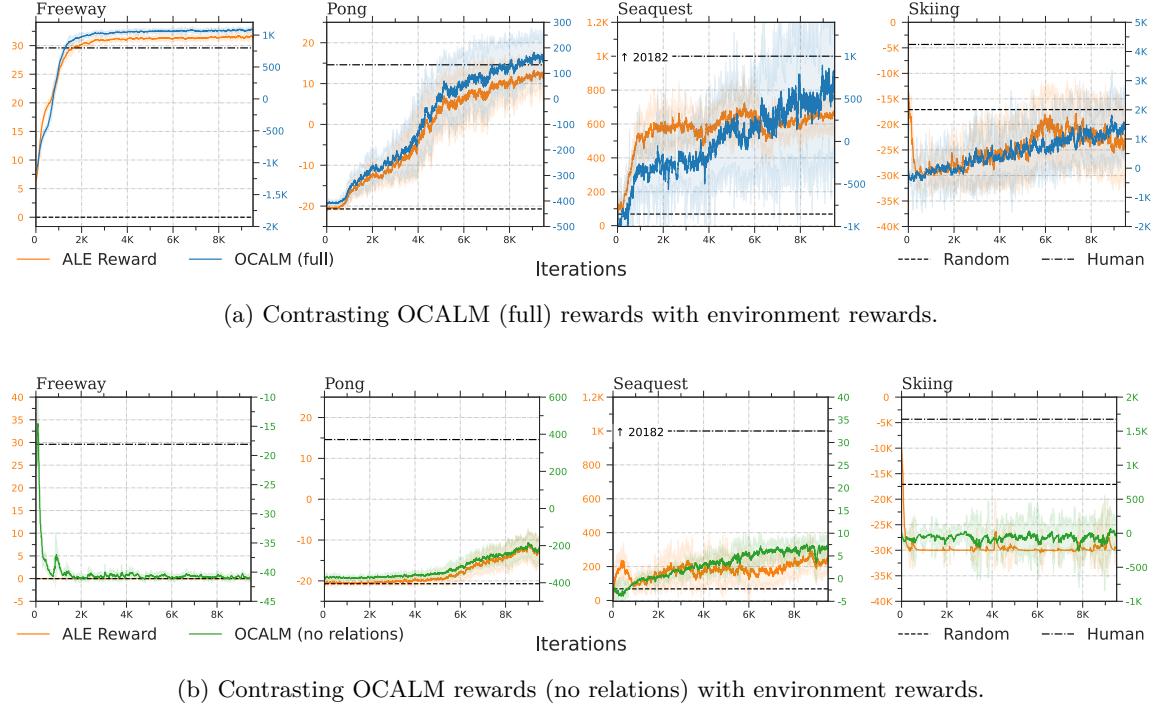


Figure 3: **OCALM generates meaningful reward functions that correlate with the intended game rewards.** These figures show the performance of agents trained on OCALM-derived rewards, measured on both the OCALM-derived reward and the environment reward. The scales of rewards differ, therefore the axes are scaled to better visualize the correlation. Both plot for the same game share the same axis range for better comparability. The results indicate that (1) the reward functions generated by OCALM correspond to objectives learnable by an RL agent, and (2) the OCALM-derived rewards correlate with the environment rewards. All experiments were averaged over 3 seeds, with standard deviations shown as shaded areas.

to prevent overfitting (Machado et al., 2018). The results are presented as figures here, refer to Appendix A.2 for numerical results.

To generate our reward function, we assume access to object-centric state descriptions of the game state. To focus on description-based reward derivation, we use representations from the Object-Centric Atari (OCAtari) framework (Delfosse et al., 2023a). While a learned object detector could extract objects from raw input (Redmon et al., 2016; Lin et al., 2020), we use OCAtari for simplicity.

OCALM generates reward signals allowing to master the game (Q1). We first test whether OCALM generates rewards that correspond to learnable tasks. For this purpose, we track the learning curves of agents trained on OCALM-derived rewards and verify that agents improve over time, e.g., learn to maximize the reward. Figure 3 shows that this is generally the case, with an exception for Freeway when using the ablated variant of OCALM (no relations) (see Figure 3b). For all other games, and for all games when using the full OCALM pipeline, the agents improve over time when measured on the OCALM-derived reward.

Without the relational inductive bias, OCALM fails to generate learnable rewards for Freeway. This is due to a bug in the generated reward function (see Appendix A.4.1, Listing 11), which fails to identify the player-controlled chicken. Although it is quite possible that the relational inductive bias helps to avoid such bugs through mechanisms similar to chain-of-thought reasoning, they cannot entirely be prevented. More research is necessary to understand the impact of the

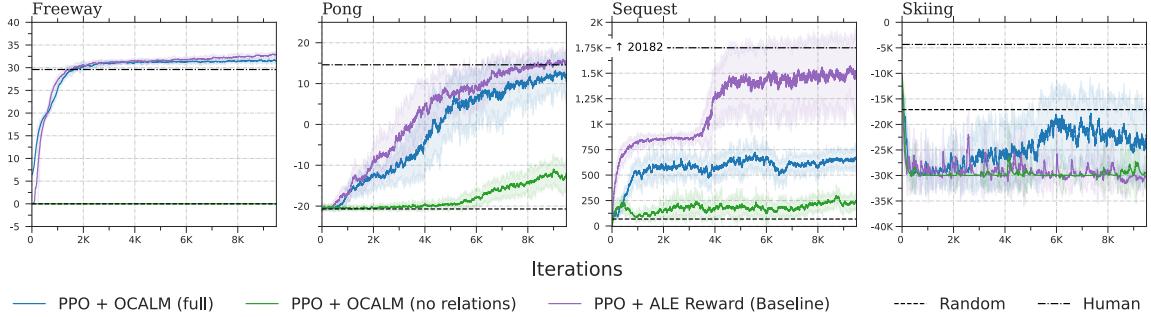


Figure 4: OCALM agents can master different Atari environments. Comparing the performance of agents trained on OCALM-derived rewards to agents trained on the true game score. All experiments were averaged over 3 seeds, with standard deviations shown as shaded areas.

relational inductive bias on the failure rate of generated reward functions. Iterative refinement could help further alleviate this issue, but generating successful reward functions in a single shot remains a significant computational advantage.

OCALM-based agents can master different Atari environments without access to the true game score (Q2). Figure 4 shows the performance of agents trained on OCALM-derived rewards compared to those trained on the true game score. Performance is measured on the true game score in both cases, which OCALM agents cannot access during training. Our goal is not to exceed the baseline agents' performance, but to show that OCALM agents can master environments without access to true rewards. Even though the OCALM-derived reward functions differ from the environment reward, we observe that OCALM agents, when using relational prompting, reliably improve their performance over the course of training when measured on the true game score. This further confirms that the reward functions generated by OCALM are correlated with the true game score, as discussed in the previous paragrpah. For Freeway and Seaquest in particular, OCALM agents were able to reach competitive performance compared to the baseline agents, without requiring access to the true game score. Although OCALM agents do not match the baseline's performance in Pong, they still show significant learning progress, again indicating the reward function generated by OCALM correlates with the environment reward.

Relational prompting of OCALM agents improves reward quality (Q3). Figure 4 shows that agents trained on OCALM-derived rewards with the relational inductive bias (denoted OCALM (full)) generally outperform those without it (denoted OCALM (no relations)). This is particularly evident in Freeway and Seaquest, where OCALM (full) agents reach performance competitive with the baseline, while OCALM (no relations) agents fail to learn the task. In Pong and Skiing, OCALM (full) agents perform equivalently to OCALM (no relations) agents, indicating the relational inductive bias is not beneficial in all cases, but also does not harm performance. Note that the OCALM (no relations) variant also skips the reward scaling step, which could be another contributing factor to the performance difference. Qualitatively, when inspecting the reward functions generated by OCALM (Appendix A.4), we observe that the relational inductive bias helps to capture more complex concepts, such as the distance to the nearest obstacle in Skiing Appendix A.4.4, which in turn can lead to better-shaped reward functions.

OCALM generates interpretable reward functions (Q4). The reward functions generated by OCALM (*cf.* Appendix A.4) are based on high-level objects and relations, documented with comments, making them easy to interpret and understand. Relational prompting further aids in generating interpretable reward functions by introducing easy-to-understand relational concepts, which add an abstraction layer to the reward function. Examples include collision detection, a relation generated for all games, and easily understandable concepts such as `has_reached_top` in Freeway (Appendix A.4.1), and more complex relations including multiple objects such as `detect_score_event` in Pong (Appendix A.4.2) or `check_gate_passage` in Skiing (Appendix A.4.4).

Table 1: Relating OCALM to the most closely related work, EUREKA (Ma et al., 2024) and Text2Reward (Xie et al., 2024). In contrast to our work, EUREKA and Text2Reward use multiple iterations to refine the reward function. All three approaches provide the LLM with additional information about the environment and the task. While EUREKA and Text2Reward evaluate on joint control tasks (locomotion and manipulation), we evaluate on relational tasks (Atari games).

Approach	1-Shot	Add. Context	Relational	Evaluation
EUREKA	No	Source code	No	Joint control
Text2Reward	No	Symb. state abstr.	No	Joint control
OCALM (ours)	Yes	Symb. state abstr.	Yes	Relational tasks

4 Limitations

In our evaluations, we use the integrated object extractor of OCAtari which provides ground truth data. Such extractors can also be optimized using supervised (Redmon et al., 2016) or self-supervised (Lin et al., 2020; Delfosse et al., 2023c) object detection methods. We additionally rely on the language models ability to generate a reward function in a single shot. While our relational inductive bias helps, the LLM may miss crucial information such as the frequency of certain events, which is important to tune the relative scales of different reward components. Related works (Ma et al., 2024; Xie et al., 2024) rely on iterative refinement of the reward function, which could further enhance OCALM’s performance. Nonetheless, the relational inductive bias enables OCALM to frequently learn successful reward functions in a single shot, a significant computational efficiency advantage.

5 Related Work

OCALM lies at the intersection of several research areas, including reinforcement learning from human feedback (RLHF), language-guided RL, explainable RL and relation extraction. We discuss the most relevant work in these areas below.

Reward learning has been studied in various forms and based on different sources, such as demonstrations (Arora & Doshi, 2021) and human preferences (Kaufmann et al., 2023). While these approaches can be very effective, they often require a large amount of human supervision, which can be costly and time-consuming. Our method, by combining human guidance given through the task description with the extensive world knowledge of LLMs, helps to alleviate this issue. Particularly closely related to ours, RL-VLM-F (Wang et al., 2024) is a notable approach that learns a reward model from pairwise comparisons judged by a vision-language model based on a natural language task description. Similar to OCALM, this leverages the prior knowledge of the vision-language model. In contrast to our work, however, Wang et al. (2024) and most other reward learning methods learn black-box reward models in the form of neural networks, which are not interpretable.

LLM-written reward functions have been studied by Ma et al. (2024), who propose EUREKA, and Xie et al. (2024), who propose Text2Reward. This is the most closely related work to ours. Table 1 highlights the most relevant differences between our method and theirs. Like our approach, EUREKA and Text2Reward use an LLM to generate reward functions for RL agents. They assume access to a natural language task description and an environment description, specifying the structure of observations. EUREKA assumes descriptions are given in the form of incomplete source code, while Text2Reward requires class definitions that define the components of the state. Both approaches evaluate the generated reward functions on robotic manipulation and locomotion tasks.

EUREKA and Text2Reward work in an iterative setting, where the reward function is refined based on feedback from the environment or a human expert. This can help further improve the reward function, but also requires more computational resources, time and supervision (either from a human expert or a success signal). Since the focus of our study is on the benefits of object-centric reward

specifications, we leave the iterative refinement for future work and instead focus on improving single-shot performance.

In contrast to these prior works, we focus on relational reasoning environments, which require the agent to reason about multiple objects and their interactions. We leverage a relational inductive bias for improved one-shot performance, reducing the need for iterative refinement and human supervision. We further evaluate on the prominent Atari Learning Environment (Bellemare et al., 2013), the most used benchmark for reinforcement learning Delfosse et al. (2023a), and show the importance of object-centric inductive biases for learning reward functions in this setting.

RL from natural language task descriptions has been studied extensively, e.g., by Nair et al. (2021). While these approaches are similar to ours in that they use natural language to specify the task, they typically do not leverage the world knowledge of LLMs, do not learn interpretable reward functions, and do not use relational inductive biases.

Explainable RL (XRL) is subfield of explainable AI (XAI) (Milani et al., 2023; Dazeley et al., 2023; Krajna et al., 2022). XRL aims to offer insights into the behavior of RL agents, aiding in realigning agents. OCALM helps in this endeavor by providing inherently interpretable reward functions, which can be inspected and verified by experts. This can be used to align the reward functions with certain societal values, such as more pacific gameplays in *e.g.* shooting games.

Relation extraction has been studied in many forms, including prior task-knowledge integration (Reid et al., 2022) and neural guidance from a pretrained fully deep agent (Delfosse et al., 2023b). In contrast to these works, we use an LLM to extract relations between objects in the environment, which are then used to derive reward functions. Particularly relevant is the work by Wu et al. (2023), who extract relevant relations using LLMs with access to an instruction manual. This differs from our work in that they use the extracted relations to supplement existing rewards instead of entirely replacing the environment reward function.

6 Conclusion

We have presented OCALM, a novel approach for deriving inherently interpretable reward functions for RL agents from natural language task descriptions. Our method leverages the extensive world knowledge of LLMs and the object-centric, relational nature of the environment to generate symbolic reward functions that can be inspected and verified by experts. We have shown that OCALM can be used to train RL agents on Atari games, demonstrating that the derived reward functions are effective in guiding the agent to learn the desired behavior. OCALM agents utilize the abstracted knowledge of LLMs alongside explicit relational concepts to derive effective and inherently interpretable reward functions for complex RL environments.

Broader Impact Statement

We here develop RL agents with transparent, human-defined objectives, improving RL accessibility to non-experts. We thus reduce the barrier to entry for non-experts, helping to ensure that the objectives of the agents are aligned with the user’s intentions. A malicious user can, however, utilize such approaches for developing agents with harmful objectives, thereby potentially leading to a negative impact on further users or society as a whole. Even so, the inspectable nature of transparent approaches will allow to identify such potentially harmful misuses, or hidden misalignment.

References

- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete Problems in AI Safety, 2016. arXiv preprint.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight Experience

Replay. In *Advances in Neural Information Processing Systems (NIPS)*, volume 30. Curran Associates, Inc., 2017.

Saurabh Arora and Prashant Doshi. A Survey of Inverse Reinforcement Learning: Challenges, Methods and Progress. *Artificial Intelligence*, 297:103500, 2021. doi: [10.1016/j.artint.2021.103500](https://doi.org/10.1016/j.artint.2021.103500).

M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The Arcade Learning Environment: An Evaluation Platform for General Agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013. doi: [10.1613/jair.3912](https://doi.org/10.1613/jair.3912).

Paul Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep Reinforcement Learning from Human Preferences. In *Advances in Neural Information Processing Systems (NIPS)*, volume 30. Curran Associates, Inc., 2017.

Richard Dazeley, Peter Vamplew, and Francisco Cruz. Explainable reinforcement learning for broad-XAI: A conceptual framework and survey. *Neural Computing and Applications*, 35(23):16893–16916, 2023. doi: [10.1007/s00521-023-08423-1](https://doi.org/10.1007/s00521-023-08423-1).

Edward L. Deci and Richard M. Ryan. *Intrinsic Motivation and Self-Determination in Human Behavior*. Springer Science & Business Media, 2013. ISBN 978-1-4899-2271-7.

Quentin Delfosse, Jannis Blüml, Bjarne Gregori, Sebastian Sztwiertnia, and Kristian Kersting. OCAtari: Object-Centric Atari 2600 Reinforcement Learning Environments, 2023a. arXiv preprint.

Quentin Delfosse, Hikaru Shindo, Devendra Dhami, and Kristian Kersting. Interpretable and Explainable Logical Policies via Neurally Guided Symbolic Abstraction. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, 2023b.

Quentin Delfosse, Wolfgang Stammer, Thomas Rothenbächer, Dwarak Vittal, and Kristian Kersting. Boosting Object Representation Learning via Motion and Object Continuity. In *Machine Learning and Knowledge Discovery in Databases: Research Track*. Springer Nature Switzerland, 2023c. doi: [10.1007/978-3-031-43421-1_36](https://doi.org/10.1007/978-3-031-43421-1_36).

Quentin Delfosse, Sebastian Sztwiertnia, Mark Rothermel, Wolfgang Stammer, and Kristian Kersting. Interpretable Concept Bottlenecks to Align Reinforcement Learning Agents, 2024. arXiv preprint.

Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Kinal Mehta, and João G. M. Araújo. CleanRL: High-quality Single-file Implementations of Deep Reinforcement Learning Algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022. ISSN 1533-7928.

Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. A Survey of Reinforcement Learning from Human Feedback, 2023. arXiv preprint.

W. Bradley Knox, Alessandro Allievi, Holger Banzhaf, Felix Schmitt, and Peter Stone. Reward (Mis)design for autonomous driving. *Artificial Intelligence*, 316:103829, 2023. doi: [10.1016/j.artint.2022.103829](https://doi.org/10.1016/j.artint.2022.103829).

Agneza Krajna, Mario Brcic, Tomislav Lipic, and Juraj Doncevic. Explainability in reinforcement learning: Perspective and position, 2022. arXiv preprint.

Zhixuan Lin, Yi-Fu Wu, Skand Vishwanath Peri, Weihao Sun, Gautam Singh, Fei Deng, Jindong Jiang, and Sungjin Ahn. SPACE: Unsupervised Object-Oriented Scene Representation via Spatial Attention and Decomposition. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.

Lirui Luo, Guoxi Zhang, Hongming Xu, Yaodong Yang, Cong Fang, and Qing Li. INSIGHT: End-to-End Neuro-Symbolic Visual Reinforcement Learning with Language Explanations, 2024. arXiv preprint.

Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-Level Reward Design via Coding Large Language Models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.

Marlos C. Machado, Marc G. Bellemare, Erik Talvitie, Joel Veness, Matthew Hausknecht, and Michael Bowling. Revisiting the Arcade Learning Environment: Evaluation Protocols and Open Problems for General Agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018. doi: [10.1613/jair.5699](https://doi.org/10.1613/jair.5699).

Stephanie Milani, Nicholay Topin, Manuela Veloso, and Fei Fang. Explainable Reinforcement Learning: A Survey and Comparative Review. *ACM Computing Surveys*, 2023. doi: [10.1145/3616864](https://doi.org/10.1145/3616864).

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing Atari with Deep Reinforcement Learning, 2013. arXiv preprint.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015. doi: [10.1038/nature14236](https://doi.org/10.1038/nature14236).

Suraj Nair, Eric Mitchell, Kevin Chen, Brian Ichter, Silvio Savarese, and Chelsea Finn. Learning Language-Conditioned Robot Behavior from Offline Data and Crowd-Sourced Annotation. In *Proceedings of the Conference on Robot Learning (CoRL)*. PMLR, 2021.

Pierre-Yves Oudeyer and Frederic Kaplan. How can we define intrinsic motivation? In *Proceedings of the Eight International Conference on Epigenetic Robotics: Modeling Cognitive Development in Robotic Systems*. Lund University Cognitive Studies, Lund: LUCS, Brighton, 2008.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, 2022.

David Raposo, Sam Ritter, Adam Santoro, Greg Wayne, Theophane Weber, Matt Botvinick, Hado van Hasselt, and Francis Song. Synthetic Returns for Long-Term Credit Assignment, 2021. arXiv preprint.

Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You Only Look Once: Unified, Real-Time Object Detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. doi: [10.1109/CVPR.2016.91](https://doi.org/10.1109/CVPR.2016.91).

Machel Reid, Yutaro Yamada, and Shixiang Shane Gu. Can Wikipedia Help Offline Reinforcement Learning?, 2022. arXiv preprint.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms, 2017. arXiv preprint.

K. W. Spence. The role of secondary reinforcement in delayed reward learning. *Psychological Review*, 54(1):1–8, 1947. doi: [10.1037/h0056533](https://doi.org/10.1037/h0056533).

Mark Towers, Jordan K. Terry, Ariel Kwiatkowski, John U. Balis, Gianluca de Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, Arjun KG, Markus Krimmel, Rodrigo Perez-Vicente, Andrea Pierré, Sander Schulhoff, Jun Jet Tai, Andrew Tan Jin Shen, and Omar G. Younis. Gymnasium. Zenodo, 2023. URL <https://zenodo.org/records/8127026>.

Hado van Hasselt, Arthur Guez, and David Silver. Deep Reinforcement Learning with Double Q-Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016. doi: [10.1609/aaai.v30i1.10295](https://doi.org/10.1609/aaai.v30i1.10295).

Yufei Wang, Zhanyi Sun, Jesse Zhang, Zhou Xian, Erdem Biyik, David Held, and Zackory Erickson. RL-VLM-F: Reinforcement Learning from Vision Language Foundation Model Feedback, 2024. arXiv preprint.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc V. Le, and Denny Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, 2022.

Yue Wu, Yewen Fan, Paul Pu Liang, Amos Azaria, Yuanzhi Li, and Tom M. Mitchell. Read and Reap the Rewards: Learning to Play Atari with the Help of Instruction Manuals. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, 2023.

Tianbao Xie, Siheng Zhao, Chen Henry Wu, Yitao Liu, Qian Luo, Victor Zhong, Yanchao Yang, and Tao Yu. Text2Reward: Reward Shaping with Language Models for Reinforcement Learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.

A Appendix

As mentioned in the main body, the appendix contains additional materials and supporting information for the following aspects: the hyperparameters used in this work (Appendix A.1), details on the prompts used for the LLM (Appendix A.3) as well as the generated reward functions (Appendix A.4), and numerical results (Appendix A.2).

A.1 Hyperparameters and Experimental Details

In this section, we list the hyperparameters used during the training and optimization of our models. For our experiments, we adopted the parameter set proposed by Huang et al. (2022) for our PPO agents, as detailed in Table 2.

Table 2: **Hyperparameter Configuration for Training Settings (PPO)**. This table provides a comprehensive overview of the essential hyperparameters utilized in our experimental section.

Hyperparameter	Value	Hyperparameter	Value
batch size	1024	Clipping Coef.	0.1
γ	0.99	KL target	None
minibatch size	256	GAE λ	0.95
seeds	42,73,91	input representation	4x84x84
total timesteps	10M	gym version	0.28.1
learning rate	0.00025	pytorch version	1.12.1
optimizer	Adam		
more information	https://docs.cleanrl.dev/rl-algorithms/c51/		

OCALM-based agents use the same PPO hyperparameter values as agents trained on ALE rewards. All agents use ConvNets (Mnih et al., 2015) with ReLU activation functions for policy and value networks. We utilized a decreasing learning rate of 2.5×10^{-4} over 10 million steps. We use the Atari environments in version v5 provided by Gymnasium Towers et al. (2023), following best-practices recommended by Machado et al. (2018). To accelerate training, we used 8 parallel game environments.

To mitigate noise and fluctuations, we use exponential moving average (EMA) smoothing in Figure 3 and Figure 4. We use an effective window size of 50, resulting in a smoothing factor $\alpha = 2/(1+50) \approx 0.039$ used in the following formula:

$$EMA_t = (1 - \alpha) \cdot EMA_{t-1} + \alpha \cdot y_t . \quad (1)$$

To manage irregular training intervals due to rewards are not always being reported in the same timestep, we ignore missing values when computing the average, relying on the EMA smoothing to provide a continuous curve. For the error bands we used the standard deviation of your data within a rolling window.

A.2 Numerical Results

In this section, we provide additional numerical results for the experiments conducted in this work.

Table 3: Numerical results for the experiments we conducted, including random and human baselines from van Hasselt et al. (2016) for comparison. Standard deviations are provided where available. Our agents use PPO and *ALE v5* and have been trained using 10 million frames. The results reported are the in-game rewards from the ALE/emulator, not from OCALM. Note, van Hasselt et al. (2016) predate the *v5* version of the ALE environments used by us, which is based on the best practices outlined by Machado et al. (2018). However, this should not change the values much since these changes have less influence on humans or the random agent.

Game	PPO ALE Reward (Baseline)	PPO OCALM (full) (Ours)	PPO OCALM (no relations)	Random van Hasselt et al.	Human van Hasselt et al.
Freeway	33.8 ± 0.2	32.35 ± 0.25	0.00	0.00	29.6
Pong	17.5 ± 0.5	16.4 ± 1.4	-15.8 ± 3.5	-20.7	14.6
Seaquest	1132.4 ± 271.4	672.2 ± 28.3	243 ± 86.3	68.4	20182
Skiing	-23921.3 ± 10528.6	-28577.7 ± 2842.3	-30000	-17098	-4336

A.3 LLM Prompting Details

In our experiments we used the LLM gpt-4-turbo¹ with seed 42 and top_k = 0. We further defined a system prompt that was used for both, direct and relational multi-turn prompting (Listing 2).

Listing 2: System prompt provided to the LLM.

```
1 You are a helpful assistant that creates reward functions for reinforcement learning
  researchers.
```

For direct prompting, we asked the LLM to create a reward function directly, given a game instruction and the game object classes. The game instructions are a few sentences that describe the objective of the game (based on the documentation of the Gymnasium environment collection (Towers et al., 2023)) and the game objects are Python classes provided by the OCAtari framework (Delfosse et al., 2023a), further described in Appendix A.3.2.

Listing 3: Prompt for reward function based on game instructions and game objects provided to the LLM.

```
1 We want to create a object centric reward function to train a reinforcement learning
  agent to play the game <GAME>. Here is a description of the game and its objects:
2
3 <PARENT GAME OBJECT CLASS>
4
5 <GAME OBJECT CLASSES>
6
7 The game instructions are the following:
8 <INSTRUCTIONS>
9
10 Please provide a Python file with a reward function that uses a list of objects of
    type GameObject as input that will help the agent to play the game, i.e.:
11 '''python
12 def reward_function(game_objects) -> float:
13     ...
14     return reward
15 '''
16
17 Do not use undefined variables or functions. Do not give any textual explanations,
    just generate the python code. If you give an explanation, please provide it in
    the form of a comment in the code.
```

For relational multi-turn prompting we first ask the LLM to provide functions that might be relevant for understanding the state and events of the game Listing 4. Based on these functions the model is then asked to create a reward function with Listing 5. As a last step, the model is asked to adjust the rewards so that they are on a scale from -1 to 1 Listing 6.

¹<https://platform.openai.com/docs/models/gpt-4-turbo-and-gpt-4>

Listing 4: Prompt for helpful functions based on game instructions and game objects.

```

1 We want to create a reward function for playing the Atari game <GAME>. As a first
  step we want to collect functions that are helpful for understanding events that
  are happening in the game that could be relevant for the reward, i.e., items
  colliding. In the following there will be existing game objects given, please
  generate functions that can be used to understand the game state. Please don't
  use undefined variables or functions.
2
3 Here is a description of the game and its objects:
4
5 <PARENT GAME OBJECT CLASS>
6
7 <GAME OBJECT CLASSES>
8
9 The game instructions are the following:
10 <INSTRUCTIONS>
```

Listing 5: Prompt for reward function based on identified functions from before.

```

1 Now please create a object centric reward function to train a reinforcement learning
  agent to play the game <GAME>. The reward function uses a list of objects of
  type GameObject as input, i.e.:
2 '''python
3 def reward_function(game_objects) -> float:
4     ...
5     return reward
6 '''
7 You can use the identified functions from before. Please don't use other undefined
  variables or functions.
```

,

Listing 6: Prompt for rescaling reward values.

```

1 "Thank you. Now please adjust the rewards so that the rewards are in the range [-1,
  1]."
```

A.3.1 The Object Properties used for OCALM

In this paper, we used different object properties as the inputs to the LLM-written reward functions. The object-centric environment context is given by the classes provided by the OCATari framework (Delfosse et al., 2023a), i.e., the parent game object class² and the game-specific objects³. The game objects related to the score were omitted (since we are assuming a reward-free environment). An overview of the object properties used by OCALM is provided in Table 4 and concrete implementation details can be found in Appendix A.3.2.

Name	Definition	Description
class	NAME	object class (e.g. "Agent", "Ball", "Ghost")
position	x, y	position on the screen
position history	$x_t, y_t, x_{t-1}, y_{t-1}$	position and past position on the screen
orientation	o	object's orientation if available
RGB	R, G, B	RGB values

Table 4: Descriptions of object properties used by OCALM.

²https://github.com/k4ntz/OC_Atari/blob/master/ocatari/ram/game_objects.py

³e.g., https://github.com/k4ntz/OC_Atari/blob/master/ocatari/ram/pong.py

A.3.2 Example of Game Objects

As described in the previous section, the game objects are Python classes provided by the (MIT licensed) OCAtari framework (Delfosse et al., 2023a). We provide the parent class for game objects and an example of game objects for Pong here for illustration purposes. Note that we have elided parts of the parent class in the listing for brevity, indicated by `#elided#`. Refer to https://github.com/k4ntz/OC_Atari/blob/v0.1.0/ocatari/ram/game_objects.py for the full parent class and https://github.com/k4ntz/OC_Atari/blob/v0.1.0/ocatari/ram/pong.py for the source of the Pong example.

Listing 7: The parent classes for game objects.

```
1 class GameObject:
2 """
3 The Parent Class of every detected object in the Atari games (RAM Extraction
4 mode)
5
6 #elided#
7 """
8
9 GET_COLOR = False
10 GET_WH = False
11
12 def __init__(self):
13     self.rgb = (0, 0, 0)
14     self._xy = (0, 0)
15     self.wh = (0, 0)
16     self._prev_xy = None
17     self._orientation = None
18     self.hud = False
19
20 def __repr__(self):
21     return f"{self.__class__.__name__} at ({self._xy[0]}, {self._xy[1]}),\n{self.wh}"
22
23 @property
24 def category(self):
25     return self.__class__.__name__
26
27 @property
28 def x(self):
29     return self._xy[0]
30
31 @property
32 def y(self):
33     return self._xy[1]
34
35 #elided
36
37 def _save_prev(self):
38     self._prev_xy = self._xy
39
40 # @x.setter
41 # def x(self, x):
42 #
43 #     self._xy = x, self.xy[1]
44
45 # @y.setter
46 # def y(self, y):
47 #     self._xy = self.xy[0], y
48
49 @property
50 def orientation(self):
51     return self._orientation
52
53 orientation.setter
54 def orientation(self, o):
```

```

54         self._orientation = o
55
56     @property
57     def center(self):
58         return self._xy[0] + self.wh[0]/2, self._xy[1] + self.wh[1]/2
59
60     def is_on_top(self, other):
61         """
62             Returns ``True`` if this and another gameobject overlap.
63
64             :return: True if objects overlap
65             :rtype: bool
66         """
67         return (other.x <= self.x <= other.x + other.w) and \
68                (other.y <= self.y <= other.y + other.h)
69
70     def manathan_distance(self, other):
71         """
72             Returns the manathan distance between the center of both objects.
73
74             :return: True if objects overlap
75             :rtype: bool
76         """
77         c0, c1 = self.center, other.center
78         return abs(c0[0] - c1[0]) + abs(c0[1]- c1[1])
79
80     def closest_object(self, others):
81         """
82             Returns the closest object from others, based on manathan distance between
83             the center of both objects.
84
85             :return: (Index, Object) from others
86             :rtype: int
87         """
88         if len(others) == 0:
89             return None
90         return min(enumerate(others), key=lambda item:
91             self.manathan_distance(item[1]))
92
93
94 class ValueObject(GameObject):
95     """
96         This class represents a game object that incorporates any notion of a value.
97         For example:
98             * the score of the player (or sometimes Enemy).
99             * the level of useable/deployable resources (oxygen bars, ammunition bars, power
100                gauges, etc.)
101                * the clock/timer
102
103        :ivar value: The value of the score.
104        :vartype value: int
105
106
107    def __init__(self):
108        super().__init__()
109        self._value = 0
110        self._prev_value = None
111
112    @property
113    def value(self):
114        return self._value
115
116    @value.setter
117    def value(self, value):
118        self._value = None if value is None else int(value)
119
120    @property
121    def prev_value(self):

```

```

119     if self._prev_value is not None:
120         return self._prev_value
121     else:
122         return self._value
123
124     def _save_prev(self):
125         super().___save_prev()
126         self._prev_value = self._value
127
128     @property
129     def value_diff(self):
130         return self.value - self.prev_value

```

Listing 8: Example of game objects for Pong.

```

1 class Player(GameObject):
2 """
3     The player figure i.e., the movable bar at the side.
4 """
5
6     def __init__(self):
7         super().__init__()
8         self.xy = 0, 0
9         self.wh = 4, 15
10        self.rgb = 92, 186, 92
11
12
13 class Enemy(GameObject):
14 """
15     The enemy bar on the opposite side.
16 """
17
18     def __init__(self):
19         super().__init__()
20         self.xy = 0, 0
21         self.wh = 4, 15
22         self.rgb = 213, 130, 74
23
24
25 class Ball(GameObject):
26 """
27     The game ball.
28 """
29
30     def __init__(self):
31         super().__init__()
32         self.xy = 0, 0
33         self.wh = 2, 4
34         self.rgb = 236, 236, 236

```

A.4 Game Descriptions and Rewards

In the following, we provide the used game descriptions, our OCALM (full) reward functions and OCALM (no relations) reward functions for the games Pong, Freeway, Seaquest and Skiing.

A.4.1 Freeway

Listing 9: Game description of Freeway.

```

1 You control a chicken that has to cross ten horizontal lanes of a freeway traffic.
2 There exist two chickens in the game, you control the left chicken, that starts at
   the bottom of the road and should go to the top.
3 Cars are traveling along the horizontal lanes and you should cross without getting
   run over by a car.
4 The screen height is 160.

```

Listing 10: OCALM (full) reward function of the game Freeway.

```

5 from ocatari.ram.freeway import *
6
7
8 def detect_collision(chicken, car):
9     """
10     Determines if a collision occurs between a chicken and a car.
11
12     :param chicken: An instance of the Chicken class.
13     :param car: An instance of the Car class.
14     :return: True if there is a collision, False otherwise.
15     """
16
17     chicken_x1, chicken_y1 = chicken.x, chicken.y
18     chicken_x2, chicken_y2 = chicken.x + chicken.w, chicken.y + chicken.h
19     car_x1, car_y1 = car.x, car.y
20     car_x2, car_y2 = car.x + car.w, car.y + car.h
21
22     # Check if any corner of the chicken is within the car's boundaries
23     if (car_x1 <= chicken_x1 <= car_x2 or car_x1 <= chicken_x2 <= car_x2) and (
24         car_y1 <= chicken_y1 <= car_y2 or car_y1 <= chicken_y2 <= car_y2
25     ):
26         return True
27     return False
28
29 def has_reached_top(chicken, screen_height):
30     """
31     Checks if the chicken has reached the top of the screen.
32
33     :param chicken: An instance of the Chicken class.
34     :param screen_height: The height of the screen (int).
35     :return: True if the chicken has reached the top, False otherwise.
36     """
37
38     return chicken.y <= 0
39
40 def progress_made(chicken, screen_height):
41     """
42     Measures the progress made by the chicken towards the top of the screen.
43
44     :param chicken: An instance of the Chicken class.
45     :param screen_height: The height of the screen (int).
46     :return: A float representing the percentage of the screen crossed.
47     """
48
49     return (screen_height - chicken.y) / screen_height
50
51 def check_if_reset(chicken, initial_position):
52     """
53     Checks if the chicken has been reset to the initial position (usually after a
54     collision).
55
56     :param chicken: An instance of the Chicken class.
57     :param initial_position: The starting position of the chicken (tuple of int).
58     :return: True if the chicken is at the initial position, False otherwise.
59     """
60
61     return chicken.xy == initial_position
62
63 def find_closest_car(chicken, cars):
64     """
65     Finds the closest car to the chicken based on the Manhattan distance.
66
67     :param chicken: An instance of the Chicken class.
68     :param cars: A list of Car instances.
69     :return: The closest car instance.
    """

```

```

70     if not cars:
71         return None
72     return min(cars, key=lambda car: chicken.manathan_distance(car))
73
74
75 def reward_function(game_objects) -> float:
76     # Initialize reward
77     reward = 0.0
78
79     # Constants
80     SCREEN_HEIGHT = 160
81     COLLISION_PENALTY = -1.0 # Scaled down to fit within [-1, 1]
82     PROGRESS_REWARD = 0.1 # Scaled down to incrementally increase reward
83     SUCCESS_REWARD = 1.0 # Maximum reward for reaching the top
84
85     # Filter out chickens and cars from game_objects
86     chickens = [obj for obj in game_objects if isinstance(obj, Chicken)]
87     cars = [obj for obj in game_objects if isinstance(obj, Car)]
88
89     # Assume control of the leftmost chicken (player's chicken)
90     if chickens:
91         player_chicken = min(chickens, key=lambda c: c.x)
92
93         # Check if the chicken has reached the top
94         if has_reached_top(player_chicken, SCREEN_HEIGHT):
95             reward += SUCCESS_REWARD
96
97         # Reward based on progress towards the top
98         reward += progress_made(player_chicken, SCREEN_HEIGHT) * PROGRESS_REWARD
99
100        # Check for collisions with any car
101        for car in cars:
102            if detect_collision(player_chicken, car):
103                reward += COLLISION_PENALTY
104                break # Only penalize once per time step
105
106    # Ensure reward stays within the range [-1, 1]
107    reward = max(min(reward, 1.0), -1.0)
108
109    return reward

```

Listing 11: Direct reward function of the game Freeway.

```

110 from ocatari.ram.freeway import *
111
112 def reward_function(game_objects) -> float:
113     # Initialize reward
114     reward = 0.0
115
116     # Define constants for screen dimensions
117     SCREEN_HEIGHT = 160
118
119     # Find the player's chicken and cars in the game objects
120     player_chicken = None
121     cars = []
122
123     for obj in game_objects:
124         if obj.category == "Chicken" and obj.x < SCREEN_HEIGHT // 2: # Assuming
125             player_chicken starts on the left
126             player_chicken = obj
127         elif obj.category == "Car":
128             cars.append(obj)
129
130     # If the player's chicken is not found, return a neutral reward
131     if player_chicken is None:
132         return reward
133
134     # Reward for moving upwards, normalized by screen height

```

```

134     reward += (player_chicken.dy / SCREEN_HEIGHT)
135
136     # Penalty for moving downwards, normalized and increased in impact
137     if player_chicken.dy < 0:
138         reward -= 2 * (abs(player_chicken.dy) / SCREEN_HEIGHT)
139
140     # Check for collisions with cars
141     for car in cars:
142         if player_chicken.is_on_top(car):
143             # Large penalty for collisions, ensuring it stays within the range
144             reward -= 0.5
145
146     # Bonus for reaching the top of the screen, ensuring it stays within the range
147     if player_chicken.y <= 0:
148         reward += 0.5
149
150     # Ensure the reward is clamped between -1 and 1
151     reward = max(min(reward, 1), -1)
152
153     return reward

```

A.4.2 Pong

Listing 12: Game description of Pong.

154 In this game the agent has to knock the ball past the enemy's paddle, situated on the left (good), and avoid letting the ball go past its paddle, on the right (bad). If the ball passes the paddle of the enemy, the agent gets a point. If the ball passes past the agents paddle, the enemy gets a point. After a scored point the ball continues moving behind the paddle for multiple frames before it respawns for a new round.

Listing 13: OCALM (full) reward function of the game Pong.

```

1 from ocatari.ram.pong import *
2
3
4 def check_collision(obj1, obj2):
5     """
6     Check if two GameObjects are colliding.
7
8     :param obj1: First game object.
9     :param obj2: Second game object.
10    :return: True if objects are colliding, False otherwise.
11    """
12    return (
13        obj1.x < obj2.x + obj2.w
14        and obj1.x + obj1.w > obj2.x
15        and obj1.y < obj2.y + obj2.h
16        and obj1.y + obj1.h > obj2.y
17    )
18
19
20 def ball_passed_paddle(ball, paddle, playing_field_width):
21     """
22     Check if the ball has passed the given paddle.
23
24     :param ball: The ball object.
25     :param paddle: The paddle object (player or enemy).
26     :param playing_field_width: The width of the playing field.
27     :return: True if the ball has passed the paddle, False otherwise.
28     """
29     if paddle.category == "Player":
30         # Check if the ball has passed the player's paddle on the right
31         return ball.x > playing_field_width
32     elif paddle.category == "Enemy":
33         # Check if the ball has passed the enemy's paddle on the left

```

```

34         return ball.x + ball.w < 0
35     return False
36
37
38 def update_game_state(objects):
39     """
40     Update the game state by saving the previous positions of the objects.
41
42     :param objects: List of all game objects.
43     """
44     for obj in objects:
45         obj._save_prev()
46
47
48 def detect_score_event(ball, player_paddle, enemy_paddle, playing_field_width):
49     """
50     Detect if a scoring event has occurred.
51
52     :param ball: The ball object.
53     :param player_paddle: The player's paddle object.
54     :param enemy_paddle: The enemy's paddle object.
55     :param playing_field_width: The width of the playing field.
56     :return: 'player' if player scores, 'enemy' if enemy scores, None otherwise.
57     """
58     if ball_passed_paddle(ball, enemy_paddle, playing_field_width):
59         return "player"
60     elif ball_passed_paddle(ball, player_paddle, playing_field_width):
61         return "enemy"
62     return None
63
64
65 def reward_function(game_objects) -> float:
66     # Initialize reward
67     reward = 0.0
68
69     # Identify the ball, player paddle, and enemy paddle from game_objects
70     ball = None
71     player_paddle = None
72     enemy_paddle = None
73     playing_field_width = 160 # Assuming a standard Atari Pong width
74
75     for obj in game_objects:
76         if isinstance(obj, Ball):
77             ball = obj
78         elif isinstance(obj, Player):
79             player_paddle = obj
80         elif isinstance(obj, Enemy):
81             enemy_paddle = obj
82
83     # Check if the ball and paddles are identified correctly
84     if ball is None or player_paddle is None or enemy_paddle is None:
85         return reward # No change in reward if essential objects are missing
86
87     # Check for scoring events
88     score_event = detect_score_event(
89         ball, player_paddle, enemy_paddle, playing_field_width
90     )
91     if score_event == "player":
92         reward += 1.0 # Reward for scoring against the enemy
93     elif score_event == "enemy":
94         reward -= 1.0 # Penalty for the enemy scoring
95
96     # Check for collisions with paddles to encourage hitting the ball
97     if check_collision(ball, player_paddle) or check_collision(ball, enemy_paddle):
98         reward += 0.1 # Small reward for hitting the ball with either paddle
99
100    # Normalize the reward to be within [-1, 1]
101    reward = max(min(reward, 1), -1)

```

```

102
103     return reward

```

Listing 14: Direct reward function of the game Pong.

```

155 from ocatari.ram.pong import *
156
157 def reward_function(game_objects) -> float:
158     # Initialize reward
159     reward = 0.0
160
161     # Define identifiers for game objects
162     player_id = None
163     enemy_id = None
164     ball_id = None
165
166     # Find indices for player, enemy, and ball
167     for i, obj in enumerate(game_objects):
168         if obj.category == "Player":
169             player_id = i
170         elif obj.category == "Enemy":
171             enemy_id = i
172         elif obj.category == "Ball":
173             ball_id = i
174
175     # Get the game objects
176     player = game_objects[player_id]
177     enemy = game_objects[enemy_id]
178     ball = game_objects[ball_id]
179
180     # Check if the ball has passed the enemy paddle
181     if ball.x < enemy.x:
182         reward += 1 # Reward for scoring a point
183
184     # Check if the ball has passed the player paddle
185     if ball.x > player.x + player.w:
186         reward -= 1 # Penalty for letting the enemy score
187
188     # Normalize reward to be within the range [-1, 1]
189     reward = max(min(reward, 1), -1)
190
191     return reward

```

A.4.3 Seaquest

Listing 15: Game description of Seaquest.

```

192 You a sub (Player) able to move in all directions and fire torpedoes.
193 The goal is to retrieve as many divers as you can, while dodging and blasting enemy
194 subs and killer sharks.
195 The game begins with one sub and three waiting on the horizon. Each time you
196 increase your score by 10,000 points, an extra sub will be delivered to your
197 base.
198 Your sub will explode if it collides with anything except your divers. The sub has a
199 limited amount of oxygen that decreases at a constant rate during the game. When
200 the oxygen tank is almost empty, you need to surface and if you don't do it in
201 time, your sub will blow up and you'll lose one diver.
202 Each time you're forced to surface, with less than six divers, you lose one diver as
203 well.

```

Listing 16: OCALM (full) reward function of the game Seaquest.

```

197 from ocatari.ram.seaquest import *
198
199 def check_collision(obj1, obj2):
200     """
201         Check if two GameObjects collide based on their bounding boxes.

```

```

202     """
203     return (obj1.x < obj2.x + obj2.w and
204         obj1.x + obj1.w > obj2.x and
205         obj1.y < obj2.y + obj2.h and
206         obj1.y + obj1.h > obj2.y)
207
208 def update_game_state(objects):
209     """
210     Update positions of all game objects and check for collisions.
211     """
212     collisions = []
213     for obj in objects:
214         # Update position based on velocity
215         obj.xy = (obj.x + obj.dx, obj.y + obj.dy)
216
217         # Check for collisions with other objects
218         for other in objects:
219             if obj != other and check_collision(obj, other):
220                 collisions.append((obj, other))
221     return collisions
222
223 def manage_oxygen_and_lives(player, oxygen_bar, lives):
224     """
225     Decrease oxygen levels and manage lives based on oxygen and collisions.
226     """
227     # Decrease oxygen
228     oxygen_bar.value -= 1
229     if oxygen_bar.value <= 0:
230         player.lives -= 1
231         oxygen_bar.value = 100 # Reset oxygen after surfacing or losing a life
232
233     # Check if lives are depleted
234     if lives.value <= 0:
235         print("Game Over")
236
237 def update_score_and_divers(player, divers_collected, score):
238     """
239     Update score based on collected divers and manage divers.
240     """
241     for diver in divers_collected:
242         if check_collision(player, diver):
243             score.value += 1000 # Increment score for each diver collected
244             divers_collected.remove(diver) # Remove diver from the game
245
246 def fire_torpedo(player, torpedoes):
247     """
248     Create a new torpedo at the player's location and add it to the torpedoes list.
249     """
250     new_torpedo = PlayerMissile()
251     new_torpedo.xy = player.xy
252     torpedoes.append(new_torpedo)
253
254
255 def reward_function(game_objects) -> float:
256     reward = 0.0
257
258     # Define categories for easy identification
259     player = None
260     divers = []
261     enemies = []
262     player_missiles = []
263     enemy_missiles = []
264     oxygen_bar = None
265
266     # Classify objects
267     for obj in game_objects:
268         if isinstance(obj, Player):
269             player = obj

```

```

270     elif isinstance(obj, Diver):
271         divers.append(obj)
272     elif isinstance(obj, Shark) or isinstance(obj, Submarine):
273         enemies.append(obj)
274     elif isinstance(obj, PlayerMissile):
275         player_missiles.append(obj)
276     elif isinstance(obj, EnemyMissile):
277         enemy_missiles.append(obj)
278     elif isinstance(obj, OxygenBar):
279         oxygen_bar = obj
280
281     # Check for collisions and manage interactions
282     if player:
283         for diver in divers:
284             if check_collision(player, diver):
285                 reward += 0.1 # Scaled down reward for collecting a diver
286                 divers.remove(diver) # Assume diver is collected and removed from
287                 the game
288
289         for enemy in enemies:
290             if check_collision(player, enemy):
291                 reward -= 0.1 # Scaled down penalty for colliding with an enemy
292
293         for missile in enemy_missiles:
294             if check_collision(player, missile):
295                 reward -= 0.05 # Scaled down penalty for getting hit by an enemy
296                 missile
297
298         # Reward for hitting enemies with missiles
299         for missile in player_missiles:
300             for enemy in enemies:
301                 if check_collision(missile, enemy):
302                     reward += 0.05 # Scaled down reward for destroying an enemy
303                     enemies.remove(enemy) # Assume enemy is destroyed and removed
304                     from the game
305                     player_missiles.remove(missile) # Remove missile after hitting
306
307         # Manage oxygen levels
308         if oxygen_bar and oxygen_bar.value <= 20:
309             reward -= 0.05 # Scaled down penalty for low oxygen levels
310
311         # Encourage surfacing if oxygen is too low
312         if oxygen_bar and oxygen_bar.value <= 10:
313             reward -= 0.1 # Scaled down higher penalty for critically low oxygen
314
315     return reward

```

Listing 17: Direct reward function of the game Seaquest.

```

313 from ocatari.ram.seaquest import *
314
315 def reward_function(game_objects) -> float:
316     # Initialize reward
317     reward = 0.0
318
319     # Define constants for reward values, adjusted to be within the range [-1, 1]
320     REWARD_FOR_COLLECTING_DIVER = 0.1
321     PENALTY_FOR_COLLISION = -0.1
322     PENALTY_FOR_LOW_OXYGEN = -0.05
323     REWARD_FOR_BLASTING_ENEMY = 0.05
324     PENALTY_FOR_LOSING_DIVER_WHEN_SURFACING = -0.025
325
326     # Helper function to find an object by its class name
327     def find_objects_by_type(type_name):
328         return [obj for obj in game_objects if obj.category == type_name]
329
330     # Get specific game objects

```

```

331     player = find_objects_by_type('Player')[0] if find_objects_by_type('Player')
332     else None
333     divers = find_objects_by_type('Diver')
334     sharks = find_objects_by_type('Shark')
335     enemy_subs = find_objects_by_type('Submarine')
336     enemy_missiles = find_objects_by_type('EnemyMissile')
337     player_missiles = find_objects_by_type('PlayerMissile')
338     oxygen_bar = find_objects_by_type('OxygenBar')[0] if
339     find_objects_by_type('OxygenBar') else None
340
341     # Reward for collecting divers
342     for diver in divers:
343         if player and player.is_on_top(diver):
344             reward += REWARD_FOR_COLLECTING_DIVER
345
346     # Penalty for collisions with sharks or enemy submarines
347     for shark in sharks:
348         if player and player.is_on_top(shark):
349             reward += PENALTY_FOR_COLLISION
350
351     for enemy_sub in enemy_subs:
352         if player and player.is_on_top(enemy_sub):
353             reward += PENALTY_FOR_COLLISION
354
355     # Check for low oxygen
356     if oxygen_bar and oxygen_bar.value < 20:
357         reward += PENALTY_FOR_LOW_OXYGEN
358
359     # Reward for blasting enemy submarines with missiles
360     for missile in player_missiles:
361         for enemy_sub in enemy_subs:
362             if missile.is_on_top(enemy_sub):
363                 reward += REWARD_FOR_BLASTING_ENEMY
364
365     # Penalty for enemy missiles hitting the player
366     for missile in enemy_missiles:
367         if player and missile.is_on_top(player):
368             reward += PENALTY_FOR_COLLISION
369
370     # Penalty for surfacing with less than six divers
371     collected_divers = find_objects_by_type('CollectedDiver')
372     if len(collected_divers) < 6 and player and player.y == 0: # Assuming y=0 is
373     the surface
374     reward += PENALTY_FOR_LOSING_DIVER_WHEN_SURFACING
375
376     return reward

```

A.4.4 Skiing

Listing 18: Game description of Skiing.

```

374 You control a skier (Player), going down a slope who can move sideways.
375 The Player is at the top of the screen, staying at the same y position but the other
376 objects of the environments are moving up towards him.
377 The goal is to ski in between the horizontal pairs of flags.
378 There can be up to two pairs of poles on the screen.
379 Do not hit a tree or a flag or you'll fall and lose time.

```

Listing 19: OCALM (full) reward function of the game Skiing.

```

379 from ocatari.ram.skiing import *
380
381 def check_collision(player, objects):
382     """
383     Check if the player has collided with any of the given objects (flags or trees).
384
385     :param player: The player object.

```

```

386     :param objects: A list of game objects (flags or trees).
387     :return: True if a collision is detected, False otherwise.
388     """
389     for obj in objects:
390         if (player.x < obj.x + obj.w and
391             player.x + player.w > obj.x and
392             player.y < obj.y + obj.h and
393             player.y + player.h > obj.y):
394             return True
395     return False
396
397 def check_gate_passage(player, flag1, flag2):
398     """
399     Check if the player has passed between two flags.
400
401     :param player: The player object.
402     :param flag1: The first flag object.
403     :param flag2: The second flag object.
404     :return: True if the player has passed between the flags, False otherwise.
405     """
406     if flag1.y == flag2.y: # Ensure flags are aligned horizontally
407         gate_left = min(flag1.x, flag2.x)
408         gate_right = max(flag1.x + flag1.w, flag2.x + flag2.w)
409         player_center_x = player.x + player.w / 2
410         return gate_left <= player_center_x <= gate_right
411     return False
412
413 def distance_to_nearest_obstacle(player, obstacles):
414     """
415     Calculate the Manhattan distance from the player to the nearest obstacle.
416
417     :param player: The player object.
418     :param obstacles: A list of obstacle objects (trees or moguls).
419     :return: The Manhattan distance to the closest obstacle.
420     """
421     closest_distance = float('inf')
422     for obstacle in obstacles:
423         distance = player.manathan_distance(obstacle)
424         if distance < closest_distance:
425             closest_distance = distance
426     return closest_distance
427
428 def update_game_state(objects, dy):
429     """
430     Update the positions of all game objects to simulate the player moving down the
431     slope.
432
433     :param objects: A list of all movable game objects.
434     :param dy: The change in y-coordinate to apply to each object.
435     """
436     for obj in objects:
437         if not obj.hud:
438             obj.y += dy
439
440 def reward_function(objects):
441     """
442     Calculate the reward for the current game state based on the player's
443     interactions with game objects,
444     with rewards scaled to be within the range [-1, 1].
445
446     :param objects: A list of all game objects currently in the game.
447     :return: A numerical reward value.
448     """
449     # Initialize reward
450     reward = 0

```

```

451     # Define penalty and reward values, scaled to ensure the total reward remains
452     # within [-1, 1]
453     collision_penalty = -1
454     gate_pass_reward = 0.5
455     close_to_obstacle_penalty = -0.01
456
457     # Filter objects by type
458     player = next(obj for obj in objects if isinstance(obj, Player))
459     flags = [obj for obj in objects if isinstance(obj, Flag)]
460     trees = [obj for obj in objects if isinstance(obj, Tree)]
461     moguls = [obj for obj in objects if isinstance(obj, Mogul)]
462
463     # Check for collisions with trees or flags
464     if check_collision(player, trees + flags):
465         reward += collision_penalty
466
467     # Check for successful gate passage
468     # Assuming flags come in pairs, sorted by y and then x
469     flags_sorted = sorted(flags, key=lambda f: (f.y, f.x))
470     for i in range(0, len(flags_sorted), 2):
471         if i+1 < len(flags_sorted) and check_gate_passage(player, flags_sorted[i],
472             flags_sorted[i+1]):
473             reward += gate_pass_reward
474
475     # Penalize being too close to the nearest obstacle
476     nearest_obstacle_distance = distance_to_nearest_obstacle(player, trees + moguls)
477     # Assuming a threshold below which the player is considered too close to an
478     # obstacle
479     if nearest_obstacle_distance < 20:
480         reward += close_to_obstacle_penalty * (20 - nearest_obstacle_distance)
481
482     # Ensure the reward is within the range [-1, 1]
483     reward = max(min(reward, 1), -1)
484
485     return reward

```

Listing 20: Direct reward function of the game Skiing.

```

483 from ocatari.ram.skiing import *
484
485 def reward_function(game_objects) -> float:
486     # Initialize reward
487     reward = 0.0
488
489     # Define constants for reward/penalty values
490     FLAG_PASS_REWARD = 0.1
491     TREE_COLLISION_PENALTY = -0.3
492     FLAG_COLLISION_PENALTY = -0.2
493     MOGUL_COLLISION_PENALTY = -0.05
494
495     # Helper function to check if the player collides with any object
496     def check_collision(player, obj):
497         return (obj.x <= player.x <= obj.x + obj.w or obj.x <= player.x + player.w
498             <= obj.x + obj.w) and \
499                 (obj.y <= player.y <= obj.y + obj.h or obj.y <= player.y + player.h
500                 <= obj.y + obj.h)
501
502     # Extract player and other objects
503     player = None
504     flags = []
505     trees = []
506     moguls = []
507
508     for obj in game_objects:
509         if isinstance(obj, Player):
510             player = obj
511         elif isinstance(obj, Flag):
512             flags.append(obj)

```

```
511     elif isinstance(obj, Tree):
512         trees.append(obj)
513     elif isinstance(obj, Mogul):
514         moguls.append(obj)
515
516     # Check for collisions and successful flag passes
517     if player:
518         # Check for collisions with trees
519         for tree in trees:
520             if check_collision(player, tree):
521                 reward += TREE_COLLISION_PENALTY
522
523         # Check for collisions with flags
524         for flag in flags:
525             if check_collision(player, flag):
526                 reward += FLAG_COLLISION_PENALTY
527
528         # Check for collisions with moguls
529         for mogul in moguls:
530             if check_collision(player, mogul):
531                 reward += MOGUL_COLLISION_PENALTY
532
533     # Check if player passes between flags (assuming flags come in pairs)
534     if len(flags) >= 2:
535         # Sort flags by x to pair them
536         sorted_flags = sorted(flags, key=lambda f: f.x)
537         for i in range(0, len(sorted_flags) - 1, 2):
538             flag1 = sorted_flags[i]
539             flag2 = sorted_flags[i+1]
540             # Check if player is between the flags
541             if flag1.x < player.x < flag2.x:
542                 reward += FLAG_PASS_REWARD
543
544     return reward
```