# On the Sample Efficiency of Abstractions and Potential-Based Reward Shaping in Reinforcement Learning

**Giuseppe Canonaco[1,*], Leo Ardon[2], Alberto Pozanco[1] and Daniel Borrajo[1]**

[1]AI Research Dept. JPMorganChase Madrid, ES
[2]AI Research Dept. JPMorganChase London, UK

**Abstract.** The use of Potential-Based Reward Shaping (PBRS) has shown great promise in the ongoing research effort to tackle sample inefficiency in Reinforcement Learning (RL). However, choosing the right potential function remains an open challenge. Additionally, RL techniques are usually constrained to use a finite horizon for computational limitations, which introduces a bias when using PBRS. In this paper, we first build some theoretically-grounded intuition on why selecting the potential function as the optimal value function of the task at hand produces performance advantages. We then analyse the bias induced by finite horizons in the context of PBRS producing novel insights. Finally, leveraging abstractions as a way to approximate the optimal value function of the given task, we assess the sample efficiency and performance impact of PBRS on four environments including a goal-oriented navigation task and three Arcade Learning Environments (ALE) games. Remarkably, experimental results show that we can reach the same level of performance as CNN-based solutions with a simple fully-connected network.

## 1 Introduction

Reinforcement Learning (RL) is able to achieve impressive results at the cost of requiring a huge amount of samples [34, 39, 29]. This sample inefficiency prevents RL from being more widely applied in the real world because of the cost associated with experience collection. Thus, the increasing demand of algorithms that can counteract or mitigate this issue.

In the RL literature, many different methods tackle the above-mentioned problem [37, 30, 38]. However, a promising way to compensate for the widely known sample inefficiency of RL is through Potential-Based Reward Shaping (PBRS) [28]. One property of this technique is that if we select the potential function to be the optimal value function of the problem at hand, the resulting reshaped problem will have an optimal value function that is zero everywhere, thus simplifying the learning process [28]. Unfortunately, the literature lacks theoretically-grounded motivations explaining why PBRS should entail sample efficiency for such a choice in RL problems.

Additionally, the optimal value function for most tasks is not available in practice, so computing estimates of that function becomes the only viable option. In case some knowledge about the task we are facing is available (e.g., via subject-matter experts), one way to approximate its optimal value function is through abstractions [8]. An abstraction here represents a mapping of the original task to a simpler one, such that solving the simpler one can accelerate the learning phase in the original task itself.

Due to computational limitations, we also have to practically impose a maximum number of steps allowed to solve an RL problem (finite horizon). This practical consideration has the disadvantage of introducing a bias in the context of PBRS as highlighted by recent literature [15, 17, 8]. Importantly, the introduced bias may disrupt the policy ordering in the reshaped Markov Decision Process (MDP) [31] practically turning policies that were optimal into suboptimal ones. Therefore, the need for a theoretical analysis investigating what can be preserved when applying PBRS under finite horizons.

In this paper we try to cover all the above mentioned gaps in the literature. In order to do so, we provide:

- A theoretically-grounded intuition on why PBRS entails performance benefits when we select the potential function as the optimal value function of the task at hand.
- A formal proof for the absence of bias when using PBRS in finite-horizon goal-oriented MDPs (under some conditions on the potential if the MDP is stochastic) addressing an issue opened by Cipollone et al. [8].
- A novel investigation of the above-mentioned bias for general MDPs, producing conditions under which we can preserve the total policy ordering (or only the optimal policy) w.r.t. the original infinite-horizon MDP.
- Leveraging abstractions as a tool to approximate the optimal value function of the task at hand, we will first provide an experimental back up of our theoretical findings in goal-oriented MDPs, outperforming Cipollone et al. [8]. Then, we will show how we can make a significant difference in terms of sample efficiency in the Arcade Learning Environment (ALE) benchmark [4] by accepting the bias induced by PBRS in general MDPs. Our approach achieves performances comparable or superior to the ones of Convolutional Neural Networks (CNNs) [27], but only using fully-connected architectures. To the best of our knowledge, this is unprecedented.

## 2 Background

In the RL setting, an agent interacts with the environment to learn a policy. After perceiving the current state of the environment, the agent executes an action that makes the environment transition into a

---

* Corresponding Author. Email: giuseppe.canonaco@jpmorgan.com

next state. Finally, the agent collects the reward associated with the experienced transition and the cycle starts again.

## 2.1 The Reinforcement Learning Problem

Problems in the RL paradigm are usually modeled through MDPs. An MDP is described as a tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma, \rho\}$, where $\mathcal{S}$ is the state space, $\mathcal{A}$ is the action space, $\mathcal{P}(s'|s, a)$ is the Markovian state transition function, $r(s, a, s') \in [0, \bar{R}]$ is the reward function, $\gamma \in [0, 1)$ is the discount factor, and $\rho$ is the initial state distribution. Goal-oriented MDPs characterize a sub-family of MDPs $\{\mathcal{S}, \mathcal{S}_G, \mathcal{A}, \mathcal{P}, r, \gamma, \rho\}$, where $\mathcal{S}_G \subset \mathcal{S}$ is a set of goal states, and the reward function $r(s, a, s') = 1$ if and only if $s \notin \mathcal{S}_G \wedge s' \in \mathcal{S}_G$; it is 0 otherwise.

Techniques to solve an RL problem can be divided into Value-based and Policy-based algorithms. The former strive to learn the optimal (action-)value function $V^*(s)$ $(Q^*(s, a))$, whereas the latter aim at directly learning the optimal policy. $Q^*(s, a)$ is defined as the expected obtained return starting in state $s$, executing action $a$, and following the optimal policy thereafter. Instead, $V^*(s)$ can be obtained starting from $Q^*(s, a)$ and selecting $a$ according to the optimal policy, or, equivalently, selecting $a$ greedily w.r.t. $Q^*(s, a)$. This means that $V^*(s) = \max_a Q^*(s, a)$ and that the optimal policy can be obtained by acting greedily on the environment under the guidance of $Q^*(s, a)$.

A current standard way to learn the above-mentioned functions is through the use of Deep Q-Networks (DQNs) [27]. In this case, given a dataset of collected experience $D = \langle s_t, a_t, r_t, s_{t+1} \rangle_{t=1}^N$, a DQN minimizes the following loss function during learning, leveraging samples of experience drawn uniformly at random from $D$:

$$\mathbb{E}_{(s,a,r,s') \sim D} \left[ \left( r + \gamma \max_{a'} Q_{\theta^-}(s', a') - Q_\theta(s, a) \right)^2 \right],$$

where $Q_\theta$ and $Q_{\theta^-}$ are action-value functions parameterized by the weights $\theta$ and $\theta^-$ that are represented by CNNs [25]. They respectively represent the actual action-value function and the target action-value function. The former is always updated while the latter only every $M$ steps.

## 2.2 Potential-Based Reward Shaping

The idea behind PBRS [28] is to steer the learning process towards paths that reach optimal behaviour faster by providing additional reward feedback through potential functions. This is accomplished by running the learning algorithm on $\mathcal{M}' = \{\mathcal{S}, \mathcal{A}, \mathcal{P}, r', \gamma, \rho\}$, a different MDP, called the reshaped MDP, where $r'(s, a, s') = r(s, a, s') + F(s, a, s')$ is the only change over the original MDP. Let us define a potential function as $\phi : \mathcal{S} \to \mathbb{R}$. If $F(s, a, s') = \gamma \phi(s') - \phi(s)$, $\phi$ is bounded, and the horizon length, $H$, is infinite, we are guaranteed that the optimal policy in $\mathcal{M}'$ is optimal also in $\mathcal{M}$ (this holds also for $\epsilon$-optimal policies) [28]. $\phi$ can be interpreted as a guide, telling the agent how valuable a certain state is. Furthermore, if we use this reshaping approach, the following relationships between (action-)value functions hold:

$$\mathcal{V}_{\mathcal{M}'}^\pi(s) = \mathcal{V}_{\mathcal{M}}^\pi(s) - \phi(s) \tag{1}$$

$$\mathcal{Q}_{\mathcal{M}'}^\pi(s, a) = \mathcal{Q}_{\mathcal{M}}^\pi(s, a) - \phi(s) \tag{2}$$

as shown by **Corollary 2** in [28]. Therefore, if $\phi = \mathcal{V}_{\mathcal{M}}^{\pi^*}$, then the optimal value function in the reshaped MDP is equivalent to zero.

In practice, RL algorithms cannot solve problems using an infinite horizon (as the previous guarantees require). Hence, we usually have to set a finite limit to $H$, $H < \infty$. This generates a bias, since the reshaped return depends on the final state reached by the agent [17]:

$$\mathcal{R}'(\tau) = \mathcal{R}(\tau) + \gamma^H \phi(s_H) - \phi(s_0), \tag{3}$$

where $\tau = s_0, a_0, s_1, \ldots, a_{H-1}, s_H$ is a trajectory and $\mathcal{R}(\tau) = \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t, s_{t+1})$ the discounted return. Since the last visited state depends on the policy, optimizing $\mathcal{R}'(\tau)$ is not the same as optimizing $\mathcal{R}(\tau)$.

## 3 Abstractions for PBRS

As expected, the potential function choice has an impact on PBRS efficiency. If we use the optimal value function as the potential function, it represents an excellent choice, because it makes the optimal value function of the reshaped MDP equivalent to zero (an easy function to be learned). But, it is unrealistic to assume we already have the optimal value function. Hence, it is common to use handcrafted heuristics in an attempt to approximate the optimal value function of the original task [28].

Inspired by the use of abstractions for heuristics computation in heuristic search and automated planning [10, 19], an alternative approach consists of defining an MDP, $\mathcal{M}_\alpha$, that is an abstraction of the original task [8] formalized by an aggregation function $\alpha$ whose definition can be found later. $\mathcal{M}_\alpha$'s optimal value function can then be used as the potential function for PBRS. This idea works under the intuition that the abstraction should be a simplified version of the original task, i.e., with a much smaller state/action space. Thus, generating the optimal value function of $\mathcal{M}_\alpha$ will see a strong reduction in computation. Besides, it provides an approximation of the original task's optimal value function. Furthermore, in many tasks, it is much easier for a human to provide knowledge via the definition of a simpler task rather than through the definition of a specific potential function. Take as an example specifying the potential function for a general ALE task that, in its RAM state representation, has $2^{1024}$ states. Defining a potential function over that many states in a meaningful way is a non-trivial endeavor. However, simplifying the game itself through an abstraction that neglects enemies or other complex dynamics of the game is a much easier task that can greatly reduce the number of states/actions to consider.

To formally define reshaping via abstractions, let $\mathcal{S}_\alpha$ represent the state space of $\mathcal{M}_\alpha$, with $|\mathcal{S}_\alpha| \ll |\mathcal{S}|$ and $\alpha : \mathcal{S} \to \mathcal{S}_\alpha$ be an aggregation function, mapping each state of the original task into a state in $\mathcal{S}_\alpha$. Also, let $\mathcal{V}_{\mathcal{M}_\alpha}^*(s_\alpha), s_\alpha \in \mathcal{S}_\alpha$ be the optimal value function of $\mathcal{M}_\alpha$. Then, we have the following reshaped reward in $\mathcal{M}'$, using $\mathcal{V}_{\mathcal{M}_\alpha}^*(s_\alpha)$ as the potential function ($\phi(s) = \mathcal{V}_{\mathcal{M}_\alpha}^*(\alpha(s))$):

$$r'(s, a, s') = r(s, a, s') + \gamma \mathcal{V}_{\mathcal{M}_\alpha}^*(\alpha(s')) - \mathcal{V}_{\mathcal{M}_\alpha}^*(\alpha(s)). \tag{4}$$

## 4 The Benefit of PBRS

In the previous section, we have given intuition on why the optimal value function constitutes a good choice for the potential and seen how to get an approximation of it. Let us now go deeper into why such a choice should provide a benefit for our algorithms. To do this, we begin by analyzing the performance impact of PBRS on Value Iteration (VI), which, to the best of our knowledge, has not been considered by related literature.

**Proposition 1** (VI Performance). Given an MDP $\mathcal{M}$, an initialization point $\mathcal{V}_0$, a potential function $\phi$ that is an approximation of $\mathcal{V}_{\mathcal{M}}^*$ such that $||\mathcal{V}_{\mathcal{M}}^* - \phi - \mathcal{V}_0||_\infty \leq ||\mathcal{V}_{\mathcal{M}}^* - \mathcal{V}_0||_\infty$, where $||\cdot||_\infty$ is the infinity norm, applying VI starting from $\mathcal{V}_0$ in both the original and reshaped MDP will result in the minimum number of iterations to guarantee $\epsilon$-optimality in the reshaped MDP to be lower than in the original MDP.

*Proposition 1* (proof in Appendix A.1) tells us that improvements are achievable in the context of Value Iteration whenever we are able to provide a good estimate of the optimal value function through the potential function using PBRS. However, this grants us just computational savings, and we would like to convert it into a greater sample efficiency gain in the context of a general RL setting.

For what concerns value-based algorithms like Q-Learning [43] and SARSA [32], Wiewiora [44] showed that, under the assumption of learning through the same experience, PBRS is equivalent to initializing the action-value function with the potential function. This can be interpreted as a simplistic form of Transfer Learning [37, 24]. Hence, whenever we initialize the action-value function with a potential function that represents a good starting estimate of the optimal solution, we will observe a greater sample efficiency due to a more favorable exploration-exploitation trade-off. On the other hand, in the context of the Policy Gradient Theorem (PGT) [36]:

$$\nabla_\theta J_\theta = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) \mathcal{Q}_{\mathcal{M}}^{\pi_\theta}(s,a), \quad (5)$$

where $d^{\pi_\theta}(s) = \sum_{t=0}^\infty \gamma^t Pr\{s_t = s|s_0, \pi_\theta\}$, baselines are commonly used to reduce the variance of the gradient estimator. The baseline version of PGT defines $\nabla_\theta J_\theta$ as:

$$\sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) (\mathcal{Q}_{\mathcal{M}}^{\pi_\theta}(s,a) - b(s)),$$

where $b(s)$ is the baseline. A baseline is a free parameter not affecting the estimate of the gradient, indeed $\sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) b(s) = 0$, that can reduce its variance [11].[1]

**Proposition 2** (Baseline Equivalence). In the context of the PGT, using PBRS (Grześ' approach in the finite-horizon case) is equivalent to using a baseline $b(s) = \phi(s)$.

*Proof.* Applying the PGT to the reshaped MDP $\mathcal{M}'$ via Eq. (5):

$$\nabla_\theta J_\theta = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) \mathcal{Q}_{\mathcal{M}'}^{\pi_\theta}(s,a).$$

Now, using Eq. (2), we get $\nabla_\theta J_\theta$ equivalent to:

$$\sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) (\mathcal{Q}_{\mathcal{M}}^{\pi_\theta}(s,a) - \phi(s)).$$

If we combine the above expression with $b(s) = \phi(s)$, it yields the baseline version of the PGT. In the case of finite horizons, we are forced to use a non-stationary potential [14] as proposed by Grześ [17] (which is unbiased), so that $\phi(s_H) = 0$ and the Monte-Carlo estimate of $\mathcal{Q}_{\mathcal{M}'}^{\pi_\theta}(s,a)$ can be decomposed as above.[2]  □

---

[1] $\sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) b(s) = 0$ can be obtained via $\nabla_\theta \pi_\theta(a|s) = \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s)$ and $\sum_a \nabla_\theta \pi_\theta(a|s) b(s) = b(s) \nabla_\theta \sum_a \pi_\theta(a|s) = 0$.

[2] $\phi(s_H) = 0$ only for the states encountered at the end of a trajectory. If the same state is encountered before, its original potential is used.

It is widely known that using value function estimates as baselines allows Policy Gradient algorithms to suffer less variance in their gradients [33]. This, in turn, implies a greater sample efficiency because of the more precise direction estimate toward which the algorithm will move during its optimization process. Since *Proposition 2* shows the equivalence between baselines and PBRS, we can deduce that a good value function estimate provided by the potential will help us improve sample efficiency, achieving better performance. The main advantage of using PBRS is that it is not mandatory to estimate the value function from data collected by the agent. Instead, the estimate may come from a simplified representation of the problem the agent is tackling. This enables a different way to provide knowledge to the agent while learning and does not require the representation to be perfect, relieving the user of the burden of providing accurate models.

## 5 Bias Analysis

Now, we focus on the bias induced by PBRS in the episodic setting, i.e., whenever we set a finite horizon due to computational limitations. We will start by considering goal-oriented MDPs, and then we will focus on general MDPs to understand the conditions that preserve policy ordering between the original problem and the reshaped one.

### 5.1 Goal-Oriented MDPs

An important family of goal-oriented MDPs is represented by those having a deterministic transition function, hence, we start our analysis from them.

**Theorem 3** (Policy Ordering in Deterministic Goal-Oriented Episodic MDPs). In the episodic context, let the horizon be $H$, $\mathcal{M}$ be a deterministic goal-oriented MDP, and $\mathcal{M}'$ be its reshaped version obtained using a bounded potential $\phi : \mathcal{S} \to [0, \Phi]$. Moreover, let $\phi(s_G) = \Phi \ \forall \ s_G \in \mathcal{S}_G$. Then, the policy ordering in $\mathcal{M}$ is preserved in the reshaped MDP $\mathcal{M}'$ for all the deterministic policies reaching the goal in at most $H$ steps.

*Proof.* Let $\pi^*$ be the optimal policy in the original MDP (the policy that requires the least number of steps to reach the goal), and $p_\pi(\tau) = \rho(s_0)\pi(a_0|s_0)\Pi_{t=1}^{H-1}\mathcal{P}(s_{t+1}|s_t, a_t)\pi(a_t|s_t) \ \forall \pi$ be the density of the distribution over the set of all trajectories induced by the policy and the transition dynamics (notice that this is a $\delta$ over one single trajectory because the MDP and the policy are deterministic in the analysed case). Then, to show that the policy ordering is preserved, we want to prove:

$$\mathbb{E}_{\tau \sim p_{\pi^*}}\left[\mathcal{R}'(\tau)\right] > \mathbb{E}_{\tau \sim p_\pi}\left[\mathcal{R}'(\tau)\right] \ \forall \pi \neq \pi^* \quad (6)$$

$$\mathbb{E}_{\tau \sim p_{\pi^*}}\left[\mathcal{R}(\tau) + \gamma^{h(\tau)}\phi(s_G) - \cancel{\phi(s_0)}\right] >$$
$$\mathbb{E}_{\tau \sim p_\pi}\left[\mathcal{R}(\tau) + \gamma^{h(\tau)}\phi(s_h^\pi) - \cancel{\phi(s_0)}\right] \ \forall \pi \neq \pi^*, \quad (7)$$

where we have used $\mathcal{R}'(\tau) = \mathcal{R}(\tau) + \gamma^H \phi(s_H) - \phi(s_0)$ [17] in (7), and $h(\tau) \leq H$ represents the number of steps from the initial state to the goal in $\tau$ that depends on the policy and environment dynamics. Assuming $\phi(s_G) = \Phi \ \forall \ s_G \in \mathcal{S}_G$, then choosing $\pi = \pi^{**}$ (the second-best policy w.r.t. $\mathcal{M}$) will maximize the right-hand side of the inequality and imply $s_h^{\pi^{**}} = s_G'$, so we can substitute $\phi(s_h^\pi) = \phi(s_h^{\pi^{**}}) = \Phi$:

$$\mathbb{E}_{\tau \sim p_{\pi^*}}\left[\mathcal{R}(\tau) + \gamma^{h(\tau)}\Phi\right] > \mathbb{E}_{\tau \sim p_{\pi^{**}}}\left[\mathcal{R}(\tau) + \gamma^{h(\tau)}\Phi\right]. \quad (8)$$

Now, (8) is always true because:

$$\mathbb{E}_{\tau \sim p_{\pi^*}} [\mathcal{R}(\tau)] > \mathbb{E}_{\tau \sim p_{\pi^{**}}} [\mathcal{R}(\tau)]$$

and

$$\mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \gamma \gamma^{h(\tau)-1} \Phi \right] > \mathbb{E}_{\tau \sim p_{\pi^{**}}} \left[ \gamma \gamma^{h(\tau)-1} \Phi \right]$$

$$\mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t, s_{t+1}) \right] >$$

$$\mathbb{E}_{\tau \sim p_{\pi^{**}}} \left[ \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t, s_{t+1}) \right] \quad (9)$$

$$\mathbb{E}_{\tau \sim p_{\pi^*}} [\mathcal{R}(\tau)] > \mathbb{E}_{\tau \sim p_{\pi^{**}}} [\mathcal{R}(\tau)],$$

where in (9) the reward is one in any goal as per the definition of goal-oriented MDP in Section 2.

Furthermore, iterating this process for any policy starting from the best and the second best up to the ones requiring $H-1$ and $H$ steps, we see that the ordering is preserved in the reshaped MDP. For what concerns all the other policies requiring more than $H$ steps to achieve the goal, let us restart from (7) selecting $\pi^H$ instead of $\pi^*$, and $\pi^{H+k}$ for any $k$ instead of $\pi$:

$$\mathbb{E}_{\tau \sim p_{\pi^H}} \left[ \mathcal{R}(\tau) + \gamma^H \Phi \right] > \mathbb{E}_{\tau \sim p_{\pi^{H+k}}} \left[ \mathcal{R}(\tau) + \gamma^H \phi(s_H) \right].$$

This inequality is true because $\mathbb{E}_{\tau \sim p_{\pi^H}} [\mathcal{R}(\tau)] > \mathbb{E}_{\tau \sim p_{\pi^{H+k}}} [\mathcal{R}(\tau)] = 0$ and $\Phi \geq \phi(s_H)$. $\square$

In a deterministic goal-oriented MDP, whenever we fix a horizon length for computational reasons, we are setting a limit on the sensitivity of our agent. Therefore, policies reaching the goal in $H + 1$ steps cannot be distinguished from those that reach the goal in more than $H + 1$ steps. Bearing this in mind, Theorem 3 tells us that, as long as the goal can be reached in at most $H$ steps, the solution to the reshaped problem is equivalent to that of the original problem.

Even though the family of MDPs considered by Theorem 3 has great relevance in practice, it is not able to model problems where there is any sort of environmental stochastic dynamic involved. For this reason, we proceed with our analysis in the context of general goal-oriented MDPs as well. We start by defining the probability of not having reached the goal in $H$ steps for an agent following policy $\pi$:

$$\Pr_{\tau \sim p_\pi} \{s_H \notin S_G\} = \mathbb{E}_{\tau \sim p_\pi} \left[ \mathbb{I}_{\{s_H \notin S_G\}} \right], \quad (10)$$

where $p_\pi(\tau) = \rho(s_0)\pi(a_0|s_0)\Pi_{t=1}^{H-1}\mathcal{P}(s_{t+1}|s_t, a_t)\pi(a_t|s_t)$ is the density of the distribution over trajectories induced by the policy and the transition dynamics. For convenience in notation, we define the event $E = \{s_H \notin S_G\}$. Furthermore, we assume the following:

**Assumption 1** (Minimal Performance Gap). Given $\pi_A \in \Pi_H$, where $\Pi_H$ is the set of all policies having a non-zero expected return within $H$ steps in $\mathcal{M}$, and $\pi_B$ a general policy in $\mathcal{M}$, the following performance gap is at least $\epsilon > 0$:

$$\inf_{(\pi_A, \pi_B) \in \mathcal{Z}} \mathbb{E}_{\tau \sim p_{\pi_A}} [\mathcal{R}(\tau)] - \mathbb{E}_{\tau \sim p_{\pi_B}} [\mathcal{R}(\tau)] > \epsilon,$$

where $\mathcal{Z} = \{\pi_A, \pi_B : \pi_A \in \Pi_H \wedge \mathbb{E}_{\tau \sim p_{\pi_A}} [\mathcal{R}(\tau)] > \mathbb{E}_{\tau \sim p_{\pi_B}} [\mathcal{R}(\tau)]\}$.

Now we are ready to state the existence of a potential function that is order preserving for all the policies in the relevant set $\Pi_H$ through the following theorem.

**Theorem 4** (Policy Ordering in Goal-Oriented Episodic MDPs). In the episodic context, let the horizon be $H$, $\mathcal{M}$ be a goal-oriented MDP, $\mathcal{M}'$ be its reshaped version obtained using a bounded potential $\phi : \mathcal{S} \to [0, \Phi]$. Moreover, let $\phi(s_G) = \Phi \geq C$, $\forall s_G \in \mathcal{S}_G$ with $C$ being:

$$\sup_{(\pi_A, \pi_B) \in \mathcal{Z}} \frac{\gamma^{H-1} \left( \mathbb{E}_{\tau \sim p_{\pi_B}} [\phi(s_H)\mathbb{I}_E] - \mathbb{E}_{\tau \sim p_{\pi_A}} [\phi(s_H)\mathbb{I}_E] \right)}{\mathbb{E}_{\tau \sim p_{\pi_A}} [\mathcal{R}(\tau)] - \mathbb{E}_{\tau \sim p_{\pi_B}} [\mathcal{R}(\tau)]},$$

where $\mathcal{Z}$ is defined as in Assumption 1. Then, the policy ordering in $\mathcal{M}$ is preserved in $\mathcal{M}'$ for all the policies in $\Pi_H$.

As can be seen from the above theorem (proof in Appendix A.2), in the context of general goal-oriented MDPs, we are not as lucky as with the deterministic case. Nevertheless, we are able to state a condition over the potential's maximum value, $\Phi$, that allows preserving the order for the set of relevant policies $\Pi_H$. This condition involves: the performance gap at the denominator; and two expected values depending on $\Pr_{\tau \sim p_{\pi_B}} E$ and $\Pr_{\tau \sim p_{\pi_A}} E$, which are the first and second terms at the numerator, respectively. Under Assumption 1, $C$ is bounded since the smallest denominator we can have is $\epsilon$ and the numerator is the difference between the expected values of a bounded function under two different distributions. Observe that assumptions on the performance gap between policies are common in the performance analysis of RL algorithms [3] and that, letting $H$ tend to infinity, $C$ goes to zero recovering the result of Ng et al. [28] since $\Phi$ is unconstrained.

The above theorems, requiring the potential function to be maximal in the goal states, avoid the drawback of Grześ [17][3]. Indeed, his approach, choosing $\phi(s_H) = 0$ in Eq. (3), leaves all trajectories not ending up in a goal without feedback on how far they terminated w.r.t. to it, impairing the exploration advantages of PBRS [8]. Finally, the above theorems have significant practical implications. Theorem 3 allows PBRS to be used without concern for the bias due to Eq. (3) in deterministic goal-oriented MDPs. Theorem 4, on the other hand, states the existence of a potential function that does not introduce bias in the context of stochastic goal-oriented MDPs, and, even though $C$ in practice cannot be known, in Section 6.1, we will see how even a small value will result in optimal performance for our agent. Together these two theorems may avoid computationally expensive approaches like the off-policy scheme proposed by Cipollone et al. [8]. Indeed, their primary focus is on goal-oriented MDPs and they require updating two different agents (effectively doubling computational complexity) to avoid a bias that is not actually there.

### 5.2 General MDPs

In the context of general MDPs, we can establish a condition that preserves policy ordering across the two problems.

**Theorem 5** (Policy Ordering in Reshaped Episodic MDPs). Let $\mathcal{M}$ be an MDP and $\mathcal{M}'$ be its reshaped version obtained using a bounded potential $\phi : \mathcal{S} \to [0, \Phi]$. Then, there exists a horizon

$$H > \log_\gamma \inf_{\substack{\pi_A, \pi_B : \\ \mathbb{E}_\rho[V^{\pi_A}(s)] > \mathbb{E}_\rho[V^{\pi_B}(s)]}} \frac{\mathbb{E}_\rho [V^{\pi_A}(s)] - \mathbb{E}_\rho [V^{\pi_B}(s)]}{\Phi + \frac{\bar{R}}{1-\gamma}}$$

such that the ordering in the reshaped MDP is preserved w.r.t. the infinite-horizon version of $\mathcal{M}$.

---

[3] Notice that, after reaching the goal, the potential can be set to zero to avoid any sort of convergence issue for the returns. Equivalently, in practice, the simulation ends there.

The condition stated in Theorem 5 (proof in Appendix A.3) is quite restrictive, because preserving the total ordering of policies is a strong equivalence requirement. Indeed, at the numerator, we have the performance gap between two arbitrary policies in the infinite-horizon version of $\mathcal{M}$, which, even if is greater than zero, may be very small, requiring an incredibly long horizon $H$. Therefore, let us consider the least restrictive version of total ordering, that is, preserving only the optimal policy.

**Corollary 6** (Preserving the Optimal Policy in Reshaped Episodic MDPs). Let $\mathcal{M}$ be an MDP and $\mathcal{M}'$ be its reshaped version obtained using a bounded potential $\phi : \mathcal{S} \rightarrow [0, \Phi]$. Then, there exists a horizon

$$H > \log_\gamma \frac{\mathbb{E}_\rho\left[V^{\pi^*}(s)\right] - \mathbb{E}_\rho\left[V^{\pi^{**}}(s)\right]}{\Phi + \frac{\bar{R}}{1-\gamma}},$$

where $\pi^{**}$ is the second-best policy, such that the optimal policy in the reshaped MDP is equivalent to the one in the infinite-horizon version of $\mathcal{M}$.

*Proof.* It is enough to select $\pi_A = \pi^*$ in Theorem 5. This implies that $\pi_B$ must be $\pi^{**}$. □

Now, the horizon length depends on the performance gap between the optimal policy and the second-best policy, a concept often used in the literature [3]. This is much less stringent than what stated in Theorem 5.

# 6 Evaluation

In this section, we experimentally evaluate the impact of abstractions coupled with PBRS on sample efficiency. This will be accomplished in the context of a goal-oriented navigation task and some complex ALE games.

## 6.1 Grid World

To corroborate the theoretical findings about goal-oriented MDPs described in the previous section, we evaluate the performance of Q-Learning, Abstraction-based PBRS (A-PBRS) Q-Learning, and Off-Policy Abstraction-based PBRS (OPA-PBRS) Q-Learning against each other. A-PBRS Q-Learning represents the Q-Learning algorithm with the reshaping described in Equation 4. OPA-PBRS Q-Learning is the off-policy approach proposed by Cipollone et al. [8]. This comparison was done over the 8-rooms goal-oriented navigation task [8] depicted in Figure 1. The agent has to move from the top-left room (red) to the right-most one (green) in order to accomplish the task. The abstraction involves considering only the rooms and their connections rather than the connections between the cells (see Appendix B.1). Further experimental details may be found in Appendix C.

This task has a stochastic transition function, so it falls under the scope of Theorem 4. Additionally, we cannot compute the theoretical value of $\Phi$. However, even setting $\Phi = 1$ allows us not only to achieve the optimum, but also to beat OPA-PBRS Q-Learning as can be seen in Figure 2a. Q-Learning performance is reported for completeness.
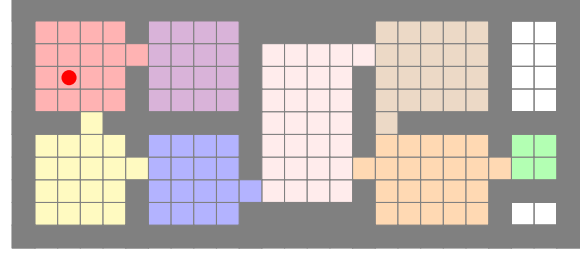


**Figure 1**: 8-rooms task. The red dot is the agent, the goal is in green, and the gray blocks represent walls. Each color is a macro-state exhaustively defining the aggregation function.

## 6.2 Arcade Learning Environment

We also show how the Abstraction-based PBRS approach behaves in much more challenging tasks, ALE games [4]. We have chosen three tasks of increasing level of difficulty and significantly different among each other in terms of dynamics: *Freeway*, *Q\*Bert*, and *Venture*. Each abstraction for these problems has been developed by simplifying the game neglecting enemies or entire areas/dynamics of the game itself. In the following section, we provide a complete and detailed description of *Freeway* and its abstraction, whereas we defer the reader to Appendix B.2 and B.3 for *Q\*Bert*, *Venture*, and their respective abstractions due to space constraints.

To tackle the above-mentioned tasks, we will use DQN (an overview of the algorithm can be found in Section 2). In the experiments, we will compare DQN based on fully-connected architectures against the same algorithm, but this time applied on a reshaped environment, where the reshaping is based on Equation 4. This is done to evaluate the contribution of the abstraction on sample efficiency. Due to the complexity of the tackled tasks, we will try to adhere as much as possible to theory. However, deviations from it will be allowed whenever these make the implementation feasible or they yield better performance (implementation details in Section 6.2.1 and Appendix B.2 and B.3). The previously mentioned algorithms, using fully-connected neural networks as value function estimators, have been trained on the RAM version of the above mentioned ALE games. This implies that the state space of those environments will be represented by a 128-dimensional vector whose components may take up to 256 values encoding various information about the current state of the game (e.g., coordinates of the agent, coordinates of the enemies, etc.). To further assess the contribution over sample efficiency and performance achievable through the abstractions, we will additionally compare the fully-connected version of DQN against its CNN version first introduced by Mnih et al. [27]. DQN based on CNN architectures will be trained straight from images (the more classical representation of the state space in ALE).

### 6.2.1 Abstracting Freeway

In *Freeway* (Figure 4a in the Appendix), the agent needs to cross a road avoiding the cars driving along the 10 different lanes. Every time the agent bumps into a car, it will get thrown backward a few steps, whereas upon crossing the entire road it will get one point and restart from the beginning. The game lasts 2 minutes and 16 seconds and the overall objective is to cross the entire road as many times as possible.

*Freeway* represents the simplest task we tackled in the context of the ALE benchmark because fully-connected neural networks are able to achieve more than 85% of the CNNs performance on their
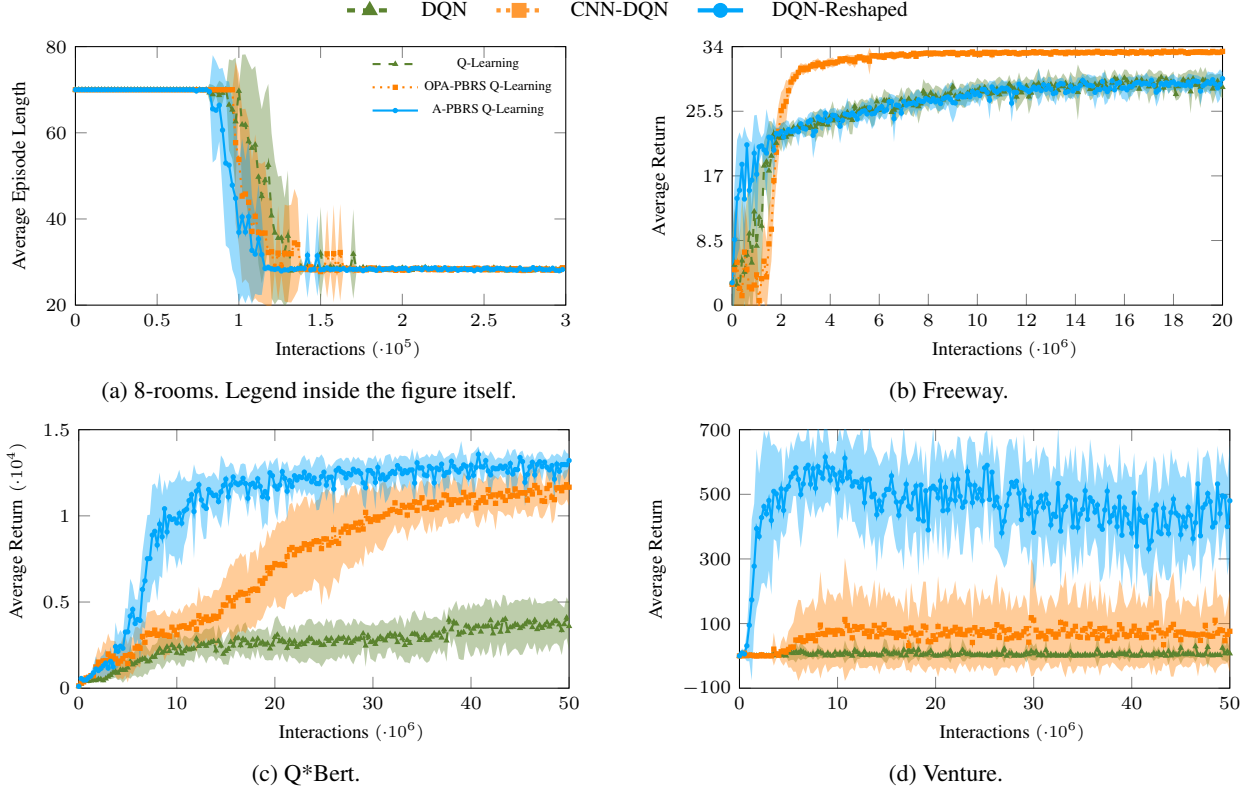
(a) 8-rooms. Legend inside the figure itself.

(b) Freeway.

(c) Q*Bert.

(d) Venture.

**Figure 2**: Average return $\pm 1$ standard deviation achieved by the algorithms computed using 10 independent runs.

own. In this task, we propose a very simple abstraction where we neglect the presence of the cars and we model a simple empty road to be crossed by the agent. The agent starts at position 6 (only the $y$ coordinate is considered because the agent cannot move horizontally) and has to cross position 180 by moving up, down, or doing nothing (1,2, and 0, respectively). Movement actions will increment or decrement the position of the agent by 1 in the respective direction. Since there are no cars in this abstraction, it is pointless to maximize the amount of times the agent crosses the road; it is sufficient to consider just one crossing. This abstraction has 177 states and its aggregation function is the following:

$$\alpha_s(s) = (s_{14}, 0)$$
$$\alpha_{s'}(s, s') = (s'_{14}, s'_{103} - s_{103}),$$

where we have used the $14^{th}$ and $103^{rd}$ coordinates ($y$ and score, respectively [1]) of the RAM vector representing the state space of the original task.[4] This aggregation function allows us to deal with a smaller abstraction because we do not have to consider multiple crossings of the whole road, but only one. Finally, this proposed abstraction is solved through VI and its solution is used to perform reshaping on the original task as per Equation 4.

### 6.2.2 Results

In Figures 2b, 2c, and 2d, we report the learning curves associated to the above described algorithms in the context of *Freeway*, *Q\*Bert*, and *Venture*, respectively. DQN with fully-connected neural networks and without reshaping is termed DQN, its equivalent version leveraging Abstraction-based PBRS is called DQN-Reshaped, and DQN based on CNN is called CNN-DQN.

---

[4] $\alpha_s$ aggregates $s$, whereas $\alpha_{s'}$ aggregates $s'$.

In the context of *Freeway*, as shown in Figure 2b, we observe higher sample efficiency of DQN-Reshaped w.r.t. DQN even though it is not able to achieve the same final performance as CNN-DQN. This fact is due to the usage of a very simplistic abstraction stemming from the assumption of no cars on the road (an exhaustive description of the abstraction may be found in Section 6.2.1). Unfortunately, accounting for the cars could not be done in the context of this work because of the complex dynamic triggered whenever the agent bumps into a car that could not be reverse engineered from interactions nor from the assembly code of the game.

In *Q\*Bert*, Figure 2c, flat fully-connected neural networks, when coupled with PBRS and abstractions (exhaustive description in Appendix B.2), are able to achieve the same performances of CNNs with approximately one fourth the amount of interactions. Additionally, DQN-Reshaped manages to have the best average performance at the end of the learning process.

In the context of *Venture* (see Figure 2d), not only we have greater sample efficiency, but reshaping and abstractions (exhaustive description in Appendix B.3) allow fully-connected architectures to achieve the best performance by a great margin w.r.t. all the other solutions.

Finally, the quality of the abstraction, in terms of how well it approximates the original task, plays an important role. Indeed, we can observe a significant difference between the improvement in performance achieved in the context of *Freeway* w.r.t. *Q\*Bert* or *Venture*. This is due to the different amount of knowledge provided by a simplistic abstraction (in *Freeway*) compared to a more complex one (in *Q\*Bert* and *Venture*). Furthermore, the longer the sequence of actions to be executed in order to experience reward, the higher the benefit of reshaping. As a matter of fact, experiencing reward in *Freeway* and *Q\*Bert* is quite easy, whereas, in *Venture*, the agent is required to execute a long plan in order to get any positive feedback from the

environment. This is reflected by the fact that abstractions in *Venture* provide the best improvement. All the experimental details may be found in Appendix C.[5]

## 7 Related Works

Ng et al. [28] were the first authors to formalize theoretical properties of reward shaping showing that PBRS preserves optimal and near-optimal policies. Wiewiora [44] proves that, under the same set of collected experience, for Q-learning, SARSA, and other TD-based algorithms, a suitable action-value function initialization is equivalent to PBRS. The first work formally discussing the bias introduced by reshaping when dealing with finite horizons can be found in [15]. In the context of POMDPs, they show that the evaluation of policies becomes the same as the horizon length tends to infinity. Furthermore, they provide a condition under which two policies, one optimizing the reshaped POMDP and one the original, will be different from each other. This is done to understand when the optimal policy in the finite-horizon reshaped POMDP differs from the optimal policy of the original finite-horizon POMDP. In contrast, our Theorem 5 produces a condition to control the bias w.r.t. to the infinite-horizon MDP, the true problem of interest. Besides, Theorem 3 and 4 establish that, performing reshaping when the horizon is finite does not change the policy ordering w.r.t. the starting finite-horizon problem for the goal-oriented case (under some assumptions if the MDP is stochastic). Grześ [17], instead, tackle the above-mentioned bias by setting the potential function at the end of the trajectory to zero, hence leveraging non-stationary potential functions [14], whereas Forbes et al. [16] using terms that do not depend on the future actions of the agent.

Along another direction, Cheng et al. [7] and Burden et al. [5] learn their potential. The former learns it beforehand with offline data using Monte-Carlo Regression, whereas the latter through direct interactions with the environment. Instead, to get the potential, we do not need data or the original environment to interact with while training.

There is not much theoretical analysis in the literature making explicit the performance improvements attainable with reward shaping. One exception is the work of Gupta et al. [18], where they provide a regret expression for a new version of UCBVI called UCBVI-Shaped. However, they assume to have a value function estimator upper-bounding the optimal value function when multiplied by a parameter $\beta$ and their approach is meant for finite state and action problems. Fully self-supervised reward-shaping techniques [45, 26] may not be able to deal with sparse reward settings because they lack an explicit exploration incentive. Therefore, to mitigate this limitation, Devidze et al. [13] propose to learn an intrinsic reward function combined with a count-based exploration bonus aiming at accelerating the agent's endeavors in maximizing the extrinsic reward. Similarly, Wang et al. [40] propose an intrinsic reward to encourage exploration. This is achieved defining the discrepancy between states as differences of potentials, where the potential function is a learned inverse dynamic bi-simulation metric. Along a different direction, Hu et al. [20] study how to learn the best way to utilize a reshaping function in order to solve the RL problem at hand. They do this through a bi-level optimization approach, where at the first level they optimize the policy to maximize the reshaped rewards, and at the second level they optimize the reshaping function for actual reward maximization. Differently w.r.t. these works we study the bias introduced by PBRS and explicitly leverage abstractions to choose potentials.

The work presented in this paper also draws inspiration from abstractions and how they are used in classical planning [23]. The idea is to automatically build a simplified version of the problem before the search starts. This is obtained by abstracting away all but a (small) part of the task, resulting in a new problem that can be optimally solved fast [19]. The plans that solve these abstract problems are then used as heuristics that help guiding the search algorithm toward the goal in the actual problem. Closely related to the above-mentioned way of using abstractions is the concept of abstractions for MDPs [9, 35]. However, they usually deal with finite state and action MDPs. Furthermore, they directly employ the solutions obtained from the abstractions onto the real problem like Wang et al. [41, 42].

Finally, Cipollone et al. [8] propose an off-policy RL scheme to leverage hierarchies of abstractions and improve sample efficiency. They mainly focus on goal-oriented MDPs for which they provide a theoretical analysis on the exploration quality of the behavioral policy in their proposed solution that depends on the size of the abstraction state space. Hence, with sufficiently large abstractions (e.g., on the order of $10^5$ states) easily achievable in the context of ALE games (see Section 6), the bound is not helpful. Furthermore, thanks to Theorem 3 and 4, we can avoid using their off-policy scheme, which requires twice the computational complexity to maintain convergence to the optimum. In the context of general MDPs, their approach still represents a sound way to preserve convergence guarantees. However, whenever we cannot afford the computational burden of training two separate agents, thanks to Theorem 5, we can try to control the amount of bias we introduce in the reshaped MDP.

## 8 Conclusions and Future Work

In this paper, we have tackled the problem of sample efficiency in RL through abstraction and PBRS, an idea that was recently introduced in the literature by Cipollone et al. [8]. Contrary to them, we have embraced the bias of PBRS in the context of finite horizons, investigated its effects on policy ordering for goal-oriented and general MDPs, and assessed the contribution to sample efficiency of abstractions coupled with PBRS in highly complex domains such as the ALE benchmark. Given the impressive sample efficiency achieved, devising ways to learn abstractions automatically from interactions becomes an important avenue for future developments. This could be accomplished by taking inspiration from the Hierarchical Reinforcement Learning literature [22]. Moreover, being able to build abstractions over state spaces represented by images is another interesting future direction for this work that could be pursued through learning an action schema directly from images [2]. Finally, the idea of using abstractions as a guide to solve the original task is deeply rooted into the planning literature and the concept of heuristics; it would be worth investigating their connections [6].

## Disclaimer

---

[5] Code will be made available upon request.

# References

[1] A. Anand, E. Racah, S. Ozair, Y. Bengio, M.-A. Côté, and R. D. Hjelm. Unsupervised state representation learning in atari. In *Advances in Neural Information Processing Systems*, 2019.

[2] M. Asai and A. Fukunaga. Classical planning in deep latent space: Bridging the subsymbolic-symbolic boundary. In *AAAI Conference on Artificial Intelligence*, 2018.

[3] P. Auer, T. Jaksch, and R. Ortner. Near-optimal regret bounds for reinforcement learning. In *Advances in Neural Information Processing Systems*, 2008.

[4] 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*, 2013.

[5] J. Burden, S. K. Siahroudi, and D. Kudenko. Latent property state abstraction for reinforcement learning. In *AAMAS Workshop on Adaptive Learning Agents*, 2021.

[6] G. Canonaco, A. Pozanco, and D. Borrajo. Projection abstractions in planning under the lenses of abstractions for mdps. *arXiv preprint arXiv:2412.02615*, 2024.

[7] C.-A. Cheng, A. Kolobov, and A. Swaminathan. Heuristic-guided reinforcement learning. In *Advances in Neural Information Processing Systems*, 2021.

[8] R. Cipollone, G. De Giacomo, M. Favorito, L. Iocchi, and F. Patrizi. Exploiting multiple abstractions in episodic rl via reward shaping. In *AAAI Conference on Artificial Intelligence*, 2023.

[9] E. Congeduti and A. Frans. A cross-field review of state abstraction for markov decision processes. In *34th Benelux Conference on Artificial Intelligence (BNAIC) and the 30th Belgian Dutch Conference on Machine Learning (Benelearn)*, 2022.

[10] J. C. Culberson and J. Schaeffer. Pattern databases. *Computational Intelligence*, 1998.

[11] M. P. Deisenroth, G. Neumann, J. Peters, et al. A survey on policy search for robotics. *Foundations and Trends® in Robotics*, 2013.

[12] C. D'Eramo, D. Tateo, A. Bonarini, M. Restelli, and J. Peters. Mushroomrl: Simplifying reinforcement learning research. *Journal of Machine Learning Research*, 2021.

[13] R. Devidze, P. Kamalaruban, and A. Singla. Exploration-guided reward shaping for reinforcement learning under sparse rewards. In *Advances in Neural Information Processing Systems*, 2022.

[14] S. M. Devlin and D. Kudenko. Dynamic potential-based reward shaping. In *International Conference on Autonomous Agents and Multi-Agent Systems*, 2012.

[15] A. Eck, L.-K. Soh, S. Devlin, and D. Kudenko. Potential-based reward shaping for finite horizon online pomdp planning. *Autonomous Agents and Multi-Agent Systems*, 2016.

[16] G. C. Forbes, N. Gupta, L. Villalobos-Arias, C. M. Potts, A. Jhala, and D. L. Roberts. Potential-based reward shaping for intrinsic motivation. In *International Conference on Autonomous Agents and Multi-Agent Systems*, 2024.

[17] M. Grześ. Reward shaping in episodic reinforcement learning. In *International Conference on Autonomous Agents and Multi-Agent Systems*, 2017.

[18] A. Gupta, A. Pacchiano, Y. Zhai, S. M. Kakade, and S. Levine. Unpacking reward shaping: Understanding the benefits of reward engineering on sample complexity. In *Advances in Neural Information Processing Systems*, 2022.

[19] P. Haslum, A. Botea, M. Helmert, B. Bonet, S. Koenig, et al. Domain-independent construction of pattern database heuristics for cost-optimal planning. In *AAAI Conference on Artificial Intelligence*, 2007.

[20] Y. Hu, W. Wang, H. Jia, Y. Wang, Y. Chen, J. Hao, F. Wu, and C. Fan. Learning to utilize shaping rewards: A new approach of reward shaping. In *Advances in Neural Information Processing Systems*, 2020.

[21] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.

[22] M. Klissarov and M. C. Machado. Deep laplacian-based options for temporally-extended exploration. In *International Conference on Machine Learning*, 2023.

[23] C. A. Knoblock. Automatically generating abstractions for planning. *Artificial Intelligence*, 1994.

[24] A. Lazaric. Transfer in reinforcement learning: a framework and a survey. *Reinforcement Learning: State-of-the-Art*, 2012.

[25] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 1998.

[26] F. Memarian, W. Goo, R. Lioutikov, S. Niekum, and U. Topcu. Self-supervised online reward shaping in sparse-reward environments. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2021.

[27] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 2015.

[28] A. Y. Ng, D. Harada, and S. Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *International Conference on Machine Learning*, 1999.

[29] OpenAI, C. Berner, G. Brockman, B. Chan, V. Cheung, P. Dębiak, C. Dennison, D. Farhi, Q. Fischer, S. Hashme, C. Hesse, R. Józefowicz, S. Gray, C. Olsson, J. Pachocki, M. Petrov, H. P. de Oliveira Pinto, J. Raiman, T. Salimans, J. Schlatter, J. Schneider, S. Sidor, I. Sutskever, J. Tang, F. Wolski, and S. Zhang. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.

[30] M. Papini, D. Binaghi, G. Canonaco, M. Pirotta, and M. Restelli. Stochastic variance-reduced policy gradient. In *International Conference on Machine Learning*, 2018.

[31] M. L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.

[32] G. A. Rummery and M. Niranjan. *On-line Q-learning using connectionist systems*. University of Cambridge, Department of Engineering Cambridge, UK, 1994.

[33] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.

[34] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 2018.

[35] R. A. Starre, M. Loog, E. Congeduti, and F. A. Oliehoek. An analysis of model-based reinforcement learning from abstracted observations. *Transactions on Machine Learning Research*, 2023.

[36] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, 1999.

[37] M. E. Taylor and P. Stone. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 2009.

[38] A. Touati, J. Rapin, and Y. Ollivier. Does zero-shot reinforcement learning exist? In *International Conference on Learning Representations*, 2022.

[39] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 2019.

[40] Y. Wang, M. Yang, R. Dong, B. Sun, F. Liu, et al. Efficient potential-based exploration in reinforcement learning using inverse dynamic bisimulation metric. In *Advances in Neural Information Processing Systems*, 2024.

[41] Z. Wang, X. Xiao, Z. Xu, Y. Zhu, and P. Stone. Causal dynamics learning for task-independent state abstraction. In *International Conference on Machine Learning*, 2022.

[42] Z. Wang, C. Wang, X. Xiao, Y. Zhu, and P. Stone. Building minimal and reusable causal state abstractions for reinforcement learning. In *AAAI Conference on Artificial Intelligence*, 2024.

[43] C. J. Watkins and P. Dayan. Q-learning. *Machine Learning*, 1992.

[44] E. Wiewiora. Potential-based shaping and q-value initialization are equivalent. *Journal of Artificial Intelligence Research*, 2003.

[45] Z. Zheng, J. Oh, and S. Singh. On learning intrinsic rewards for policy gradient methods. In *Advances in Neural Information Processing Systems*, 2018.

# A  Proofs

## A.1  Proof of Proposition 1

*Proof.* Applying VI in the original MDP $\mathcal{M}$, where $\mathcal{V}_{n,\mathcal{M}}^{src}$ is the output after $n$ iterations, yields the following bound by the $\gamma$-contraction property:

$$||\mathcal{V}_{\mathcal{M}}^* - \mathcal{V}_{n,\mathcal{M}}^{src}||_\infty \leq \gamma^n ||\mathcal{V}_{\mathcal{M}}^* - \mathcal{V}_0||_\infty. \tag{11}$$

Applying VI in the reshaped MDP $\mathcal{M}'$, where $\mathcal{V}_{n',\mathcal{M}}^{dst}$ is the output after $n'$ iterations brought back into the original MDP through Eq. (1), yields the following bound by the $\gamma$-contraction property:

$$||\mathcal{V}_{\mathcal{M}}^* - \phi - \mathcal{V}_{n',\mathcal{M}}^{dst} + \phi||_\infty \leq \gamma^{n'} ||\mathcal{V}_{\mathcal{M}}^* - \phi - \mathcal{V}_0||_\infty. \tag{12}$$

Now, imposing inequalities (11) and (12) both smaller than $\epsilon$, we get:

$$||\mathcal{V}_{\mathcal{M}}^* - \mathcal{V}_{n,\mathcal{M}}^{src}||_\infty \leq \gamma^n ||\mathcal{V}_{\mathcal{M}}^* - \mathcal{V}_0||_\infty \leq \epsilon$$
$$||\mathcal{V}_{\mathcal{M}}^* - \mathcal{V}_{n',\mathcal{M}}^{dst}||_\infty \leq \gamma^{n'} ||\mathcal{V}_{\mathcal{M}}^* - \phi - \mathcal{V}_0||_\infty \leq \epsilon.$$

Using the fact that $\gamma < 1$, we can rewrite the inequalities above as:

$$n \geq \log_\gamma \left( \frac{\epsilon}{||\mathcal{V}_{\mathcal{M}}^* - \mathcal{V}_0||_\infty} \right)$$
$$n' \geq \log_\gamma \left( \frac{\epsilon}{||\mathcal{V}_{\mathcal{M}}^* - \phi - \mathcal{V}_0||_\infty} \right)$$

then, since we assumed $||\mathcal{V}_{\mathcal{M}}^* - \phi - \mathcal{V}_0||_\infty \leq ||\mathcal{V}_{\mathcal{M}}^* - \mathcal{V}_0||_\infty$, which can be made true through $\phi$ being an approximation of $\mathcal{V}_{\mathcal{M}}^*$, we have:

$$\log_\gamma \left( \frac{\epsilon}{||\mathcal{V}_{\mathcal{M}}^* - \phi - \mathcal{V}_0||_\infty} \right) \leq \log_\gamma \left( \frac{\epsilon}{||\mathcal{V}_{\mathcal{M}}^* - \mathcal{V}_0||_\infty} \right)$$

hence, VI applied on the reshaped problem will require a smaller minimum number of iterations to guarantee an $\epsilon$-optimal solution. Notice that, the superscripts *src* and *dst* are used to disambiguate the two solutions of VI that may be different. $\square$

## A.2  Proof of Theorem 4

*Proof.* Let $\pi^*$ be the optimal policy in the original MDP, and $p_\pi(\tau) = \rho(s_0)\pi(a_0|s_0)\Pi_{t=1}^{H-1}\mathcal{P}(s_{t+1}|s_t, a_t)\pi(a_t|s_t)$ $\forall \pi$ be the density of the distribution over the set of all trajectories induced by the policy and the transition dynamics.

To show that the policy ordering is preserved, we want to prove:

$$\mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \mathcal{R}'(\tau) \right] > \mathbb{E}_{\tau \sim p_\pi} \left[ \mathcal{R}'(\tau) \right] \ \forall \pi \neq \pi^*, \pi \in \Pi_H \tag{13}$$

$$\mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \mathcal{R}(\tau) + \gamma^{h(\tau)}\phi(s_h^{\pi^*}) - \cancel{\phi(s_0)} \right] >$$
$$\mathbb{E}_{\tau \sim p_\pi} \left[ \mathcal{R}(\tau) + \gamma^{h(\tau)}\phi(s_h^\pi) - \cancel{\phi(s_0)} \right] \ \forall \pi \neq \pi^*, \pi \in \Pi_H, \tag{14}$$

where we have used $\mathcal{R}'(\tau) = \mathcal{R}(\tau) + \gamma^H\phi(s_H) - \phi(s_0)$ [17] in (14), and $h(\tau) \leq H$ represents the random variable associated to the number of steps from the initial state to the goal in $\tau$ that depends on the policy and environment dynamics. We know that $\mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \mathcal{R}(\tau) \right] > \mathbb{E}_{\tau \sim p_\pi} \left[ \mathcal{R}(\tau) \right] \ \forall \pi \neq \pi^*$, so we only have to prove:

$$\mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \gamma^{h(\tau)}\phi(s_h^{\pi^*}) \right] > \mathbb{E}_{\tau \sim p_\pi} \left[ \gamma^{h(\tau)}\phi(s_h^\pi) \right] \ \forall \pi \neq \pi^*, \pi \in \Pi_H.$$

Then:

$$\mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \gamma^{h(\tau)}\phi(s_h^{\pi^*})\mathbb{I}\{s_h^{\pi^*} \in S_G\} \right] + \mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \gamma^H\phi(s_H)\mathbb{I}\{s_H \notin S_G\} \right] >$$
$$\mathbb{E}_{\tau \sim p_\pi} \left[ \gamma^{h(\tau)}\phi(s_h^\pi)\mathbb{I}\{s_h^\pi \in S_G\} \right] + \mathbb{E}_{\tau \sim p_\pi} \left[ \gamma^H\phi(s_H)\mathbb{I}\{s_H \notin S_G\} \right],$$

where $\mathbb{I}\{s_h^\pi \in S_G\}$ is the indicator function selecting only the trajectories that end in a goal state in at most $H$ steps, and $\mathbb{I}\{s_H \notin S_G\} = \mathbb{I}_E$ ($E$ defined in Section 5) selects all the trajectories not ending in the goal within $H$ steps. Now, since $\mathbb{E}_{\tau \sim p_\pi} \left[ \gamma^{h(\tau)}\phi(s_h^\pi)\mathbb{I}\{s_h^\pi \in S_G\} \right] = \gamma\mathbb{E}_{\tau \sim p_\pi} \left[ \mathcal{R}(\tau) \right] \Phi \ \forall \pi$ because in goal-oriented MDPs $\mathbb{E}_{\tau \sim p_\pi} \left[ \gamma^{h(\tau)-1}\mathbb{I}\{s_h^\pi \in S_G\} \right] = \mathbb{E}_{\tau \sim p_\pi} \left[ \gamma^{h(\tau)-1}r(s_{h(\tau)-1}, a_{h(\tau)-1}, s_{h(\tau)}) \right]$, then:

$$\gamma \left( \mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \mathcal{R}(\tau) \right] - \mathbb{E}_{\tau \sim p_\pi} \left[ \mathcal{R}(\tau) \right] \right) \Phi > \gamma \left( \mathbb{E}_{\tau \sim p_\pi} \left[ \gamma^{H-1}\phi(s_H)\mathbb{I}_E \right] - \mathbb{E}_{\tau \sim p_{\pi^*}} \left[ \gamma^{H-1}\phi(s_H)\mathbb{I}_E \right] \right)$$

$$\Phi > \frac{\gamma^{H-1}\left(\mathbb{E}_{\tau\sim p_\pi}\left[\phi(s_H)\mathbb{I}_E\right] - \mathbb{E}_{\tau\sim p_{\pi^*}}\left[\phi(s_H)\mathbb{I}_E\right]\right)}{\mathbb{E}_{\tau\sim p_{\pi^*}}\left[\mathcal{R}(\tau)\right] - \mathbb{E}_{\tau\sim p_\pi}\left[\mathcal{R}(\tau)\right]}.$$

The above is true because we chose $\Phi \geq C$. We can repeat all the above reasoning for $\pi^{**}$ instead of $\pi^*$ such that $\pi \notin \{\pi^*, \pi^{**}\}$, and so on up until we exhaust all the policies in $\Pi_H$.

Now taking $\pi_i \in \Pi_H$ and $\pi_o \notin \Pi_H$, we have to prove:

$$\mathbb{E}_{\tau\sim p_{\pi_i}}\left[\mathcal{R}(\tau) + \gamma^{h(\tau)}\phi(s_h^{\pi_i})\right] > \mathbb{E}_{\tau\sim p_{\pi_o}}\left[\gamma^H\phi(s_H)\right] \ \forall \pi_o \notin \Pi_H$$

that holds true if:

$$\Phi > \frac{\gamma^{H-1}\left(\mathbb{E}_{\tau\sim p_{\pi_o}}\left[\phi(s_H)\mathbb{I}_E\right] - \mathbb{E}_{\tau\sim p_{\pi_i}}\left[\phi(s_H)\mathbb{I}_E\right]\right)}{\mathbb{E}_{\tau\sim p_{\pi_i}}\left[\mathcal{R}(\tau)\right]}$$

that is smaller than the chosen $C$ for any $\pi_i, \pi_o$. $\qquad\square$

## A.3  Proof of Theorem 5

*Proof.* Given the definition of $p_\pi(\tau)$ (see Section 5), let $p_\pi^\infty(\tau)$ be its limit for $H$ going to infinity. Furthermore, let $\mathbb{E}_{p_{\pi_A}}\left[\mathcal{R}'(\tau_A)\right]$ and $\mathbb{E}_{p_{\pi_B}}\left[\mathcal{R}'(\tau_B)\right]$ be the performance of two generic policies in the reshaped MDP such that $\mathbb{E}_\rho\left[V^{\pi_A}(s)\right] > \mathbb{E}_\rho\left[V^{\pi_B}(s)\right]$ in the original MDP. Then, for the ordering to be preserved in the reshaped MDP, the following must hold:

$$\mathbb{E}_{p_{\pi_A}}\left[\mathcal{R}'(\tau_A)\right] > \mathbb{E}_{p_{\pi_B}}\left[\mathcal{R}'(\tau_B)\right] \ \forall \pi_A, \pi_B$$

$$\mathbb{E}_{p_{\pi_A}}\left[\mathcal{R}(\tau_A) + \gamma^H\phi(s_H) - \cancel{\phi(s_0)}\right] > \mathbb{E}_{p_{\pi_B}}\left[\mathcal{R}(\tau_B) + \gamma^H\phi(s_H) - \cancel{\phi(s_0)}\right] \ \forall \pi_A, \pi_B \qquad (15)$$

$$\mathbb{E}_{p_{\pi_A}}\left[\mathcal{R}(\tau_A)\right] - \mathbb{E}_{p_{\pi_B}}\left[\mathcal{R}(\tau_B)\right] > \mathbb{E}_{p_{\pi_B}}\left[\gamma^H\phi(s_H)\right] - \mathbb{E}_{p_{\pi_A}}\left[\gamma^H\phi(s_H)\right] \ \forall \pi_A, \pi_B,$$

where in (15) we have used: $\mathcal{R}'(\tau) = \mathcal{R}(\tau) + \gamma^H\phi(s_H) - \phi(s_0)$. The potential is bounded, $\phi \in [0, \Phi]$. Hence, given $\mathcal{R}(\tau^\infty) = \sum_{t=H}^\infty \gamma^t r(s_t, a_t, s_{t+1})$, we have:

$$\mathbb{E}_{p_{\pi_A}}\left[\mathcal{R}(\tau_A)\right] - \mathbb{E}_{p_{\pi_B}}\left[\mathcal{R}(\tau_B)\right] > \gamma^H\Phi \ \forall \pi_A, \pi_B \qquad (16)$$

$$\mathbb{E}_{p_{\pi_A}}\left[\mathcal{R}(\tau_A)\right] - \mathbb{E}_{p_{\pi_B}}\left[\mathcal{R}(\tau_B)\right] + \mathbb{E}_{p_{\pi_A}^\infty}\left[\mathcal{R}(\tau_A^\infty)\right] - \mathbb{E}_{p_{\pi_B}^\infty}\left[\mathcal{R}(\tau_B^\infty)\right] > \gamma^H\Phi + \mathbb{E}_{p_{\pi_A}^\infty}\left[\mathcal{R}(\tau_A^\infty)\right] - \mathbb{E}_{p_{\pi_B}^\infty}\left[\mathcal{R}(\tau_B^\infty)\right]$$

$$\mathbb{E}_\rho\left[V^{\pi_A}(s)\right] - \mathbb{E}_\rho\left[V^{\pi_B}(s)\right] > \gamma^H\Phi + \mathbb{E}_{p_{\pi_A}^\infty}\left[\mathcal{R}(\tau_A^\infty)\right] \qquad (17)$$

$$\mathbb{E}_\rho\left[V^{\pi_A}(s)\right] - \mathbb{E}_\rho\left[V^{\pi_B}(s)\right] > \gamma^H\Phi + \frac{\gamma^H}{1-\gamma}\bar{R} \qquad (18)$$

$$\gamma^H < \inf_{\substack{\pi_A, \pi_B: \\ \mathbb{E}_\rho[V^{\pi_A}(s)] > \mathbb{E}_\rho[V^{\pi_B}(s)]}} \frac{\mathbb{E}_\rho\left[V^{\pi_A}(s)\right] - \mathbb{E}_\rho\left[V^{\pi_B}(s)\right]}{\Phi + \frac{\bar{R}}{1-\gamma}}$$

$$H > \log_\gamma \inf_{\substack{\pi_A, \pi_B: \\ \mathbb{E}_\rho[V^{\pi_A}(s)] > \mathbb{E}_\rho[V^{\pi_B}(s)]}} \frac{\mathbb{E}_\rho\left[V^{\pi_A}(s)\right] - \mathbb{E}_\rho\left[V^{\pi_B}(s)\right]}{\Phi + \frac{\bar{R}}{1-\gamma}},$$

where, in (16), we have used the fact that the potential is bounded within $[0, \Phi]$ to upper bound the right-hand side. Instead, in (17) and (18), we have used the fact that the reward function is bounded within $[0, \bar{R}]$ to upper bound the right-hand side. $\qquad\square$

# B  Abstractions

## B.1  Grid World

In this section, we describe the 8-rooms abstraction, reported as a graph in Figure 3. Each node is associated with a room with the topology of the graph and the coloring faithfully representing Figure 1 (this also defines the aggregation function). The connections between nodes are bidirectional and represent the available actions. Additionally, any action has a 90% chance of success and a 10% failure probability. This means that, with probability 0.9, the agent ends up in the intended target node, with probability 0.1, instead, it will stay where it is. Closely following Cipollone et al. [8], this stochastic behavior is implemented in the original 8-rooms environment of Figure 1 as well, but, in this case, the agent has a 4% failure probability upon executing any action, and, when it fails, one among all the available actions is executed at random. Finally, from the practical point of view, it is worth pointing out that, in order for the potential to be maximum in the goal state as required by Theorem 4, we add in the abstraction a fictitious done action that gives us a reward of 1 and transitions us from $\alpha(s_G)$ into a fictitious terminal state of the abstraction itself (this also simplifies the implementation when dealing with multiple goals in the abstraction because all have to transition to the same terminal state).
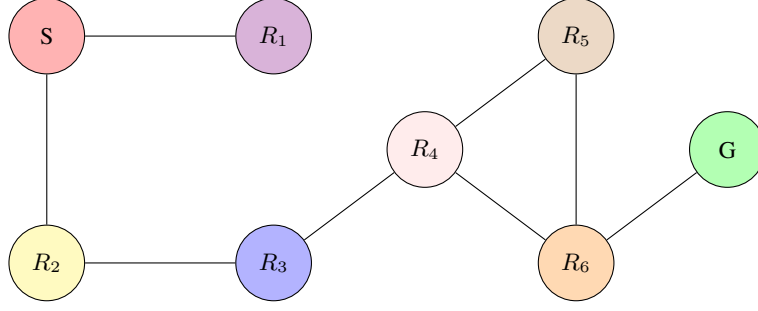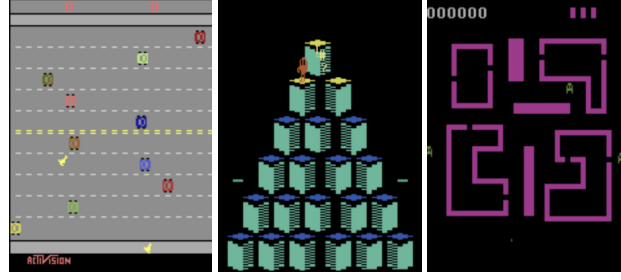
**Figure 3**: Grid World abstraction.



(a) Freeway.          (b) Q*Bert.          (c) Venture.

**Figure 4**: The tackled ALE games.

## B.2   Q*Bert

In *Q*Bert* (Figure 4b), the agent can jump onto some tiles over a grid. By jumping onto tiles, the agent is able to change their color. The goal of the agent is to change every tile on the grid to a target color of interest. This will make it succeed and go to the next level. The agent is awarded 25 points for each tile turned to the target color. Furthermore, it is awarded a bonus of 3100 points on level completion. Unfortunately, while the agent attempts to accomplish its task, there are enemies attempting at its life or undoing its tile-coloring work.

### B.2.1   Abstracting Q*Bert

In order to abstract away what the agent has to accomplish in *Q*Bert*, we neglect all the enemies and consider only the positions of the tiles on the grid. In a nutshell, the agent can move between tiles, changing their color upon visiting them. There are 21 positions the agent can be in, four movement actions (left, right, up, and down), and one action that simulates the death of the agent, bringing it back to the initial position without undoing the coloring that has been done so far. The goal of the agent is to visit every tile in order to color it with the target color. The reward in this abstraction will be 0 for every step except upon succeeding with the task, where the reward is 1. This abstraction has $1,172,830$ states. We solved it using VI, and the $\epsilon$-optimal value function is used to guide the agent through reshaping in the original task.

### B.2.2   The Abstraction Details

The nodes in Figure 5 represent the tiles the agent has to color. The agent starts at tile 1 that has coordinates $(77, 25)$ (see Table 1) and it can move according to the connections reported in the graph of Figure 5. For instance, if the agent is at node 5, it can move to 3 by executing the *up* action, move to 2 by executing *left*, move to 9 by executing *right*, move to 8 by executing *down*, and, finally, *die* there by moving back directly to 1. To avoid cluttering, the dying action has been reported only for node 6 in Figure 5; it applies to all the nodes, and it does not change the color of the landing tile (which is always 1). All the other movement actions will change the color of the landing tile, and they have been reported only for node 5 for the sake of readability (the pattern is the same for every node except for the borders where missing arcs mean that the associated actions are not available for execution). The state space of the abstraction is then represented by the $x$ and $y$ coordinates of the agent corresponding to the coordinates reported in Table 1 and 21 variables representing the colors of the tiles that start all at the same source color and need to be transitioned all to the target by the agent. Formally, being $s$ the state vector of the original task, we have:

$$\bar{s} = (s_{43}, s_{67}, s_{21}, s_{52}, s_{54}, s_{83}, s_{85}, s_{87}, s_{98}, s_{100}, s_{102}, s_{104}, s_1, s_3, s_5, s_7, s_9, s_{32}, s_{34}, s_{36}, s_{38}, s_{40}, s_{42})$$

$$= (\underbrace{x, y,}_{\text{Coordinates}} \underbrace{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21}_{\text{Node IDs reported in Figure 5}})$$
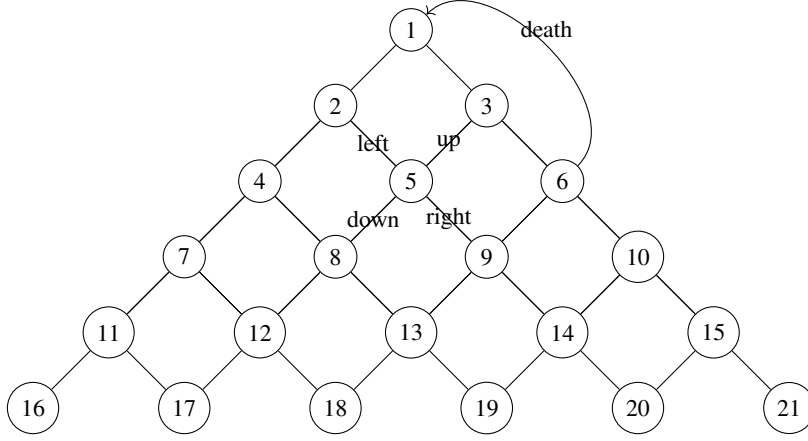
**Figure 5**: Q*Bert Transition Graph. The Node ID correspond to the position in $x, y$ coordinates reported in Table 1.

**Table 1**: Node ID to position mapping.

| Node ID | $(x, y)$ | Node ID | $(x, y)$ |
|---------|----------|---------|----------|
| 1 | (77, 25) | 12 | (53, 137) |
| 2 | (65, 53) | 13 | (77, 137) |
| 3 | (93, 53) | 14 | (105, 137) |
| 4 | (53, 81) | 15 | (129, 137) |
| 5 | (77, 81) | 16 | (16, 165) |
| 6 | (105, 81) | 17 | (41, 165) |
| 7 | (41, 109) | 18 | (65, 165) |
| 8 | (65, 109) | 19 | (93, 165) |
| 9 | (93, 109) | 20 | (117, 165) |
| 10 | (117, 109) | 21 | (141, 165) |
| 11 | (29, 137) | | |

### B.2.3 The Coloring Across Different Levels

In the context of the original task, on every level, we have a different source-target color pair that can be detected automatically at the first few interactions of the agent with the new level. This is done by monitoring the transitions the agent is experiencing. Whenever we go to a new level ($57^{th}$ component of the RAM vector), we grab the color of tile 1, then as soon as we transition from 1 to either 2 or 3, and the color of one of those two tiles has changed w.r.t. the one we grabbed in 1, we save the source-target color pair. This automatic mechanism allows us to make our abstraction agnostic to the color change across different levels, so that the abstraction is kept fixed throughout the whole learning process and does not depend on the level the agent is at. Finally, notice that, handling all possible different source-target color pairs not only enables us to have a single abstraction for all the levels, but it also enables a pre-processing step that translates all the source-target color pairs into a single one to help the network better leverage similarities across levels. In Figure 2c, we have reported results using abstraction and pre-processing combined, whereas in Figure 6, we report an ablation study to better assess the contribution of the single components: only abstraction and only pre-processing. As we may see, with reshaping, we are able to greatly improve performances w.r.t. the basic DQN solution. The same may be stated for what concerns leveraging the pre-processing at the cost of a greater variance. The two techniques combined yield the best performance surpassing even CNN-based DQN as stated in Section 6.2.2.
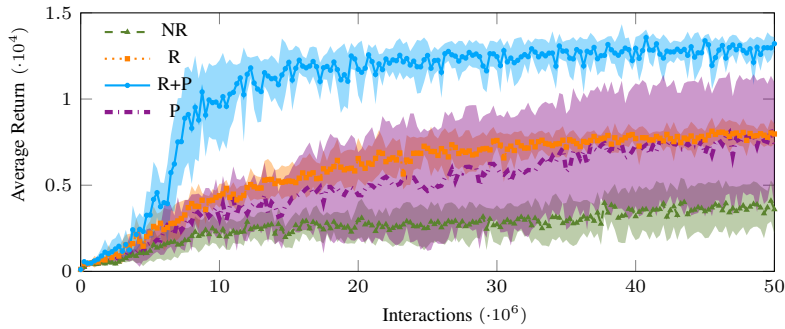


**Figure 6**: Average return achieved by simple fully-connected DQN without Reshaping (NR), with Reshaping (R), with Pre-processing (P), with Reshaping and Pre-processing (R+P) $\pm 1$ standard deviation computed using 10 independent runs in *Q*Bert*.

**Algorithm 1** Q*Bert F function
___

1: **Input:** $\langle \bar{s}, a, \bar{s}' \rangle$ the transition, $\mathcal{G}$ the abstraction's transition graph, $\phi = \mathcal{V}_\epsilon$ the potential function
2: **if** $\bar{s}_{x,y}$ not in $\mathcal{G}$ **then**
3:     **return** $0$
4: **end if**
5: **if** $\bar{s}_{x,y} == \bar{s}'_{x,y}$ **then**
6:     **return** $0$
7: **end if**
8: **if** $a$ not in $\mathcal{G}[s]$ **then**
9:     **if** $\bar{s}$ in $\phi$ **then**
10:         **return** $-\phi(\bar{s})$
11:     **else**
12:         **return** $0$
13:     **end if**
14: **else if** $\bar{s}$ in $\phi$ and **next_state**$(\bar{s}, \mathcal{G}[\bar{s}_{x,y}][a])$ in $\phi$ **then**
15:     $F = \gamma\phi(\textbf{next\_state}(\bar{s}, \mathcal{G}[\bar{s}_{x,y}][a])) - \phi(\bar{s})$
16:     **if** F==0 **then**
17:         **return** $0.1$
18:     **end if**
19:     **return** $F$
20: **else**
21:     **return** $0$
22: **end if**
___

### B.2.4   The Aggregation Function

In order to complete the map between $\mathcal{S}$ and $\bar{\mathcal{S}}$, we need to account for the transition dynamics between the nodes connected via the actions *up*, *right*, *down*, and *left* in Figure 5.[6] In our abstraction that transition is instantaneous, but in the real task it is not. Let us see an example:

$$\underbrace{(77,25),}_{\bar{s}^1_{x,y}} \underbrace{(78,24),(82,20),(86,22),(89,27),(89,35),(89,43),(89,51),(93,53)}_{\bar{s}^2_{x,y}}.$$

This represents the transitions undergone by the $x$ and $y$ coordinates of the agent when moving from tile 1 to tile 3 in the original task. This is the same transition, but at level 2:

$$\underbrace{(77,25),}_{\bar{s}^1_{x,y}} \underbrace{(81,21),(85,21),(89,25),(89,33),(89,41),(89,49),(93,53)}_{\bar{s}^2_{x,y}}.$$

Furthermore, to make things even harder, they keep on changing throughout the different levels. Since this transition is mandatory, and the agent cannot deviate from it, we are going to define the aggregation function taking the endpoints as representatives and assigning the in-between states as shown above. Formally, let $\sigma_1, \ldots, \sigma_k = \sigma \in \varsigma$ be a possible sequence of succeeding $x$ and $y$ coordinates allowed by the game between $(77, 25)$ and $(93, 53)$, then the aggregation function for this family of transitions will map $(77, 25) \rightarrow (77, 25)$ and $\{\sigma, (93, 53)\} \rightarrow (93, 53)\ \forall \sigma \in \varsigma$, and so on for every other node in Figure 5. In order to induce the above described partitioning, we are going to use the action chosen by the agent. Indeed, given the agent is at $(77, 25)$, if it executes the action *right*, it will observe the appropriate $\sigma \in \varsigma$ and then eventually $(93, 53)$. This implies that we are able to properly reshape between $(77, 25)$ and the first $x$ and $y$ coordinate pair, $\sigma_1$, in the $\sigma$ sequence. The reshaping of all the other transitions, namely, $\langle \sigma_2, \sigma_3 \rangle, \ldots, \langle \sigma_k, (93, 53) \rangle$, is set to 0.

The pseudo-code that shows an implementation of the $F$ function (see Section 2.2) for *Q\*Bert* is shown in Algorithm 1, where we are assuming that $\mathcal{G}$ is an hash-map representing the transition graph reported in Figure 5 without the death arcs. $\mathcal{V}_\epsilon$ is the solution produced by VI applied to the abstraction. Finally, the function **next_state**, given the current abstract state and the $x$ and $y$ coordinates the agents ends up at after executing $a$, returns the abstract next state by appropriately coloring the required tile in $\bar{s}$. In Line 1, we do not reshape if the starting state coordinates are out of our abstraction. In Line 5, we do not reshape if the agent stays still; this is done not to incentivize continuous movement. In Lines 8 to 13, if the agent moved out of the abstraction from a known state, we penalize it with $-\phi(s)$. This is because it is jumping out of the grid and consequently dying. Otherwise, if the state is unknown, we do not reshape (Line 12). In Line 14, we handle the transition between one abstract tile and another. If the value of $F$ is 0 we boost it to 0.1 (this happens every time we take an action that is optimal according to $\mathcal{V}_\epsilon$), otherwise, we simply return $F$.[7] In all other remaining cases 0 is returned.

Changing the reshaping function as described above deviates from theory and it may introduce bias, but as we have seen in the experiments of Figure 2c, the achieved performance is better than CNN-based DQN. Additionally, it is worth noticing that, in Algorithm 1, for the sake of readability, we did not report the pseudo-code to handle all the possible source-target color pairs. Finally, we normalize the reward of the original task dividing it by 500 (the highest reward obtainable in a single step). This does not change the task from the optimization point of view and make everything on a similar scale w.r.t. reshaping values.

___

[6] Note that this mapping induces a partition over the original state space where each $\bar{s}$ is a class representative.
[7] Notice that our abstraction is a goal-oriented MDP solved with the same discount factor used in the original task, see Section C.
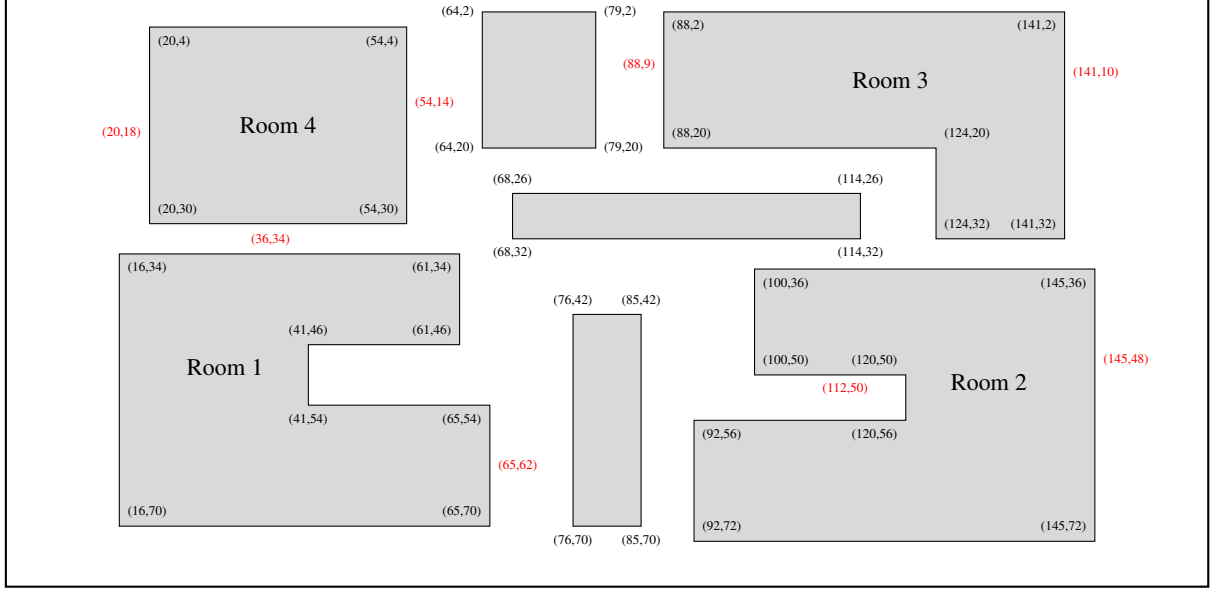
**Figure 7**: Venture abstraction: main hall.

## B.3 Venture

In *Venture* (Figure 4c), the agent is placed into a main hall that must be explored in order to find some rooms where treasures are hidden. In order to progress, the agent has to enter each room, collect its treasure, and exit the room without being killed by the enemies and pitfalls present in every room (main hall included). The agent will receive 200 points for each collected treasure and 100 points for each killed enemy in the treasure rooms (these last points are awarded only if the killing happens after the treasure has been collected). This reward system makes the task a sparse reward problem, which is much harder to deal with for RL agents.

### B.3.1 Abstracting Venture

In the context of *Venture*, the objective of the agent is twofold: exploration to find the rooms' doors in the main hall and treasure collection in the rooms. Whenever the agent enters a room, collects its treasure, and exits without being killed, the room gets locked and can no longer be accessed. Therefore, we can define a simple abstraction modeling only the problem of finding the rooms' doors. This can be done by considering the position of the agent, whether a room is locked or not, the room the agent is currently in, and the boundaries of the main hall. Again, as we did in *Q*Bert*, we will neglect the enemies that chase the agent in the main hall. In the context of the above-described abstraction, visiting a room via one of its doors makes that room transition into the previously-mentioned locked state. Hence, we are abstracting away what the agent must do inside the rooms. The available actions are the movements along the eight main cardinal directions, which will change the components of the agent's positional vector in the appropriate way (increment/decrement by one unit or no change). The goal is to visit all the rooms. When this happens, the agent receives a reward of 1, 0 everywhere else. This abstraction has 106, 929 states and is again solved $\epsilon$-optimally with VI. The obtained value function is used to do reshaping as stated in Equation 4. The aggregation function is:

$$\alpha(s) = (s_{85}, s_{26}, s_{90}, s_{17}),$$

which means we are considering only the $85^{th}, 26^{th}, 90^{th}$, and $17^{th}$ components ($x$ and $y$ coordinates, room id, and locked rooms, respectively [1]) of the original RAM vector.

In the same spirit as the abstraction above, we could build a second layer of abstraction for the rooms within the main hall. Unfortunately, neglecting the monsters does not allow us to build a good abstraction due to their aggressive behaviour. Indeed, reshaping in these rooms without accounting for the monsters would provide misleading feedback to the agent because the shortest path to the treasure and back out of the room heavily depends on the enemies' positions. Thus, we lean on the agent to find a way to kill the monsters, collect the treasure, and exit.

### B.3.2 The Abstraction Details

In Figure 7, we can find a complete description of the main-hall abstraction with all the rooms, the obstacles, and, in red, the doors to access the rooms. The agent starts in position (68, 76). Furthermore, notice that whenever we are over a door position, by playing the appropriate action, we will enter the room and end up over the corresponding door position on the other side. For instance, if the agent is at (65, 62), by playing the action *left*, it will end up being at (129, 63) in Room 1 (see Table 2). From there, the agent will only be able to go back by playing *right* effectively locking the room, as if it was solved. In Algorithm 2, we report the pseudo-code describing an implementation of the reshaping mechanism in the context of *Venture*. Here, the only potential sources of bias are due to Line 3 of Algorithm 2, where we return 0 if the agent stays still to avoid incentives to continuous movement, and in Algorithm 3, where we return 0 if one of the two abstract states is out of the abstraction. This was done because during preliminary experimental studies it appeared to give better performances.

**Table 2**: Main hall door to room door position mapping

| Main Hall Door | Action | Room Door |
|:---:|:---:|:---:|
| (65, 62) | *left* | (129, 63) |
| (36, 34) | *bottom* | (58, 11) |
| (112, 50) | *up* | (62, 18) |
| (145, 48) | *left* | (129, 13) |
| (141, 10) | *left* | (129, 15) |
| (88, 9) | *right* | (31, 15) |
| (54, 14) | *left* | (117, 39) |
| (20, 18) | *right* | (43, 39) |

---

**Algorithm 2** Venture F function

1: **Input:** $\langle s, a, s' \rangle$ the transition, $\phi^0 = \mathcal{V}_\epsilon^0$ the potential function associated with the main hall
2: **if** $s_{85,26} == s'_{85,26}$ **then**
3:     **return** $0$
4: **end if**
5: **if** $s_{90} == 8$ **then**
6:     **return** $F_8\left(s, a, s', \phi^0\right)$
7: **end if**
8: **return** $0$

---

**Algorithm 3** Venture $F_8$ function

1: **Input:** $\langle s, a, s' \rangle$ the transition, $\phi = \mathcal{V}_\epsilon^0$ the potential function associated to the main hall
2: $\bar{s} = s_{85,26,90,17}$, $\bar{s}' = s'_{85,26,90,17}$
3: **if** $\bar{s}$ in $\phi$ and $\bar{s}'$ in $\phi$ **then**
4:     $F = \gamma\phi(\bar{s}') - \phi(\bar{s})$
5:     **return** $10^3 F$
6: **end if**
7: **return** $0$

---

## C   Experimental Details

### C.1   Value Iteration

In this section, we will describe all the VI configurations to solve the above-described abstractions. For all the abstractions, we have an $\epsilon$ of $10^{-7}$, whereas the discount factor $\gamma$ is 0.98 for *Freeway* and *Venture*, 0.99 for *Q\*Bert*, and 0.9 for the Grid World.

### C.2   Deep Q-Networks

Regarding the ALE benchmark, we used Mushroom-RL [12] implementation of DQN with the hyperparameters reported in Table 3. Notice that, for *Freeway*, the **test exploration rate** has been lowered to 0.01 because it yielded better performance in evaluation for the CNN architecture. Furthermore, being *Freeway* a very simple task to solve, we lowered the **evaluation frequency** to 100000 and the **test samples** to 50000. This enables more frequent probing of the agent's performance, allowing a higher resolution for the learning curve, which, in turn, allows us to better visualize and understand the learning dynamics, especially in the early epochs where it matters the most.

Adam [21] is used as optimizer with a learning rate of 0.0001 and $(\beta_1, \beta_2) = (0.9, 0.999)$. In the case of RAM states, we have fully-connected neural networks with the following structure: one input layer of $128 \cdot 4$ neurons, two hidden layers of 128 and 256 neurons (first and second hidden layers, respectively), both with ReLU activation functions, and an output layer whose size depends on the number of actions (*Freeway* has 3, *Q\*Bert* 6, and *Venture* 18). In the case of image states, we have an input shape of $(4, 84, 84)$, and the CNN structure is as follows: 3 convolutional layers and 2 linear layers. The first convolutional layer has 4 input channels, 32 output channels, a kernel size of 8, and stride of 4; the second one has 32 input channels, 64 output channels, a kernel size of 4, and stride of 2; the third one has 64 input channels, 64 output channels, a kernel size of 3, and stride of 1. The first linear layer has 3136 neurons and the second 512. The output is the number of actions, which depends on the task and it is already stated above. Every layer has a ReLU activation function except the output layer that does not have any.

### C.3   Q-Learning

In this section, we report all the hyperparameters used to train Q-Learning in the context of the 8-rooms task. The learning rate starts at 0.85 for vanilla Q-Learning, 0.7 for our solution, and 1 for Cipollone et al. [8] solution. The different starting points have been optimized for each algorithm independently using a grid search over the following values: $[1, 0.85, 0.7, 0.55, 0.4, 0.25, 0.1]$. They all decay to 0.01. The parameter $\epsilon$ of the $\epsilon$-greedy policy starts at 1 and decays to 0.1. Both the learning rate and the exploration parameter decay over 300000 time steps (the whole learning process). Finally, we have an episode maximum length of 70 and a discount factor of 0.98.

**Table 3**: DQN Hyperparameters

| Parameter | Value | Description |
|---|---|---|
| $\gamma$ | 0.99 | The discount factor |
| initial replay size | 50000 | Initial number of samples in the replay memory |
| max replay size | 500000 | Replay memory max capacity |
| batch size | 32 | Number of samples for each fit of the network |
| history Length | 4 | Number of frames a state is made of |
| target update frequency | 10000 | Number of collected samples before a new update of the target network |
| train frequency | 4 | Number of collected samples before a new fit of the neural network |
| final exploration frame | 1000000 | The number of samples while the exploration rate decays |
| initial exploration rate | 1 | The initial exploration rate value |
| final exploration rate | 0.1 | The final exploration rate value |
| max no-op actions | 30 | Maximum amount of no-ops executed at the beginning of each episode |
| max steps | 50000000 | Total number of interactions |
| evaluation frequency | 250000 | Total number of interactions before evaluating the policy |
| test samples | 125000 | Interactions during policy evaluation |
| test exploration rate | 0.05 | Exploration rate used while evaluating the policy |