

A Survey Analyzing Generalization in Deep Reinforcement Learning

Ezgi Korkmaz

*University College London
London, United Kingdom*

Abstract

Reinforcement learning research obtained significant success and attention with the utilization of deep neural networks to solve problems in high dimensional state or action spaces. While deep reinforcement learning policies are currently being deployed in many different fields from medical applications to large language models, there are still ongoing questions the field is trying to answer on the generalization capabilities of deep reinforcement learning policies. In this paper, we will formalize and analyze generalization in deep reinforcement learning. We will explain the fundamental reasons why deep reinforcement learning policies encounter overfitting problems that limit their generalization capabilities. Furthermore, we will categorize and explain the manifold solution approaches to increase generalization, and overcome overfitting in deep reinforcement learning policies. From exploration to adversarial analysis and from regularization to robustness our paper provides an analysis on a wide range of subfields within deep reinforcement learning with a broad scope and in-depth view. We believe our study can provide a compact guideline for the current advancements in deep reinforcement learning, and help to construct robust deep neural policies with higher generalization skills.

1. Introduction

The performance of reinforcement learning algorithms (Watkins, 1989; Sutton, 1984, 1988) has been boosted with the utilization of deep neural networks as function approximators (Mnih et al., 2015). Currently, it is possible to learn deep reinforcement learning policies that can operate in large state and/or action space MDPs (Silver et al., 2017; Vinyals et al., 2019). This progress consequently resulted in building reasonable deep reinforcement learning policies that can play computer games with high dimensional state representations (e.g. Atari, StarCraft), solve complex robotics control tasks, design algorithms (Mankowitz et al., 2023; Fawzi et al., 2022), guide large language models (OpenAI, 2023; Google Gemini, 2023), and play some of the most complicated board games (e.g. Chess, Go) (Schrittwieser et al., 2020). However, deep reinforcement learning algorithms also experience several problems caused by their overall limited generalization capabilities. Some studies demonstrated these problems via adversarial perturbations introduced to the state observations of the policy (Huang et al., 2017; Kos and Song, 2017; Korkmaz, 2022; Korkmaz and Brown-Cohen, 2023), several focused on exploring the fundamental issues with function approximation, estimation biases in the state-action value function (Thrun and Schwartz, 1993; van Hasselt, 2010), or with new architectural design ideas (Wang et al., 2016). The fact that we are not able to completely explore the entire MDP for high dimensional state representation MDPs, even with deep neural networks as function approximators, is one of the root problems that

limits generalization. On top of this, some portion of the problems are directly caused by the utilization of deep neural networks and thereby the intrinsic problems inherited from their utilization (Goodfellow et al., 2015; Szegedy et al., 2014; Korkmaz, 2022, 2024b).

In order to address open questions on generalization in deep reinforcement learning, there needs to be some commonly agreed standard of what is meant by generalization. Currently, different aspects of generalization are considered in various subfields either working on the fundamental questions regarding or the applications of deep reinforcement learning. We take the point of view in this paper that these various aspects can, and should, be described and studied in a unified way. In particular, we argue that the various approaches to generalization can be succinctly classified based on which part of the Markov Decision Process is expected to vary. We make this classification formal and unify how much current work on generalization in deep reinforcement learning fits clearly into the classification we introduce. In this paper we will focus on generalization in deep reinforcement learning and the underlying causes of the limitations deep reinforcement learning research currently faces. In particular, we will try to answer the following questions:

- *How can we formalize the concept of generalization in deep reinforcement learning?*
- *What is the role of exploration in overfitting for deep reinforcement learning?*
- *What are the causes of overestimation bias observed in state-action value functions?*
- *What has been done to overcome the overfitting problems that deep reinforcement learning algorithms have encountered so far, and to enable deep neural policies to generalize to non-stationary complex environments?*

To answer these questions we will go through research connecting several subfields in reinforcement learning on the problems and corresponding proposed solutions regarding generalization. In this paper we introduce a formal definition of generalization and categorization of the different methods used to both achieve and assess generalization, and use it to systematically summarize and consolidate the current body of research. We further describe the issue of value function overestimation, and the role of exploration in overfitting in reinforcement learning. Furthermore, we explain new emerging research areas that can potentially target these questions in the long run including meta-reinforcement learning and lifelong learning. The objective of the paper is to introduce a formal generalization definition and provide a compact overview and unification of the current advancements and limitations in the field.

2. Preliminaries on Deep Reinforcement Learning

The aim in deep reinforcement learning is to learn a policy via interacting with an environment in a Markov Decision Process (MDP) that maximize expected cumulative discounted rewards. An MDP is represented by a tuple $\mathcal{M} = (S, A, \mathcal{P}, r, \rho_0, \gamma)$, where S represents the state space, A represents the action space, $r : S \times A \rightarrow \mathbb{R}$ is a reward function, $\mathcal{P} : S \times A \rightarrow \Delta(S)$ is a transition probability kernel, ρ_0 represents the initial state distribution, and γ represents the discount factor. The objective in reinforcement learning is to learn a policy $\pi : S \times A \rightarrow \mathbb{R}$ which maps states to probability distributions on actions in order to maximize the expected cumulative reward $R = \mathbb{E} \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t)$ where

$a_t \sim \pi(s_t, \cdot)$, $s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)$. The temporal difference updates achieves this objective by updating the value function $V(s)$ (Sutton, 1984, 1988)

$$V(s_t) \leftarrow V(s_t) + \alpha[r(s_{t+1}, a) + \gamma V(s_{t+1}) - V(s_t)] \quad (1)$$

While Equation 1 represents the one-step temporal difference update, i.e. TD(0), it is further possible to consider multi-step TD which focuses on multi-step return, i.e. TD(λ). In Q -learning the goal is to learn the optimal state-action value function (Watkins, 1989)

$$Q^*(s, a) = r(s, a) + \sum_{s' \in \mathcal{S}} \mathcal{P}(s' | s, a) \max_{a' \in \mathcal{A}} Q^*(s', a'). \quad (2)$$

This is achieved via iterative Bellman update (Bellman, 1957; Bellman and Dreyfus, 1959) which updates $Q(s_t, a_t)$ by

$$Q(s_t, a_t) + \alpha[\mathcal{R}_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)].$$

Thus, the optimal policy is determined by choosing the action $a^*(s) = \arg \max_a Q(s, a)$ in state s . The optimal Bellman operator is (Bellman, 1957)

$$\mathcal{B}Q(s, a) := \mathbb{E}[r(s, a)] + \gamma \mathbb{E}_{\mathcal{P}}[\max_{a'} Q(s', a')]$$

In high dimensional state space or action space MDPs the optimal policy is decided via a function-approximated state-action value function represented by a deep neural network. The loss function in deep reinforcement learning is the quadratic difference between that target network and the current state-action value function.

$$\mathcal{L}_i(\theta_i) = \mathbb{E}_{e \sim \mathcal{D}}[(r(s, a) + \gamma \max_{a'} Q(s', a', \theta_{\text{target}}) - Q(s, a, \theta_i))^2]$$

where \mathcal{D} is the experience replay buffer (Lin, 1993) in which the experiences $e = \{s_t, a_t, r_t, s_{t+1}\}$ sampled from $\mathcal{D} = \{e_1, e_2, \dots, e_N\}$. The loss function is optimized by taking the gradient with respect function approximation weights

$$\begin{aligned} \theta_{i+1} = \theta_i + \alpha & (r(s_t, a_t, s_{t+1}) + \gamma \mathcal{Q}(s_{t+1}, \arg \max_a \mathcal{Q}(s_{t+1}, a; \theta_{i-1}^{\text{target}}); \theta_{i-1}^{\text{target}}) \\ & - \mathcal{Q}(s_t, a_t; \theta_i)) \nabla_{\theta_i} \mathcal{Q}(s_t, a_t; \theta_i). \end{aligned}$$

In a parallel line of algorithm families the policy itself is directly parametrized by π_θ (Sutton et al., 1999), and the gradient estimator used in learning is

$$g = \mathbb{E}_t [\nabla_\theta \log \pi_\theta(s_t, a_t) (Q(s_t, a_t) - \max_a Q(s_t, a))]$$

where $Q(s_t, a_t)$ refers to the state-action value function at time step t . The algorithms that focus on directly parameterizing the policy try to solve the following optimization problem (Schulman et al., 2015).

$$\max_{\theta} \mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}; a \sim \pi_{\theta_{\text{old}}}(s, \cdot)} \left[\frac{\pi_\theta(s, a)}{\pi_{\theta_{\text{old}}}(s, a)} Q_{\theta_{\text{old}}}(s, a) \right] \text{ subject to } \mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}} \mathcal{D}_{KL}(\pi_\theta(s, \cdot) || \pi_{\text{old}}(s, \cdot)) \leq \delta$$

3. How to Achieve Generalization?

3.1 Generic Reinforcement Learning Algorithm

To be able to understand and analyze the connection between different approaches to achieve generalization first we will provide a clear definition intended to capture the behavior of a generic reinforcement learning algorithm.

Definition 3.1 (*Generic reinforcement learning algorithm*). A reinforcement learning training algorithm \mathcal{A} learns a policy π by interacting with an MDP \mathcal{M} . We divide up the execution of \mathcal{A} into discrete time steps as follows. At each time t , the algorithm has a current policy π_t , observes a state s_t , takes an action $a_t \sim \pi_t(s_t, \cdot)$, and observes a transition to state $s'_t \sim \mathcal{P}(\cdot \mid s_t, a_t)$ with corresponding reward $r_t = r(s_t, a_t, s'_t)$. We define the history of algorithm \mathcal{A} in MDP \mathcal{M} to be the sequence $H_t = (\pi_0, s_0, a_0, s'_0, r_0), \dots, (\pi_t, s_t, a_t, s'_t, r_t)$ of all the transitions observed by the algorithm so far. We require that the policy π_t and state s_t at time t are a function only of H_{t-1} , i.e. the transitions observed so far by \mathcal{A} . At time $t = T$, the algorithm stops and outputs the policy $\pi = \pi_T$. We use the notation \mathbb{A} to denote the set of reinforcement learning training algorithms and Π to denote the set of policies π in an MDP \mathcal{M} .

Intuitively, a reinforcement learning algorithm has a current policy π_t , performs a sequence of queries (s_t, a_t) to the MDP, and observes the resulting state transitions and rewards. In order to be as generic as possible, the definition makes no assumptions about how the algorithm chooses the sequence of queries, other than that $a_t \sim \pi_t(s_t, \cdot)$. Notably, if taking action a_t in state s_t leads to a transition to state s'_t , there is no requirement that $s_{t+1} = s'_t$. Indeed, the only assumption is that s_{t+1} and π_{t+1} may depend only on H_t , the history of transitions observed so far. This allows the definition to capture deep reinforcement learning algorithms, which may choose to query states and actions in a complex way as a function of previously observed state transitions.

3.2 Base Generalization in Deep Reinforcement Learning

We next introduce a basic metric capturing how well an algorithm generalizes given a fixed amount of interaction with a given MDP.

Definition 3.2 (*Base generalization*). Given an MDP $\mathcal{M} = (S, A, P, r, \rho_0, \gamma)$, let π_T and $\hat{\pi}_T$ be policies output by training algorithms taking T steps. The base generalization \mathcal{G}^{base} is the difference between the expected discounted cumulative rewards obtained by policy π_T and $\hat{\pi}_T$ in \mathcal{M} .

$$\begin{aligned} \mathcal{G}^{base}(\pi_T, \hat{\pi}_T) = & \mathbb{E}_{a_t \sim \pi_T(s_t, \cdot)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}) \right] \\ & - \mathbb{E}_{\hat{a}_t \sim \hat{\pi}_T(\hat{s}_t, \cdot)} \left[\sum_{t=0}^{\infty} \gamma^t r(\hat{s}_t, \hat{a}_t, \hat{s}_{t+1}) \right] \end{aligned}$$

The base generalization definition captures how well an algorithm can generalize to unseen states and transitions, given only access to T interactions with the MDP \mathcal{M} . Hence, in base generalization the role of exploration is exceedingly dominant and this will be further explained in Section 5.

3.3 Algorithmic Generalization

Based on the definition of generic reinforcement learning algorithm, we will now further define the different approaches proposed to achieve generalization. At a high level, the approaches we will discuss will be divided into two classes:

- I. Techniques that solely modify the training algorithm,
- II. Techniques that directly modify the MDP (i.e. learning environment, training data) that forms the interactions of the training algorithm with the learning environment.

Our first definition formalizes the techniques that solely modify the training algorithm.

Definition 3.3 (*Algorithmic generalization*). Let \mathcal{A} be a training algorithm that takes an MDP as input and outputs a policy. Given an MDP $\mathcal{M} = (S, A, P, r, \rho_0, \gamma)$, an algorithmic generalization method $\mathcal{G}_{\mathbb{A}}$ is given by a function $F : \mathbb{A} \rightarrow \mathbb{A}$ that runs the algorithm $F(\mathcal{A})$ in the MDP \mathcal{M} .

Algorithmic generalization captures modifications to the training algorithm itself that can range from the choice of optimization methods or regularizers, to update rules for the policy.

3.4 Generalization Through Rewards

Definition 3.4 (*Rewards transforming generalization*). Let \mathcal{A} be a training algorithm that takes as input an MDP and outputs a policy. Given an MDP $\mathcal{M} = (S, A, P, r, \rho_0, \gamma)$, a rewards transforming generalization method \mathcal{G}_R is given by a sequence of functions $F_t : (\Pi \times S \times A \times S \times \mathbb{R})^t \times \mathbb{R} \rightarrow \mathbb{R}$. The method attempts to achieve generalization by running \mathcal{A} on MDP \mathcal{M} , but modifying the rewards at each time t to be $\hat{r}_t(s_t, a_t, s'_t) = F_{t-1}(H_{t-1}, r_t)$, where H_{t-1} is the history of algorithm \mathcal{A} when running with the transformed rewards.

In particular, a method under the rewards transforming generalization category runs the original algorithm to train the policy, but modifies the observed rewards. The instances of these techniques will be mentioned and explained in Section 6.2, i.e. direct function regularization, in Section 5, i.e. the role of exploration in overfitting, and in Section 8, i.e. transfer in reinforcement learning.

3.5 Generalization Through Observations

Following the definition of reward transforming generalization we define state transforming generalization which is one of the canonical approaches for achieving generalization in deep reinforcement learning. The instances of generalization through observations will be categorized and explained in detail in Section 6.1, i.e. data augmentation, and Section ??, i.e. the adversarial perspective for deep neural policy generalization.

Definition 3.5 (*State transforming generalization*). Let \mathcal{A} be a training algorithm that takes as input an MDP and outputs a policy. Given an MDP $\mathcal{M} = (S, A, P, r, \rho_0, \gamma)$, a state transforming generalization method \mathcal{G}_S is given by a sequence of functions $F_t : (\Pi \times S \times A \times S \times \mathbb{R})^t \times S \rightarrow S$. The method attempts to achieve generalization by running

\mathcal{A} on MDP \mathcal{M} , but modifying the state chosen at time t to be $\hat{s}_t = F_{t-1}(H_{t-1}, s_t)$, where H_{t-1} is the history of algorithm \mathcal{A} when running with the transformed states.

3.6 Generalization Through Environment Dynamics

Another category of algorithms that tries to achieve generalization in deep reinforcement learning focuses on achieving this objective through environment dynamics transformation. The methods focusing on generalization through environment dynamics will be referred to and explained in Section 6.2, i.e. direct function regularization.

Definition 3.6 (*Transition probability transforming generalization*). Let \mathcal{A} be a training algorithm that takes as input an MDP and outputs a policy. Given an MDP $\mathcal{M} = (S, A, P, r, \rho_0, \gamma)$, a transition probability transforming generalization method \mathcal{G}_P is given by a sequence of functions $F_t : (\Pi \times S \times A \times S \times \mathbb{R})^t \times (S \times A \times S) \rightarrow \mathbb{R}$. The method attempts to achieve generalization by running \mathcal{A} on MDP \mathcal{M} , but modifying the transition probabilities at time t to be $\hat{P}(s_t, a_t, s'_t) = F_{t-1}(H_{t-1}, s_t, a_t, s'_t)$, where H_{t-1} is the history of algorithm \mathcal{A} when running with the transformed transition probabilities.

3.7 Generalization Through Policy

The last type of generalization method we define is based on directly modifying the current policy used by the algorithm to select actions at each time step. We will explain the instances of the techniques that focus on generalization through policy in Section 5, i.e. the role of exploration in overfitting, and Section 7, i.e. meta reinforcement learning and meta gradients.

Definition 3.7 (*Policy transforming generalization*). Let \mathcal{A} be a training algorithm that takes as input an MDP and outputs a policy. Given an MDP $\mathcal{M} = (S, A, \mathcal{P}, r, \rho_0, \gamma)$, a policy transforming generalization method \mathcal{G}_π is given by a sequence of functions $F_t : (\Pi \times S \times A \times S \times \mathbb{R})^t \times S \times \Delta(A) \rightarrow \Delta(A)$. The method attempts to achieve generalization by running \mathcal{A} on MDP \mathcal{M} , but modifying the current policy by which \mathcal{A} chooses the action at time t to be $\hat{\pi}_t(s_t, \cdot) = F_{t-1}(H_{t-1}, s_t, \pi_t(s_t, \cdot))$, where H_{t-1} is the history of algorithm \mathcal{A} when running with the transformed policy.

3.8 Assessing Generalization

All the definitions so far categorize methods to modify either training algorithms and/or the MDP, i.e. learning environment, training data, in order to achieve generalization. However, many such methods for modifying training algorithms have a corresponding method which can be used to assess the generalization capabilities of a trained policy. Our final definition captures this correspondence.

Definition 3.8 (*Generalization testing*). Let $\hat{\pi}$ be a trained policy for an MDP \mathcal{M} . Let F_t be a sequence of functions corresponding to a generalization method from one of the previous definitions. The generalization testing method of F_t is given by executing the policy $\hat{\pi}$ in \mathcal{M} , but in each time step applying the modification F_t where the history H_t is given by the transitions executed by $\hat{\pi}$ so far. When both a generalization method and a generalization testing method are used concurrently, we will use subscripts to denote the generalization

method and superscripts to denote the testing method. For instance, \mathcal{G}_S^π corresponds to training with a state transforming method, and testing with a policy transforming method.

4. Roots of Overestimation in Deep Reinforcement Learning

Many reinforcement learning algorithms compute estimates for the state-action values in an MDP. Because these estimates are usually based on a stochastic interaction with the MDP, computing accurate estimates that correctly generalize to further interactions is one of the most fundamental tasks in reinforcement learning. A major challenge in this area has been the tendency of many classes of reinforcement learning algorithms to consistently overestimate state-action values. Initially the overestimation bias for Q-learning is discussed and theoretically justified by Thrun and Schwartz (1993) as a byproduct of using function approximators for state-action value estimates. In particular, Thrun and Schwartz (1993) proves that if the reinforcement learning policy overestimates the state-action values by γc during learning then the Q-learning algorithm will fail to learn optimal policy if $\gamma > \frac{1}{1+c}$.

Following this initial discussion it has been shown that several parts of the deep reinforcement learning process can cause overestimation bias. Learning overestimated state-action values can be caused by statistical bias of utilizing a single max operator (van Hasselt, 2010), coupling between value function and the optimal policy (Raileanu and Fergus, 2021; Cobbe et al., 2021), or caused by the accumulated function approximation error (Boyan and Moore, 1994).

Several methods have been proposed to target overestimation bias for value iteration algorithms. In particular, van Hasselt (2010) demonstrated that the expectation of a maximum of a random variable is not equal to maximum of the expectation of a random variable.

$$\mathbb{E}[\max_i X_i] \neq \max_i [\mathbb{E}[X_i]] \text{ where } X = \{X_1, X_2, \dots, X_N\}$$

This clear distinction shows that simple Q-learning is a biased estimator, and to solve this overestimation bias introduced by the max operator van Hasselt (2010) proposed to utilize a double estimator for the state-action value estimates. In particular, the double estimator for double Q-learning works as follows

$$Q^I(s, a) \leftarrow Q^I(s, a) + \alpha(s, a)(r(s, a) + \gamma Q^{II}(s', \max_a Q^I(s', a)) - Q^I(s', a))$$

and

$$Q^{II}(s, a) \leftarrow Q^{II}(s, a) + \alpha(s, a)(r(s, a) + \gamma Q^I(s', \max_a Q^{II}(s', a)) - Q^{II}(s', a)).$$

Later, the authors also created a version of this algorithm that can solve high dimensional state space problems (Hasselt et al., 2016). Some of the work on this line of research targeting overestimation bias for value iteration algorithms is based on simply averaging the state-action values with previously learned state-action value estimates during training time (Anschel et al., 2017). While overestimation bias was demonstrated to be a problem and discussed over a long period of time (Thrun and Schwartz, 1993; van Hasselt, 2010), recent studies also further demonstrated that actor critic algorithms also suffer from this issue (Fujimoto et al., 2018).

Table 1: Environment and algorithm details for different exploration strategies for generalization.

Citation	Method	Learning Environment	Algorithm
Mnih et al. (2015)	ϵ -greedy	Arcade Learning Environment	DQN
Bellemare et al. (2016)	Count-based	Arcade Learning Environment	A3C and DQN
Osband et al. (2016b)	RLSVI	Tetris	Tabular Q
Osband et al. (2016a)	Bootstrapped DQN	Arcade Learning Environment	DQN
Houthoofd et al. (2016)	VIME	DeepMind Control Suite	TRPO
Fortunato et al. (2018)	NoisyNet	Arcade Learning Environment	A3C and DQN
Lee et al. (2021)	SUNRISE	DCS ¹ & ALE	SAC & RDQN
Mahankali et al. (2024)	Random Latent	Arcade Learning Environment	PPO

5. The Role of Exploration in Overfitting

The fundamental trade-off of exploration vs exploitation is the dilemma that the agent can try to take actions to move towards more unexplored states by sacrificing the current immediate rewards. While there is a significant body of studies on provably efficient exploration strategies the results from these studies do not necessarily directly transfer to the high dimensional state or action MDPs. The most prominent indication of this is that, even though it is possible to use deep neural networks as function approximators for large state spaces, the agent will simply not be able to explore the full state space. The fact that the agent is able to only explore a portion of the state space simply creates a bias in the learnt value function (III, 1995).

In this section, we will go through several exploration strategies in deep reinforcement learning and how they affect policy overfitting. A quite simple version of this is based on adding noise in action selection during training e.g. ϵ -greedy exploration. Note that this is an example of a policy transforming generalization method \mathcal{G}_π in Definition 3.7 in Section 3. While ϵ -greedy exploration is widely used in deep reinforcement learning (Wang et al., 2016; Hamrick et al., 2020; Kapturowski et al., 2023), it has also been proven that to explore the state space these algorithms may take exponentially long (Kakade, 2003). Several others focused on randomizing different components of the reinforcement learning training algorithms. In particular, Osband et al. (2016b) proposes the randomized least squared value iteration algorithm to explore more efficiently in order to increase generalization in reinforcement learning for linearly parametrized value functions. This is achieved by simply adding Gaussian noise as a function of state visitation frequencies to the training dataset. Later, the authors also propose the bootstrapped DQN algorithm (i.e. adding temporally correlated noise) to increase generalization with non-linear function approximation Osband et al. (2016a). Recently, Mahankali et al. (2024) proposed to randomize the reward function to enhance exploration in high dimensional observation MDPs where policy gradient algorithms are used to explore. This study is also a clear example of the generalization through rewards as has been explained in Definition 3.4 in Section 3.

1. DeepMind Control Suite

Houthoofd et al. (2016) proposed an exploration technique centered around maximizing the information gain on the agent’s belief of the environment dynamics. In practice, the authors use Bayesian neural networks for effectively exploring high dimensional action space MDPs. Following this line of work on increasing efficiency during exploration Fortunato et al. (2018) proposes to add parametric noise to the deep reinforcement learning policy weights in high dimensional state MDPs. While several methods focused on ensemble state-action value function learning (Osband et al., 2016a), Lee et al. (2021) proposed reweighting target Q-values from an ensemble of policies (i.e. weighted Bellman backups) combined with highest upper-confidence bound action selection. Another line of research in exploration strategies focused on *count-based methods* that use the direct count of state visitations. In this line of work, Bellemare et al. (2016) tried to lay out the relationship between count based methods and intrinsic motivation, and used count-based methods for high dimensional state MDPs (i.e. Arcade Learning Environment). Yet it is worthwhile to note that most of the current deep reinforcement learning algorithms use very simple exploration techniques such as ϵ -greedy which is based on taking the action maximizing the state-action value function with probability $1 - \epsilon$ and taking a random action with probability ϵ (Mnih et al., 2015; Hasselt et al., 2016; Wang et al., 2016; Hamrick et al., 2020; Kapturowski et al., 2023).

It is possible to argue that the fact that the deep reinforcement learning policy obtained a higher score with the same number of samples by a particular type of training method \mathcal{A} compared to method \mathcal{B} is by itself evidence that the technique \mathcal{A} leads to more generalized policies. Even though the agent is trained and tested in the same environment, the explored states during training time are not exactly the same states visited during test time. The fact that the policy trained with technique \mathcal{A} obtains a higher score at the end of an episode is sole evidence that the agent trained with \mathcal{A} was able to visit further states in the MDP and thus succeed in them. Yet, throughout the paper we will discuss different notions of generalization investigated in different subfields of reinforcement learning research. While exploration vs exploitation stands out as one of the main problems in reinforcement learning policy performance most of the work conducted in this section focuses on achieving higher score in hard-exploration games (i.e. Montezuma’s Revenge) rather than aiming for a generally higher score for each game overall across a given benchmark. Thus, it is possible that the majority of work focusing on exploration so far might not be able to obtain policies that perform as well as those in the studies described in Section 6 across a given benchmark.

6. Regularization

In this section we will focus on different regularization techniques employed to increase generalization in deep reinforcement learning policies. We will go through these works by categorizing each of them under data augmentation, adversarial training, and direct function regularization. Under each category we will connect these different lines of approach to increase generalization in deep reinforcement learning to the settings we defined in Section 3.

Table 2: Environment and algorithm details for data augmentation techniques for state observation generalization. All of the studies in this section focus on state transformation methods \mathcal{G}_S defined in Section 3.

Citation	Method	Environment	Algorithm
Yarats et al. (2021)	DrQ	DCS, Arcade Learning Environment	DQN
Laskin et al. (2020b)	CuRL	DCS, Arcade Learning Environment	SAC and DQN
Laskin et al. (2020a)	RAD	DeepMind Control Suite, ProcGen	SAC and PPO
Korkmaz (2023)	Semantic Changes	Arcade Learning Environment	DDQN & A3C
Wang et al. (2020)	Mixreg	ProcGen	DQN and PPO

6.1 Data Augmentation

Several studies focus on diversifying the observations of the deep reinforcement learning policy to increase generalization capabilities. A line of research in this regard focused on simply employing versions of data augmentation techniques (Laskin et al., 2020a,b; Yarats et al., 2021) for high dimensional state representation environments. In particular, these studies involve simple techniques such as cropping, rotating or shifting the state observations during training time. While this line of work got considerable attention, a quite recent study Agarwal et al. (2021b) demonstrated that when the number of random seeds is increased to one hundred the relative performance achieved and reported in the original papers of (Laskin et al., 2020b; Yarats et al., 2021) on data augmentation training in deep reinforcement learning decreases to a level that might be significant to mention.

While some of the work on this line of research simply focuses on using a set of data augmentation methods (Laskin et al., 2020a,b; Yarats et al., 2021), other work focuses on proposing new environments to train in (Cobbe et al., 2020). The studies on designing new environments to train deep reinforcement learning policies basically aim to provide high variation in the observed environment such as changing background colors and changing object shapes in ways that are meaningful in the game, in order to increase test time generalization. In the line of robustness and test time performance, a more recent work that is also mentioned in Section 6.3 demonstrated that imperceptible semantically meaningful data augmentations can cause significant damage on the policy performance and certified robust deep reinforcement learning policies are more vulnerable to these imperceptible augmentations (Korkmaz, 2021a, 2023).

Within this category some work focuses on producing more observations by simply blending in (e.g. creating a mixture state from multiple different observations) several observations to increase generalization (Wang et al., 2020). While most of the studies trying to increase generalization by data augmentation techniques are primarily conducted in the DeepMind Control Suite or the Arcade Learning Environment (ALE) (Bellemare et al., 2013), some small fraction of these studies (Wang et al., 2020) are conducted in relatively recently designed training environments like ProcGen (Cobbe et al., 2020). In the line of research proposing learning environments Dennis et al. (2020) proposed unsupervised environment design by changing the environment parameters to asses generalization for maze structured environments by minimax training where the "adversary" creating an environ-

Table 3: Environment and algorithm details for different direct function regularization strategies for trying to overcome overfitting problems in reinforcement learning. Note that most of the methods based on direct function regularization are a form of algorithmic generalization \mathcal{G}_A to overcome overfitting as described in Section 3.

Citation	Proposed Method	Learning Environment
Igl et al. (2019)	SNI and IBAC	GridWorld and CoinRun
Vieillard et al. (2020b)	Munchausen RL	Arcade Learning Environment
Lee et al. (2020)	Network Randomization	2D CoinRun and 3D DeepMind Lab
Amit et al. (2020)	Discount Regularization	GridWorld and MuJoCo ²
Agarwal et al. (2021a)	PSM	DDMC and Rectangle Game ³
Liu et al. (2021)	BN and dropout and L_2/L_1	MuJoCo

ment for the policy to solve a task with goal and obstacles as an underspecified parameter. Cobbe et al. (2019) focuses on decoupling the training and testing set for reinforcement learning via simply proposing a new game environment CoinRun.

6.2 Direct Function Regularization

While some of the work we have discussed so far focuses on regularizing the data (i.e. state observations) as in Section 6.1, some focuses on directly regularizing the function learned with the intention of simulating techniques from deep neural network regularization like batch normalization and dropout (Igl et al., 2019). While some studies have attempted to simulate these known techniques in reinforcement learning, some focus on directly applying them to overcome overfitting. In this line of research, Liu et al. (2021) proposes to use known techniques from deep neural network regularization to apply in continuous control deep reinforcement learning training. In particular, these techniques are batch normalization (BN) (Ioffe and Szegedy, 2015), weight clipping, dropout, entropy and L_2/L_1 weight regularization. All these methods fall under the algorithmic generalization category \mathcal{G}_A as described in Section 3.

Lee et al. (2020) proposes to utilize a random network to essentially achieve a version of randomization in the input observations to increase generalization skills of deep reinforcement learning policies, and tests the proposal in the 2D CoinRun game proposed by Cobbe et al. (2019) and 3D DeepMind Lab (?). In particular, the authors essentially introduce a random convolutional layer to achieve this objective. This study is an example of an algorithmic generalization method \mathcal{G}_A described in Definition 3.3 when the single layer random network is not placed at the first layer of the deep neural network. However, when this single layer random network is placed at the first layer of the neural network, this method is essentially just introducing some noise to the state observations of the policy, thus this is an example of state transforming generalization. When this single random layer is placed

2. Low dimensional setting of MuJoCo is used for this study (Todorov et al., 2012).

3. Rectangle game is a simple video game with only two actions, "Right" and "Jump". The game has black background and two rectangles where the goal of the game is to avoid white obstacles and reach to the right side of the screen. Agarwal et al. (2021a) is the only paper we encountered experimenting with this particular game.

Table 4: Algorithm details for different direct function regularization strategies for trying to overcome overfitting problems in reinforcement learning. Note that most of the methods based on direct function regularization are a form of algorithmic generalization \mathcal{G}_A to overcome overfitting as described in Section 3.

Citation	Proposed Method	Reinforcement Learning Algorithm
Igl et al. (2019)	SNI and IBAC	Proximal Policy Optimization (PPO)
Vieillard et al. (2020b)	Munchausen RL	DQN and IQN
Lee et al. (2020)	Network Randomization	Proximal Policy Optimization (PPO)
Amit et al. (2020)	Discount Regularization	Twin Delayed DDPG (TD3)
Agarwal et al. (2021a)	PSM	Data Regularized-Q (DrQ)
Liu et al. (2021)	BN and dropout and L_2/L_1	PPO, TRPO, SAC, A2C

other than first, the method is no longer a state transforming generalization method because the states are not modified before they have been observed by the algorithm, but rather implicitly changed due to a random convolutional layer added in the architecture. We will further provide clear instances of the state transformation generalization also in Section 6.3 when the worst-case perturbation methods to target generalization in reinforcement learning policies are explained.

Some work employs contrastive representation learning to learn deep reinforcement learning policies from state observations that are close to each other (Agarwal et al., 2021a). This study leverage the temporal aspect of reinforcement learning and propose a policy similarity metric. The main goal of the paper is to lay out the sequential structure and utilize representation learning to learn generalizable abstractions from state representations. One drawback of this study is that most of the experimental study is conducted in a non-baseline environment (i.e. Rectangle game and Distracting DM Control Suite). Malik et al. (2021) studies query complexity of reinforcement learning policies that can generalize to multiple environments. The authors of this study focus on an example of the transition probability transformation setting \mathcal{G}_P in Definition 3.6, and the reward function transformation setting \mathcal{G}_R in Definition 3.4.

Another line of study in direct function generalization investigates the relationship between reduced discount factor and adding an ℓ_2 -regularization term to the loss function, i.e. weight decay (Amit et al., 2020). The authors in this work demonstrate the explicit connection between reducing the discount factor and adding an ℓ_2 -regularizer to the value function for temporal difference learning. In particular, this study demonstrates that adding an ℓ_2 -regularization term to the loss function is equal to training with a lower discount term, which the authors refer to as *discount regularization*. The results of this study however are based on experiments from tabular reinforcement learning, and the low dimensional setting of the MuJoCo environment (Todorov et al., 2012). This study is also another clear example of algorithmic generalization \mathcal{G}_A as described in Definition 3.3.

On the reward transformation for generalization setting \mathcal{G}_R defined in Definition 3.4, Vieillard et al. (2020b) adds the scaled log policy to the current rewards. To overcome overfitting some work tries to learn explicit or implicit similarity between the states to obtain a reasonable policy (Lan et al., 2021). In particular, the authors in this work try to

Table 5: Environment and algorithm details for adversarial policy regularization and attack techniques in deep reinforcement learning. Note that most of the methods based on adversarial approaches are a form of generalization assessment through state observations \mathcal{G}^S as described in Definition 3.8, and some falls under the generalization through environment dynamics \mathcal{G}^P as described in Definition 3.6.

Citation	Method	Environment	Algorithm
Huang et al. (2017)	Fast Gradient Sign (FGSM)	ALE	DQN, TRPO, A3C
Kos and Song (2017)	Fast Gradient Sign (FGSM)	ALE	DQN & IQN
Korkmaz (2022)	Adversarial Framework	ALE	DDQN & A3C
Lin et al. (2017)	Timing	ALE	A3C & DQN
Pinto et al. (2017)	Zero-sum game	MuJoCo	RARL
Gleave et al. (2020)	Adversarial Policies	MuJoCo	PPO
Korkmaz (2023)	Natural Attacks	ALE	DDQN & A3C
Huan et al. (2020)	State Adversarial-DQN	ALE and L_M^4	DDQN & PPO
Korkmaz (2024c)	Diagnostic Adversarial Volatility	ALE	DDQN

unify the state space representations by providing a taxonomy of metrics in reinforcement learning. Several studies proposed different ways to include Kullback-Leibler divergence between the current policy and the pre-updated policy to add as a regularization term in the reinforcement learning objective (Schulman et al., 2015). Recently, some studies argued that utilizing Kullback-Leibler regularization implicitly averages the state-action value estimates (Vieillard et al., 2020a).

6.3 The Adversarial Perspective for Deep Neural Policy Generalization

One of the ways to regularize the state observations is based on considering worst-case perturbations added to state observations (i.e. adversarial perturbations). This line of work starts with introducing perturbations produced by the fast gradient sign method proposed by Goodfellow et al. (2015) into deep reinforcement learning observations at test time (Huang et al., 2017; Kos and Song, 2017), and compares the generalization capabilities of the trained deep reinforcement learning policies in the presence worst-case perturbations and Gaussian noise. These gradient based adversarial methods are based on taking the gradient of the cost function used to train the policy with respect to the state observation.

$$s_{\text{adv}} = s + \epsilon \cdot \frac{\nabla_x J(s, Q(s, a))}{\|\nabla_s J(s, Q(s, a))\|_p},$$

Several other techniques have been proposed on the optimization line of the adversarial alteration of state observations. In this line of work, Korkmaz (2020) suggested a Nesterov momentum-based method to produce adversarial perturbations for deep reinforcement

4. Low dimensional state MuJoCo refers to the setting of MuJoCo where the state dimensions are not represented by pixels and dimensions of the state observations range from 11 to 117.

learning policies.

$$v_{t+1} = \mu \cdot v_t + \frac{\nabla_{s_{\text{adv}}} J(s_{\text{adv}}^t + \mu \cdot v_t, a)}{\|\nabla_{s_{\text{adv}}} J(s_{\text{adv}}^t + \mu \cdot v_t, a)\|_1}$$

$$s_{\text{adv}}^{t+1} = s_{\text{adv}}^t + \alpha \cdot \frac{v_{t+1}}{\|v_{t+1}\|_2}$$

Here $J(s_{\text{adv}}, a)$ is based on the cost function used to train the policy, s_{adv} represents the adversarial state observation, and μ is the momentum acceleration parameter. While a line of studies focused on optimization aspects of the adversarial perturbations, some studies demonstrated further the hidden linearity of deep reinforcement learning policies by revealing how these policies learn shared adversarial features across states, MDPs and across algorithms (Korkmaz, 2022).

In this work the authors investigate the root causes of this problem, and demonstrate that policy high-sensitivity directions and the perceptual similarity of the state observations are uncorrelated. Furthermore, the study demonstrates that the current state-of-the-art adversarial training techniques also learn similar high-sensitivity directions as the vanilla trained deep reinforcement learning policies.⁵ More recently, a line of work proposed theoretically founded algorithms to understand the temporal and spatial correlation of deep reinforcement learning decision making and what affects this decision making process (Korkmaz, 2024c). In particular, this study identifies what precisely affects and contributes to the decision making process of deep reinforcement learning policies from distributional shift to worst-case perturbations (i.e. adversarial), from algorithmic differences to architectural changes.

While several studies focused on improving computation techniques to optimize optimal perturbations, a line of research focused on making deep neural policies resilient to these perturbations. Pinto et al. (2017) proposed to model the dynamics between the adversary and the deep neural policy as a zero-sum game (Littman, 1994) where the goal of the adversary is to minimize expected cumulative rewards of the deep reinforcement learning policy.

$$R^{\text{agent}} = \mathbb{E}_{s_0 \sim \rho; a^{\text{agent}} \sim \pi^{\text{agent}}; a^{\text{adv}} \sim \pi^{\text{adv}}} \left[\sum_{t=0}^{T-1} r^{\text{agent}}(s, a^{\text{agent}}, a^{\text{adv}}) \right]$$

Here the adversarial policy is represented by π^{adv} , the policy of the agent represented by π^{agent} , and the rewards received by the agent represented by r^{agent} . The Nash equilibrium of the optimal rewards for this zero-sum game is

$$R^{\text{agent}^*} = \min_{\pi^{\text{adv}}} \max_{\pi^{\text{agent}}} R^{\text{agent}}(\pi^{\text{agent}}, \pi^{\text{adv}}) = \max_{\pi^{\text{agent}}} \min_{\pi^{\text{adv}}} R^{\text{agent}}(\pi^{\text{agent}}, \pi^{\text{adv}})$$

This study is a clear example of transition probability perturbation to achieve generalization $\mathcal{G}_{\mathcal{P}}$ in Definition 3.6 of Section 3. Gleave et al. (2020) approached this problem with an

5. From the security point of view, this adversarial framework is under the category of black-box adversarial attacks for which this is the first study that demonstrated that deep reinforcement learning policies are vulnerable to black-box adversarial attacks (Korkmaz, 2022). Furthermore, note that black-box adversarial perturbations are more generalizable global perturbations that can affect many different policies.

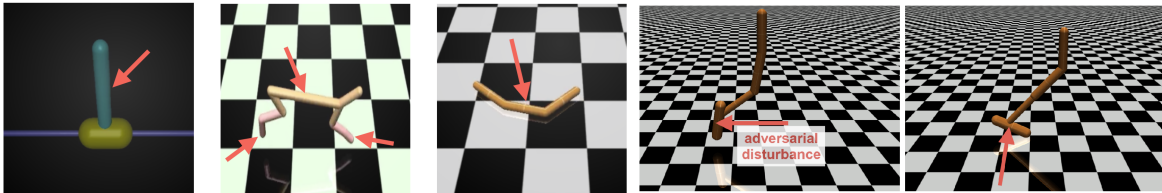


Figure 1: Robust adversarial reinforcement learning proposed in (Pinto et al., 2017). This paper proposes the zero-sum game to model the relationship between the agent and the adversary while focusing on introducing disturbances to the environment dynamics. Here the empirical studies are conducted in the MuJoCo environment.

adversary model which is restricted to take natural actions in the MDP instead of modifying the observations with ℓ_p -norm bounded perturbations. The authors model this dynamic as a zero-sum Markov game and solve it via self play Proximal Policy Optimization (PPO). Some recent studies, proposed to model the interaction between the adversary and the deep reinforcement learning policy as a state-adversarial MDP, and claimed that their proposed algorithm State Adversarial Double Deep Q-Network (SA-DDQN) learns theoretically certified robust policies against natural noise and perturbations. In particular, these certified adversarial training techniques aim to add a regularizer term to the temporal difference loss in deep Q -learning

$$\mathcal{H}(r_i + \gamma \max_a \hat{Q}_\theta(s_i, a; \theta) - Q_\theta(s_i, a_i; \theta)) + \kappa \mathcal{R}(\theta)$$

where \mathcal{H} is the Huber loss, \hat{Q} refers to the target network and κ is to adjust the level of regularization for convergence. The regularizer term can vary for different certified adversarial training techniques yet the baseline technique uses $\mathcal{R}(\theta)$

$$\max_{\hat{s} \in B(s)} \{ \max_{a \neq \arg \max_{a'} Q(s, a')} \max_{a'} Q_\theta(\hat{s}, a) - Q_\theta(\hat{s}, \arg \max_{a'} Q(s, a')), -c \}.$$

where $B(s)$ is an ℓ_p -norm ball of radius ϵ . While these certified adversarial training techniques drew some attention from the community, more recently manifold concerns have been raised on the robustness of theoretically certified adversarially trained deep reinforcement learning policies (Korkmaz, 2021c,b, 2022, 2024a). In these studies, the authors argue that adversarially trained (i.e. certified robust) deep reinforcement learning policies learn inaccurate state-action value functions and non-robust features from the environment. More importantly, recently it has been shown that certified robust deep reinforcement learning policies have worse generalization capabilities compared to vanilla trained reinforcement learning policies in high dimensional state space MDPs (Korkmaz, 2023). While this study provides a contradistinction between adversarial and natural directions that are intrinsic to the MDP, it further demonstrates that the certified adversarial training techniques block generalization capabilities of standard deep reinforcement learning policies. Furthermore note that this study is also a clear example of a state observation perturbation generalization testing method \mathcal{G}_S^S in Definition 3.8 in Section 3. For a more comprehensive view on generalization and robustness see Korkmaz (2024b).

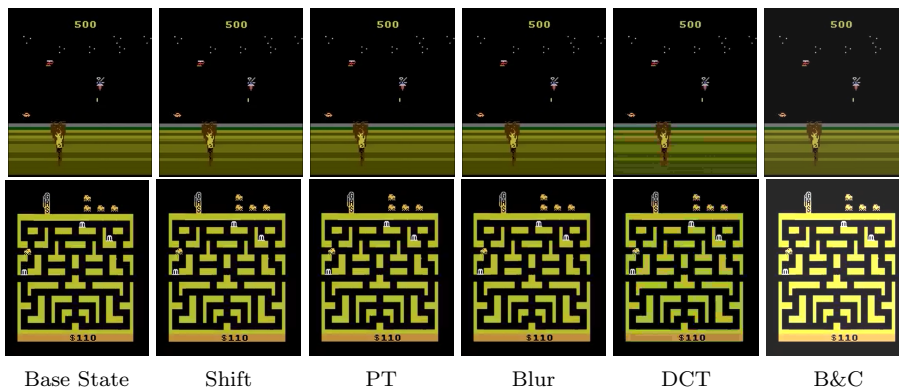


Figure 2: State transformation generalization under adversarial perspective in the Arcade Learning Environment (Korkmaz, 2023). Note that under the adversarial influence direction of research, the state transformation generalization is constrained by the imperceptibility of the transformations. Columns: base frame, shifting, perspective transformation, blurring, discrete cosine transform artifacts, brightness and contrast. Up: JamesBond. Down: BankHeist.

It is important to observe that the methods that focuses on improving generalization, i.e. robust training, described in this section rarely employ the different generalization testing methods proposed by other work. Thus, focusing narrowly on one aspect of generalization with one dimensional improvements in actuality decreases generalization on another aspect, as has been shown in the case of adversarial training (Korkmaz, 2023). Therefore we again emphasize the need to understand the significance of a concrete definition of generalization, and a unified baseline to precisely measure it.

7. Meta-Reinforcement Learning and Meta Gradients

A quite recent line of research directs its research efforts to discovering reinforcement learning algorithms automatically, without explicitly designing them, via meta-gradients (Oh et al., 2020; Xu et al., 2020). This line of study targets learning the "learning algorithm" by only interacting with a set of environments as a meta-learning problem. In particular,

$$\eta^* = \arg \max_{\eta} \mathbb{E}_{\varepsilon \sim \rho(\varepsilon)} \mathbb{E}_{\theta_0 \sim \rho(\theta_0)} [\mathbb{E}_{\theta_N} [\sum_{t=0}^{\infty} \gamma^t r_t]]$$

here the optimal update rule is parametrized by η , for a distribution on environments $\rho(\varepsilon)$ and initial policy parameters $\rho(\theta_0)$ where $\mathbb{E}_{\theta_N} [\sum_{t=0}^{\infty} \gamma^t r_t]$ is the expected return for the end of the lifetime of the agent. The objective of meta-reinforcement learning is to be able to build agents that can learn *how to learn* over time, thus allowing these policies to adapt to a changing environment or even any other changing conditions of the MDP.

Quite recently, a significant line of research has been conducted to achieve this objective, particularly Oh et al. (2020) proposes to discover update rules for reinforcement learning. This line of work also falls under the algorithmic generalization $\mathcal{G}_{\mathbb{A}}$ in Definition 3.3 defined in Section 3. Following this work Xu et al. (2020) proposed a joint meta-learning framework to learn what the policy should predict and how these predictions should be used in

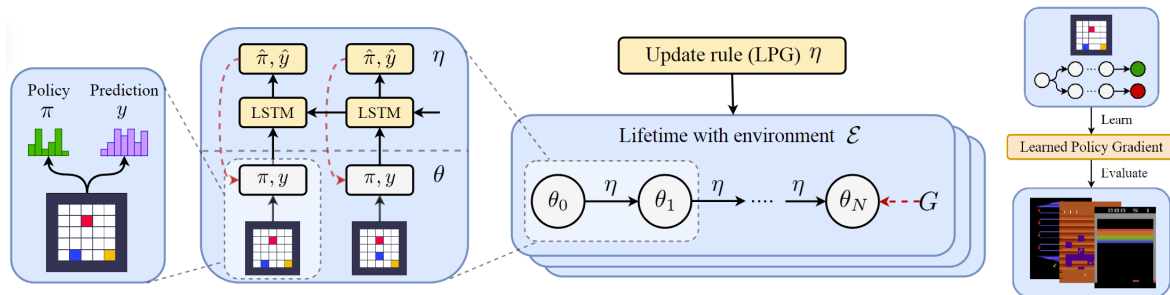


Figure 3: Meta training of the learned policy gradient that have been described in (Oh et al., 2020). Right: The learned policy gradient algorithm that has been trained in toy examples can generalize to more complex environment such as the Arcade Learning Environment.

updating the policy. Recently, Kirsch et al. (2022) proposes to use symmetry information in discovering reinforcement learning algorithms and discusses meta-generalization. There is also some work on enabling reinforcement learning algorithms to discover temporal abstractions (Veeriah et al., 2021). In particular, temporal abstraction refers to the ability of the policy to abstract a sequence of actions to achieve certain sub-tasks. As it is promised within this subfield, meta-reinforcement learning is considered to be a research direction that could enable us to build deep reinforcement learning policies that can generalize to different environments, to changing environments over time, or even to different tasks.

8. Transfer in Reinforcement Learning

Transfer in reinforcement learning is a subfield heavily discussed in certain applications of reinforcement learning algorithms, e.g. robotics. In current robotics research there is not a safe way of training a reinforcement learning agent by letting the robot explore in real life. Hence, the way to overcome this is to train policies in a simulated environment, and install the trained policies in the actual application setting. The fact that the simulation environment and the installation environment are not identical is one of the main problems for reinforcement learning application research. This is referred to as the *sim-to-real gap*.

Another subfield in reinforcement learning research focusing on obtaining generalizable policies investigates this concept through *transfer in reinforcement learning*. The consideration in this line of research is to build policies that are trained for a particular task with limited data and to try to make these policies perform well on slightly different tasks. An initial discussion on this starts with Taylor and Stone (2007) to obtain policies initially trained in a source task and transferred to a target task in a more sample efficient way. Later, Tirinzoni et al. (2018) proposes to transfer value functions that are based on learning a prior distribution over optimal value functions from a source task. However, this study is conducted in simple environments with low dimensional state spaces. Barreto et al. (2017) considers the reward transformation setting \mathcal{G}_R in Definition 3.4 from Section 3. In particular, the authors consider a policy transfer between a specific task with a reward function $r(s, a)$ and a different task with reward function $r'(s, a)$. The goal of the study is to decouple

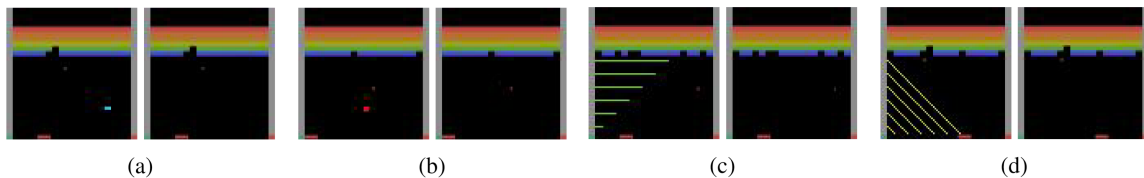


Figure 4: Transfer in reinforcement learning as has been described in (Gamrian and Goldberg, 2019) that falls under the generalization through observation category explained in Definition 3.5. The frames are taken from Breakout game in the Arcade Learning Environment. The left frames represent the target task and the right frames represents the source tasks generated via generative adversarial networks.

the state representations from the task. In the setting of state transformation for generalization \mathcal{G}_S in Definition 3.5 Gamrian and Goldberg (2019) focuses on state-wise differences between source and target task. In particular, the authors use unaligned generative adversarial networks to create target task states from source task states. In the setting of policy transformation for generalization \mathcal{G}_π in Definition 3.7 Jain et al. (2020) focuses on zero-shot generalization to a newly introduced action set to increase adaptability. While transfer learning is a promising research direction for reinforcement learning, the studies in this subfield still remain oriented only towards reinforcement learning applications, and thus the main focus on applications centered on this subfield provides a non-unified progress in research due to the lack of an established baseline in which the proposed claims and algorithms can be consistently compared.

9. Lifelong Reinforcement Learning

Lifelong learning is a subfield closely related to transfer learning that has recently drawn attention from the reinforcement learning community. Lifelong learning aims to build policies that can sequentially solve different tasks by being able to transfer knowledge between tasks. On this line of research, Lecarpentier et al. (2021) provide an algorithm for value-based transfer in the Lipschitz continuous task space with theoretical contributions for lifelong learning goals. In the setting of action transformation for generalization \mathcal{G}_π in Definition 3.7 Chandak et al. (2020) focuses on temporally varying (e.g. variations between source task and target task) the action set in lifelong learning. In lifelong reinforcement learning some studies focus on different exploration strategies. In particular, Garcia and Thomas (2019) models the exploration strategy problem for lifelong learning as another MDP, and the study uses a separate reinforcement learning agent to find an optimal exploration method for the initial lifelong learning agent. The lack of benchmarks limits the progress of lifelong reinforcement learning research by restricting the direct comparison between proposed algorithms or methods. However, quite recent work proposed a new training environment benchmark based on robotics applications for lifelong learning to overcome this issue (Wolczyk et al., 2021)⁶.

6. The state dimension for this benchmark is 12. Hence, the state space is low dimensional.

10. Inverse Reinforcement Learning

Inverse reinforcement learning focuses on learning a functioning policy in the absence of a reward function. Since the real reward function is inaccessible in this setting and the reward function needs to be learnt from observing an expert completing the given task, the inverse reinforcement learning setting falls under the reward transformation for generalization setting \mathcal{G}_R defined in Definition 3.4 in Section 3. The initial work that introduced inverse reinforcement learning was proposed by Ng and Russell (2000) demonstrating that multiple different reward functions can be constructed for an observed optimal policy. The authors of this initial study achieve this objective via linear programming,

$$\begin{aligned} \max \sum_{s \in S_\rho} \min_{a \in A} \{ & p(\mathbb{E}_{s' \sim \mathcal{P}(s, a_1 | \cdot)} \mathcal{V}^\pi(s') - \mathbb{E}_{s' \sim \mathcal{P}(s, a | \cdot)} \mathcal{V}^\pi(s')) \} \\ \text{s.t. } & |\alpha_i| \leq 1, i = 1, 2, \dots, d \end{aligned}$$

where $p(x) = x$ if $x \geq 0$, $p(x) = 2x$ otherwise and $\mathcal{V}^\pi = \alpha_1 \mathcal{V}_1^\pi + \alpha_2 \mathcal{V}_2^\pi + \dots + \alpha_d \mathcal{V}_d^\pi$. In this line of work, there has been recent progress that achieved learning functioning policies in high-dimensional state observation MDPs (Garg et al., 2021). The study achieves this by learning a soft Q -function from observing expert demonstrations, and the study further argues that it is possible to recover rewards from the learnt soft state-action value function.

11. Conclusion

In this paper we tried to answer the following questions: (i) *What are the explicit problems limiting reinforcement learning algorithms from obtaining high-performing policies that can generalize to complex environments?* (ii) *How can we unify and categorize the concept of generalization in deep reinforcement learning considering many subfields under reinforcement learning at their core focus on the same objective?* (iii) *What are the similarities and differences of these different techniques proposed by different subfields of reinforcement learning research to build reinforcement learning policies that can robustly generalize?* To answer these questions first we introduce a theoretical analysis and mathematical framework to unify and categorize the concept of generalization in deep reinforcement learning. Then we explain the connection and the significance of exploration in overfitting to a learning environment, and explain the manifold causes of overestimation bias in reinforcement learning. Starting from all the different regularization techniques in either state representations or in learnt value functions from worst-case to average-case, we provide a current layout of the wide range of reinforcement learning subfields that are essentially working towards the same objective, i.e. generalizable deep reinforcement learning policies. Finally, we provided a discussion for each category on the drawbacks and advantages of these algorithms. We believe our study can provide a compact unifying formalization on recent reinforcement learning generalization research. We believe our theoretical framework can guide current and future research to build deep reinforcement learning agents that can robustly generalize to complex environments.

References

- Agarwal, R., Machado, M. C., Castro, P. S., and Bellemare, M. G. (2021a). Contrastive behavioral similarity embeddings for generalization in reinforcement learning. In *International Conference on Learning Representations (ICLR)*.
- Agarwal, R., Schwarzer, M., Castro, P. S., Courville, A. C., and Bellemare, M. G. (2021b). Deep reinforcement learning at the edge of the statistical precipice. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W., editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 29304–29320.
- Amit, R., Meir, R., and Ciosek, K. (2020). Discount factor as a regularizer in reinforcement learning. In *International Conference on Machine Learning (ICML)*.
- Anschel, O., Baram, N., and Shimkin, N. (2017). Averaged-dqn: Variance reduction and stabilization for deep reinforcement learning. In *International Conference on Machine Learning (ICML)*.
- Barreto, A., Dabney, W., Munos, R., Hunt, J. J., Schaul, T., Silver, D., and van Hasselt, H. (2017). Successor features for transfer in reinforcement learning. In Guyon, I., von Luxburg, U., Bengio, S., Wallach, H. M., Fergus, R., Vishwanathan, S. V. N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4055–4065.
- Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. (2013). The arcade learning environment: An evaluation platform for general agents. *Journal Artificial Intelligence Research (JAIR)*, 47:253–279.
- Bellemare, M. G., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. (2016). Unifying count-based exploration and intrinsic motivation. In Lee, D. D., Sugiyama, M., von Luxburg, U., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 1471–1479.
- Bellman, R. (1957). *Dynamic programming*. Princeton University Press, Princeton.
- Bellman, R. and Dreyfus, S. (1959). *Functional approximation and dynamic programming*. Mathematical Tables and Other Aids to Computation.
- Boyan, J. A. and Moore, A. W. (1994). Generalization in reinforcement learning: Safely approximating the value function. In Tesauro, G., Touretzky, D. S., and Leen, T. K., editors, *Advances in Neural Information Processing Systems 7, [NIPS Conference, Denver, Colorado, USA, 1994]*, pages 369–376. MIT Press.
- Chandak, Y., Theodorou, G., Nota, C., and Thomas, P. S. (2020). Lifelong learning with a changing action set. In *AAAI Conference on Artificial Intelligence, AAAI*.

- Cobbe, K., Hesse, C., Hilton, J., and Schulman, J. (2020). Leveraging procedural generation to benchmark reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 2048–2056. PMLR.
- Cobbe, K., Hilton, J., Klimov, O., and Schulman, J. (2021). Phasic policy gradient. In Meila, M. and Zhang, T., editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 2020–2027. PMLR.
- Cobbe, K., Klimov, O., Hesse, C., Kim, T., and Schulman, J. (2019). Quantifying generalization in reinforcement learning. In *International Conference on Machine Learning (ICML)*.
- Dennis, M., Jaques, N., Vinitzky, E., Bayen, A. M., Russell, S., Critch, A., and Levine, S. (2020). Emergent complexity and zero-shot transfer via unsupervised environment design. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Fawzi, A., Balog, M., Huang, A., Hubert, T., Romera-Paredes, B., Barekattain, M., Novikov, A., Ruiz, F. J. R., Schrittwieser, J., Swirszcz, G., Silver, D., Hassabis, D., and Kohli, P. (2022). Discovering faster matrix multiplication algorithms with reinforcement learning. *Nature*, 610(7930):47–53.
- Fortunato, M., Azar, M. G., Piot, B., Menick, J., Osband, I., Graves, A., Mnih, V., Munos, R., Hassabis, D., Pietquin, O., Blundell, C., and Legg, S. (2018). Noisy networks for exploration. *International Conference on Learning Representations (ICLR)*.
- Fujimoto, S., van Hoof, H., and Meger, D. (2018). Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning (ICML)*.
- Gamrian, S. and Goldberg, Y. (2019). Transfer learning for related reinforcement learning tasks via image-to-image translation. In *International Conference on Machine Learning (ICML)*.
- Garcia, F. M. and Thomas, P. S. (2019). A meta-mdp approach to exploration for lifelong reinforcement learning. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R., editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 5692–5701.
- Garg, D., Chakraborty, S., Cundy, C., Song, J., and Ermon, S. (2021). Iq-learn: Inverse soft-q learning for imitation. *Neural Information Processing Systems (NeurIPS) [Spotlight Presentation]*.
- Gleave, A., Dennis, M., Wild, C., Neel, K., Levine, S., and Russell, S. (2020). Adversarial policies: Attacking deep reinforcement learning. *International Conference on Learning Representations (ICLR)*.

- Goodfellow, I., Shlens, J., and Szegedy, C. (2015). Explaining and harnessing adversarial examples. *International Conference on Learning Representations (ICLR)*.
- Google Gemini (2023). Gemini: A family of highly capable multimodal models. *Technical Report*, <https://arxiv.org/abs/2312.11805>.
- Hamrick, J., Bapst, V., SanchezGonzalez, A., Pfaff, T., Weber, T., Buesing, L., and Battaglia, P. (2020). Combining q-learning and search with amortized value estimates. In *8th International Conference on Learning Representations, ICLR*.
- Hasselt, H. v., Guez, A., and Silver, D. (2016). Deep reinforcement learning with double q-learning. *AAAI Conference on Artificial Intelligence, AAAI*.
- Houthoofd, R., Chen, X., Duan, Y., Schulman, J., Turck, F. D., and Abbeel, P. (2016). VIME: variational information maximizing exploration. In Lee, D. D., Sugiyama, M., von Luxburg, U., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 1109–1117.
- Huan, Z., Hongge, C., Chaowei, X., Li, B., Boning, M., Liu, D., and Hsiesh, C. (2020). Robust deep reinforcement learning against adversarial perturbations on state observatons. *Conference on Neural Information Processing Systems (NeurIPS)*.
- Huang, S., Papernot, N., Goodfellow, Ian an Duan, Y., and Abbeel, P. (2017). Adversarial attacks on neural network policies. *International Conference on Learning Representations (ICLR)*.
- Igl, M., Ciosek, K., Li, Y., Tschitschek, S., Zhang, C., Devlin, S., and Hofmann, K. (2019). Generalization in reinforcement learning with selective noise injection and information bottleneck. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Garnett, R., editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 13956–13968.
- III, L. C. B. (1995). Residual algorithms: Reinforcement learning with function approximation. In Prieditis, A. and Russell, S., editors, *Machine Learning, Proceedings of the Twelfth International Conference on Machine Learning, Tahoe City, California, USA, July 9-12, 1995*, pages 30–37. Morgan Kaufmann.
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Bach, F. R. and Blei, D. M., editors, *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 448–456. JMLR.org.
- Jain, A., Szot, A., and Lim, J. J. (2020). Generalization to new actions in reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 4661–4672. PMLR.

- Kakade, S. (2003). On the sample complexity of reinforcement learning. In *PhD Thesis*.
- Kapturowski, S., Campos, V., Jiang, R., Rakicevic, N., van Hasselt, H., Blundell, C., and Badia, A. P. (2023). Human-level atari 200x faster. In *The Eleventh International Conference on Learning Representations, ICLR 2023*.
- Kirsch, L., Flennerhag, S., van Hasselt, H., Friesen, A. L., Oh, J., and Chen, Y. (2022). Introducing symmetries to black box meta reinforcement learning. In *AAAI Conference on Artificial Intelligence, AAAI*.
- Korkmaz, E. (2020). Nesterov momentum adversarial perturbations in the deep reinforcement learning domain. *International Conference on Machine Learning (ICML) Workshop*.
- Korkmaz, E. (2021a). Adversarial training blocks generalization in neural policies. *International Conference on Learning Representation (ICLR) Robust and Reliable Machine Learning in the Real World Workshop*.
- Korkmaz, E. (2021b). Investigating vulnerabilities of deep neural policies. In de Campos, C. P., Maathuis, M. H., and Quaeghebeur, E., editors, *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence, UAI 2021, Virtual Event, 27-30 July 2021*, volume 161 of *Proceedings of Machine Learning Research*, pages 1661–1670. AUAI Press.
- Korkmaz, E. (2021c). Non-robust feature mapping in deep reinforcement learning. *International Conference on Machine Learning (ICML) Adversarial Machine Learning Workshop*.
- Korkmaz, E. (2022). Deep reinforcement learning policies learn shared adversarial features across mdps. *AAAI Conference on Artificial Intelligence, AAAI*.
- Korkmaz, E. (2023). Adversarial robust deep reinforcement learning requires redefining robustness. *AAAI Conference on Artificial Intelligence, AAAI*.
- Korkmaz, E. (2024a). Adversarial Robust Deep Reinforcement Learning is Neither Robust nor Safe. Conference on Neural Information Processing Systems (NeurIPS) Workshop on Statistical Foundations of LLMs and Foundation Models.
- Korkmaz, E. (2024b). Principled Analysis of Machine Learning Paradigms. PhD Thesis.
- Korkmaz, E. (2024c). Understanding and diagnosing deep reinforcement learning. In *Proceedings of the 41st International Conference on Machine Learning ICML*, Proceedings of Machine Learning Research (PMLR). PMLR.
- Korkmaz, E. and Brown-Cohen, J. (2023). Detecting adversarial directions in deep reinforcement learning to make robust decisions. In *International Conference on Machine Learning, ICML 2023*, volume 202 of *Proceedings of Machine Learning Research*, pages 17534–17543. PMLR.

- Kos, J. and Song, D. (2017). Delving into adversarial attacks on deep policies. *International Conference on Learning Representations (ICLR)*.
- Lan, C. L., Bellemare, M. G., and Castro, P. S. (2021). Metrics and continuity in reinforcement learning. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Virtual Event, February 2-9, 2021*, pages 8261–8269. AAAI Press.
- Laskin, M., Lee, K., Stooke, A., Pinto, L., Abbeel, P., and Srinivas, A. (2020a). Reinforcement learning with augmented data. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Laskin, M., Srinivas, A., and Abbeel, P. (2020b). CURL: contrastive unsupervised representations for reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 5639–5650. PMLR.
- Lecarpentier, E., Abel, D., Asadi, K., Jinnai, Y., Rachelson, E., and Littman, M. L. (2021). Lipschitz lifelong reinforcement learning. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 8270–8278. AAAI Press.
- Lee, K., Laskin, M., Srinivas, A., and Abbeel, P. (2021). SUNRISE: A simple unified framework for ensemble learning in deep reinforcement learning. In Meila, M. and Zhang, T., editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 6131–6141. PMLR.
- Lee, K., Lee, K., Shin, J., and Lee, H. (2020). Network randomization: A simple technique for generalization in deep reinforcement learning. In *International Conference on Learning Representations (ICLR)*.
- Lin, L.-J. (1993). Reinforcement learning for robots using neural networks. Technical report.
- Lin, Y.-C., Zhang-Wei, H., Liao, Y.-H., Shih, M.-L., Liu, i.-Y., and Sun, M. (2017). Tactics of adversarial attack on deep reinforcement learning agents. *IJCAI*.
- Littman, M. L. (1994). Markov games as a framework for multi-agent reinforcement learning. In Cohen, W. W. and Hirsh, H., editors, *Machine Learning, Proceedings of the Eleventh International Conference, Rutgers University, New Brunswick, NJ, USA, July 10-13, 1994*, pages 157–163. Morgan Kaufmann.
- Liu, Z., Li, X., and Darrell, T. (2021). Regularization matters in policy optimization - an empirical study on continuous control. In *International Conference on Learning Representations (ICLR)*.

- Mahankali, S., Hong, Z., Sekhari, A., Rakhlin, A., and Agrawal, P. (2024). Random latent exploration for deep reinforcement learning. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Malik, D., Li, Y., and Ravikumar, P. (2021). When is generalizable reinforcement learning tractable? In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W., editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 8032–8045.
- Mankowitz, D. J., Michi, A., Zhernov, A., Gelmi, M., Selvi, M., Paduraru, C., Leurent, E., Iqbal, S., Lespiau, J., Ahern, A., Köppe, T., Millikin, K., Gaffney, S., Elster, S., Broshear, J., Gamble, C., Milan, K., Tung, R., Hwang, M., Cemgil, T., Barekatain, M., Li, Y., Mandhane, A., Hubert, T., Schrittwieser, J., Hassabis, D., Kohli, P., Riedmiller, M. A., Vinyals, O., and Silver, D. (2023). Faster sorting algorithms discovered using deep reinforcement learning. *Nature*, 618(7964):257–263.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, a. G., Graves, A., Riedmiller, M., Fidjeland, A., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518:529–533.
- Ng, A. Y. and Russell, S. J. (2000). Algorithms for inverse reinforcement learning. In Langley, P., editor, *Proceedings of the Seventeenth International Conference on Machine Learning (ICML 2000), Stanford University, Stanford, CA, USA, June 29 - July 2, 2000*, pages 663–670.
- Oh, J., Hessel, M., Czarnecki, W. M., Xu, Z., van Hasselt, H., Singh, S., and Silver, D. (2020). Discovering reinforcement learning algorithms. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- OpenAI (2023). Gpt-4 technical report. *CoRR*.
- Osband, I., Blundell, C., Pritzel, A., and Roy, B. V. (2016a). Deep exploration via bootstrapped DQN. In Lee, D. D., Sugiyama, M., von Luxburg, U., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 4026–4034.
- Osband, I., Roy, B. V., and Wen, Z. (2016b). Generalization and exploration via randomized value functions. In *International Conference on Machine Learning (ICML)*.
- Pinto, L., Davidson, J., Sukthankar, R., and Gupta, A. (2017). Robust adversarial reinforcement learning. In Precup, D. and Teh, Y. W., editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 2817–2826. PMLR.

- Raileanu, R. and Fergus, R. (2021). Decoupling value and policy for generalization in reinforcement learning. In *International Conference on Machine Learning (ICML)*.
- Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T., Lillicrap, T. P., and Silver, D. (2020). Mastering atari, go, chess and shogi by planning with a learned model. *Nat.*, 588(7839):604–609.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M. I., and Moritz, P. (2015). Trust region policy optimization. In Bach, F. R. and Blei, D. M., editors, *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 1889–1897. JMLR.org.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., Chen, Y., Lillicrap, T. P., Hui, F., Sifre, L., van den Driessche, G., Graepel, T., and Hassabis, D. (2017). Mastering the game of go without human knowledge. *Nat.*, 550(7676):354–359.
- Sutton, R. (1984). Temporal credit assignment in reinforcement learning. PhD Thesis University of Massachusetts Amherst.
- Sutton, R. (1988). Learning to predict by the methods of temporal difference. *Machine Learning*.
- Sutton, R. S., McAllester, D. A., Singh, S., and Mansour, Y. (1999). Policy gradient methods for reinforcement learning with function approximation. In Solla, S. A., Leen, T. K., and Müller, K., editors, *Advances in Neural Information Processing Systems 12, [NIPS Conference, Denver, Colorado, USA, November 29 - December 4, 1999]*, pages 1057–1063. The MIT Press.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. (2014). Intriguing properties of neural networks. *International Conference on Learning Representations (ICLR)*.
- Taylor, M. E. and Stone, P. (2007). Cross-domain transfer for reinforcement learning. In Ghahramani, Z., editor, *Machine Learning, Proceedings of the Twenty-Fourth International Conference (ICML 2007), Corvallis, Oregon, USA, June 20-24, 2007*, volume 227 of *ACM International Conference Proceeding Series*, pages 879–886. ACM.
- Thrun, S. and Schwartz, A. (1993). Issues in using function approximation for reinforcement learning. In *Fourth Connectionist Models Summer School*.
- Tirinzoni, A., Rodríguez-Sánchez, R., and Restelli, M. (2018). Transfer of value functions via variational methods. In Bengio, S., Wallach, H. M., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 6182–6192.
- Todorov, E., Erez, T., and Tassa, Y. (2012). Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE.

- van Hasselt, H. (2010). Double q-learning. In Lafferty, J. D., Williams, C. K. I., Shawe-Taylor, J., Zemel, R. S., and Culotta, A., editors, *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada*, pages 2613–2621. Curran Associates, Inc.
- Veeriah, V., Zahavy, T., Hessel, M., Xu, Z., Oh, J., Kemaev, I., van Hasselt, H., Silver, D., and Singh, S. (2021). Discovery of options via meta-learned subgoals. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W., editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 29861–29873.
- Vieillard, N., Kozuno, T., Scherrer, B., Pietquin, O., Munos, R., and Geist, M. (2020a). Leverage the average: an analysis of KL regularization in reinforcement learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Vieillard, N., Pietquin, O., and Geist, M. (2020b). Munchausen reinforcement learning. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., Oh, J., Horgan, D., Kroiss, M., Danihelka, I., Huang, A., Sifre, L., Cai, T., Agapiou, J. P., Jaderberg, M., Vezhnevets, A. S., Leblond, R., Pohlen, T., Dalibard, V., Budden, D., Sulsky, Y., Molloy, J., Paine, T. L., Gülçehre, Ç., Wang, Z., Pfaff, T., Wu, Y., Ring, R., Yogatama, D., Wünsch, D., McKinney, K., Smith, O., Schaul, T., Lillicrap, T. P., Kavukcuoglu, K., Hassabis, D., Apps, C., and Silver, D. (2019). Grandmaster level in starcraft II using multi-agent reinforcement learning. *Nat.*, 575(7782):350–354.
- Wang, K., Kang, B., Shao, J., and Feng, J. (2020). Improving generalization in reinforcement learning with mixture regularization. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Wang, Z., Schaul, T., Hessel, M., van Hasselt, H., Lanctot, M., and de Freitas, N. (2016). Dueling network architectures for deep reinforcement learning. In Balcan, M. and Weinberger, K. Q., editors, *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 1995–2003. JMLR.org.
- Watkins, C. (1989). Learning from delayed rewards. In *PhD thesis, Cambridge*.

- Wolczyk, M., Zajac, M., Pascanu, R., Kucinski, L., and Milos, P. (2021). Continual world: A robotic benchmark for continual reinforcement learning. In Ranzato, M., Beygelzimer, A., Dauphin, Y. N., Liang, P., and Vaughan, J. W., editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 28496–28510.
- Xu, Z., van Hasselt, H. P., Hessel, M., Oh, J., Singh, S., and Silver, D. (2020). Meta-gradient reinforcement learning with an objective discovered online. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Yarats, D., Kostrikov, I., and Fergus, R. (2021). Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. In *International Conference on Learning Representations (ICLR)*.