

# A Brief Survey of Deep Reinforcement Learning

Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, Anil Anthony Bharath

**Abstract**—Deep reinforcement learning is poised to revolutionise the field of AI and represents a step towards building autonomous systems with a higher level understanding of the visual world. Currently, deep learning is enabling reinforcement learning to scale to problems that were previously intractable, such as learning to play video games directly from pixels. Deep reinforcement learning algorithms are also applied to robotics, allowing control policies for robots to be learned directly from camera inputs in the real world. In this survey, we begin with an introduction to the general field of reinforcement learning, then progress to the main streams of value-based and policy-based methods. Our survey will cover central algorithms in deep reinforcement learning, including the deep  $Q$ -network, trust region policy optimisation, and asynchronous advantage actor-critic. In parallel, we highlight the unique advantages of deep neural networks, focusing on visual understanding via reinforcement learning. To conclude, we describe several current areas of research within the field.

## I. INTRODUCTION

One of the primary goals of the field of artificial intelligence (AI) is to produce fully autonomous agents that interact with their environments to learn optimal behaviours, improving over time through trial and error. Crafting AI systems that are responsive and can effectively learn has been a long-standing challenge, ranging from robots, which can sense and react to the world around them, to purely software-based agents, which can interact with natural language and multimedia. A principled mathematical framework for experience-driven autonomous learning is reinforcement learning (RL) [135]. Although RL had some successes in the past [141, 129, 62, 93], previous approaches lacked scalability and were inherently limited to fairly low-dimensional problems. These limitations exist because RL algorithms share the same complexity issues as other algorithms: memory complexity, computational complexity, and in the case of machine learning algorithms, sample complexity [133]. What we have witnessed in recent years—the rise of deep learning, relying on the powerful function approximation and representation learning properties of deep neural networks—has provided us with new tools to overcoming these problems.

The advent of deep learning has had a significant impact on many areas in machine learning, dramatically improving the state-of-the-art in tasks such as object detection, speech recognition, and language translation [70]. The most important property of deep learning is that deep neural networks can automatically find compact low-dimensional representations (features) of high-dimensional data (e.g., images, text and audio). Through crafting inductive biases into neural network architectures, particularly that of hierarchical representations, machine learning practitioners have made effective progress in addressing the curse of dimensionality [15]. Deep learning has similarly accelerated progress in RL, with the use of deep learning algorithms within RL defining the field of “deep reinforcement learning” (DRL). The aim of this survey

is to cover both seminal and recent developments in DRL, conveying the innovative ways in which neural networks can be used to bring us closer towards developing autonomous agents. For a more comprehensive survey of recent efforts in DRL, including applications of DRL to areas such as natural language processing [106, 5], we refer readers to the overview by Li [78].

Deep learning enables RL to scale to decision-making problems that were previously intractable, i.e., settings with high-dimensional state and action spaces. Amongst recent work in the field of DRL, there have been two outstanding success stories. The first, kickstarting the revolution in DRL, was the development of an algorithm that could learn to play a range of Atari 2600 video games at a superhuman level, directly from image pixels [84]. Providing solutions for the instability of function approximation techniques in RL, this work was the first to convincingly demonstrate that RL agents could be trained on raw, high-dimensional observations, solely based on a reward signal. The second standout success was the development of a hybrid DRL system, AlphaGo, that defeated a human world champion in Go [128], paralleling the historic achievement of IBM’s Deep Blue in chess two decades earlier [19] and IBM’s Watson DeepQA system that beat the best human Jeopardy! players [31]. Unlike the handcrafted rules that have dominated chess-playing systems, AlphaGo was composed of neural networks that were trained using supervised and reinforcement learning, in combination with a traditional heuristic search algorithm.

DRL algorithms have already been applied to a wide range of problems, such as robotics, where control policies for robots can now be learned directly from camera inputs in the real world [74, 75], succeeding controllers that used to be hand-engineered or learned from low-dimensional features of the robot’s state. In a step towards even more capable agents, DRL has been used to create agents that can meta-learn (“learn to learn”) [29, 156], allowing them to generalise to complex visual environments they have never seen before [29]. In Figure 1, we showcase just some of the domains that DRL has been applied to, ranging from playing video games [84] to indoor navigation [167].

Video games may be an interesting challenge, but learning how to play them is not the end goal of DRL. One of the driving forces behind DRL is the vision of creating systems that are capable of learning how to adapt in the real world. From managing power consumption [142] to picking and stowing objects [75], DRL stands to increase the amount of physical tasks that can be automated by learning. However, DRL does not stop there, as RL is a general way of approaching optimisation problems by trial and error. From designing state-of-the-art machine translation models [168] to constructing new optimisation functions [76], DRL has already been used to approach all manner of machine learning tasks.

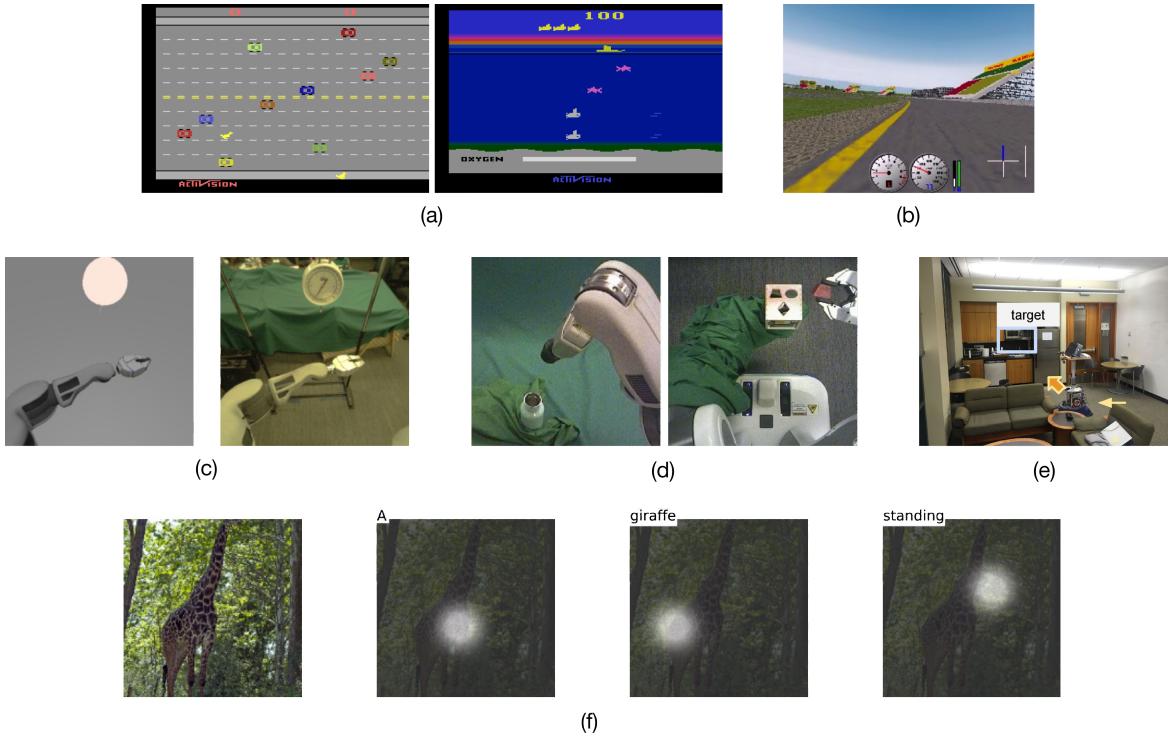


Fig. 1. A range of visual RL domains. (a) Two classic Atari 2600 video games, “Freeway” and “Seaquest”, from the Arcade Learning Environment (ALE) [10]. Due to the range of supported games that vary in genre, visuals and difficulty, the ALE has become a standard testbed for DRL algorithms [84, 95, 44, 122, 132, 157, 85]. As we will discuss later, the ALE is one of several benchmarks that are now being used to standardise evaluation in RL. (b) The TORCS car racing simulator, which has been used to test DRL algorithms that can output continuous actions [64, 79, 85] (as the games from the ALE only support discrete actions). (c) Utilising the potentially unlimited amount of training data that can be amassed in robotic simulators, several methods aim to transfer knowledge from the simulator to the real world [22, 115, 146]. (d) Two of the four robotic tasks designed by Levine et al. [74]: screwing on a bottle cap and placing a shaped block in the correct hole. Levine et al. [74] were able to train visuomotor policies in an end-to-end fashion, showing that visual servoing could be learned directly from raw camera inputs by using deep neural networks. (e) A real room, in which a wheeled robot trained to navigate the building is given a visual cue as input, and must find the corresponding location [167]. (f) A natural image being captioned by a neural network that uses reinforcement learning to choose where to look [166]. By processing a small portion of the image for every word generated, the network can focus its attention on the most salient points. Figures reproduced from [10, 79, 146, 74, 167, 166], respectively.

And, in the same way that deep learning has been utilised across many branches of machine learning, it seems likely that in the future, DRL will be an important component in constructing general AI systems [68].

## II. REWARD-DRIVEN BEHAVIOUR

Before examining the contributions of deep neural networks to RL, we will introduce the field of RL in general. The essence of RL is learning through *interaction*. An RL agent interacts with its environment and, upon observing the consequences of its actions, can learn to alter its own behaviour in response to rewards received. This paradigm of trial-and-error learning has its roots in behaviourist psychology, and is one of the main foundations of RL [135]. The other key influence on RL is optimal control, which has lent the mathematical formalisms (most notably dynamic programming [13]) that underpin the field.

In the RL set-up, an autonomous *agent*, controlled by a machine learning algorithm, observes a *state*  $s_t$  from its *environment* at timestep  $t$ . The agent interacts with the environment by taking an *action*  $a_t$  in state  $s_t$ . When the agent takes an action, the environment and the agent transition to a new state  $s_{t+1}$  based on the current state and the chosen action. The state is a sufficient statistic of the environment

and thereby comprises all the necessary information for the agent to take the best action, which can include parts of the agent, such as the position of its actuators and sensors. In the optimal control literature, states and actions are often denoted by  $x_t$  and  $u_t$ , respectively.

The best sequence of actions is determined by the *rewards* provided by the environment. Every time the environment transitions to a new state, it also provides a scalar reward  $r_{t+1}$  to the agent as feedback. The goal of the agent is to learn a *policy* (control strategy)  $\pi$  that maximises the expected *return* (cumulative, discounted reward). Given a state, a policy returns an action to perform; an *optimal policy* is any policy that maximises the expected return in the environment. In this respect, RL aims to solve the same problem as optimal control. However, the challenge in RL is that the agent needs to learn about the consequences of actions in the environment by trial and error, as, unlike in optimal control, a model of the state transition dynamics is not available to the agent. Every interaction with the environment yields information, which the agent uses to update its knowledge. This *perception-action-learning loop* is illustrated in Figure 2.

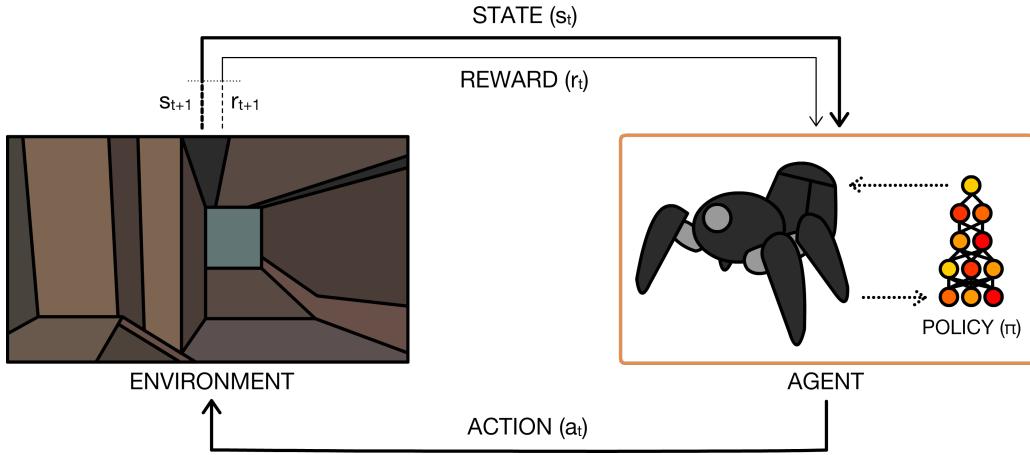


Fig. 2. The perception-action-learning loop. At time  $t$ , the agent receives state  $s_t$  from the environment. The agent uses its policy to choose an action  $a_t$ . Once the action is executed, the environment transitions a step, providing the next state  $s_{t+1}$  as well as feedback in the form of a reward  $r_{t+1}$ . The agent uses knowledge of state transitions, of the form  $(s_t, a_t, s_{t+1}, r_{t+1})$ , in order to learn and improve its policy.

### A. Markov Decision Processes

Formally, RL can be described as a Markov decision process (MDP), which consists of:

- A set of states  $\mathcal{S}$ , plus a distribution of starting states  $p(s_0)$ .
- A set of actions  $\mathcal{A}$ .
- Transition dynamics  $\mathcal{T}(s_{t+1}|s_t, a_t)$  that map a state-action pair at time  $t$  onto a distribution of states at time  $t + 1$ .
- An immediate/instantaneous reward function  $\mathcal{R}(s_t, a_t, s_{t+1})$ .
- A discount factor  $\gamma \in [0, 1]$ , where lower values place more emphasis on immediate rewards.

In general, the policy  $\pi$  is a mapping from states to a probability distribution over actions:  $\pi : \mathcal{S} \rightarrow p(\mathcal{A} = a|\mathcal{S})$ . If the MDP is *episodic*, i.e., the state is reset after each episode of length  $T$ , then the sequence of states, actions and rewards in an episode constitutes a *trajectory* or *rollout* of the policy. Every rollout of a policy accumulates rewards from the environment, resulting in the return  $R = \sum_{t=0}^{T-1} \gamma^t r_{t+1}$ . The goal of RL is to find an optimal policy,  $\pi^*$ , which achieves the maximum expected return from all states:

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E}[R|\pi]. \quad (1)$$

It is also possible to consider non-episodic MDPs, where  $T = \infty$ . In this situation,  $\gamma < 1$  prevents an infinite sum of rewards from being accumulated. Furthermore, methods that rely on complete trajectories are no longer applicable, but those that use a finite set of transitions still are.

A key concept underlying RL is the Markov property—only the current state affects the next state, or in other words, the future is conditionally independent of the past given the present state. This means that any decisions made at  $s_t$  can be based solely on  $s_{t-1}$ , rather than  $\{s_0, s_1, \dots, s_{t-1}\}$ . Although this assumption is held by the majority of RL algorithms, it is somewhat unrealistic, as it requires the states to be *fully observable*. A generalisation of MDPs are partially observable MDPs (POMDPs), in which the agent receives an

observation  $\mathbf{o}_t \in \Omega$ , where the distribution of the observation  $p(\mathbf{o}_{t+1}|s_{t+1}, a_t)$  is dependent on the current state and the previous action [56]. In a control and signal processing context, the observation would be described by a measurement/observation mapping in a state-space-model that depends on the current state and the previously applied action.

POMDP algorithms typically maintain a *belief* over the current state given the previous belief state, the action taken and the current observation. A more common approach in deep learning is to utilise recurrent neural networks (RNNs) [163, 44, 45, 85, 96], which, unlike feedforward neural networks, are dynamical systems. This approach to solving POMDPs is related to other problems using dynamical systems and state space models, where the true state can only be estimated [16].

### B. Challenges in RL

It is instructive to emphasise some challenges faced in RL:

- The optimal policy must be inferred by trial-and-error interaction with the environment. The only learning signal the agent receives is the reward.
- The observations of the agent depend on its actions and can contain strong temporal correlations.
- Agents must deal with long-range time dependencies: Often the consequences of an action only materialise after many transitions of the environment. This is known as the (temporal) *credit assignment problem* [135].

We will illustrate these challenges in the context of an indoor robotic visual navigation task: if the goal location is specified, we may be able to estimate the distance remaining (and use it as a reward signal), but it is unlikely that we will know exactly what series of actions the robot needs to take to reach the goal. As the robot must choose where to go as it navigates the building, its decisions influence which rooms it sees and, hence, the statistics of the visual sequence captured. Finally, after navigating several junctions, the robot may find itself in a dead end. There is a range of problems, from learning the consequences of actions to balancing exploration

against exploitation, but ultimately these can all be addressed formally within the framework of RL.

### III. REINFORCEMENT LEARNING ALGORITHMS

So far, we have introduced the key formalism used in RL, the MDP, and briefly noted some challenges in RL. In the following, we will distinguish between different classes of RL algorithms. There are two main approaches to solving RL problems: methods based on *value functions* and methods based on *policy search*. There is also a hybrid, *actor-critic* approach, which employs both value functions and policy search. We will now explain these approaches and other useful concepts for solving RL problems.

#### A. Value Functions

Value function methods are based on estimating the value (expected return) of being in a given state. The *state-value function*  $V^\pi(\mathbf{s})$  is the expected return when starting in state  $\mathbf{s}$  and following  $\pi$  henceforth:

$$V^\pi(\mathbf{s}) = \mathbb{E}[R|\mathbf{s}, \pi] \quad (2)$$

The optimal policy,  $\pi^*$ , has a corresponding state-value function  $V^*(\mathbf{s})$ , and vice-versa, the optimal state-value function can be defined as

$$V^*(\mathbf{s}) = \max_{\pi} V^\pi(\mathbf{s}) \quad \forall \mathbf{s} \in \mathcal{S}. \quad (3)$$

If we had  $V^*(\mathbf{s})$  available, the optimal policy could be retrieved by choosing among all actions available at  $\mathbf{s}_t$  and picking the action  $\mathbf{a}$  that maximises  $\mathbb{E}_{\mathbf{s}_{t+1} \sim \mathcal{T}(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a})}[V^*(\mathbf{s}_{t+1})]$ .

In the RL setting, the transition dynamics  $\mathcal{T}$  are unavailable. Therefore, we construct another function, the *state-action-value* or *quality function*  $Q^\pi(\mathbf{s}, \mathbf{a})$ , which is similar to  $V^\pi$ , except that the initial action  $\mathbf{a}$  is provided, and  $\pi$  is only followed from the succeeding state onwards:

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}[R|\mathbf{s}, \mathbf{a}, \pi]. \quad (4)$$

The best policy, given  $Q^\pi(\mathbf{s}, \mathbf{a})$ , can be found by choosing  $\mathbf{a}$  greedily at every state:  $\text{argmax}_{\mathbf{a}} Q^\pi(\mathbf{s}, \mathbf{a})$ . Under this policy, we can also define  $V^\pi(\mathbf{s})$  by maximising  $Q^\pi(\mathbf{s}, \mathbf{a})$ :  $V^\pi(\mathbf{s}) = \max_{\mathbf{a}} Q^\pi(\mathbf{s}, \mathbf{a})$ .

**Dynamic Programming:** To actually learn  $Q^\pi$ , we exploit the Markov property and define the function as a Bellman equation [13], which has the following recursive form:

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \mathbb{E}_{\mathbf{s}_{t+1}}[r_{t+1} + \gamma Q^\pi(\mathbf{s}_{t+1}, \pi(\mathbf{s}_{t+1}))]. \quad (5)$$

This means that  $Q^\pi$  can be improved by *bootstrapping*, i.e., we can use the current values of our estimate of  $Q^\pi$  to improve our estimate. This is the foundation of *Q-learning* [159] and the state-action-reward-state-action (SARSA) algorithm [112]:

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) \leftarrow Q^\pi(\mathbf{s}_t, \mathbf{a}_t) + \alpha \delta, \quad (6)$$

where  $\alpha$  is the learning rate and  $\delta = Y - Q^\pi(\mathbf{s}_t, \mathbf{a}_t)$  the temporal difference (TD) error; here,  $Y$  is a target as in a standard regression problem. SARSA, an *on-policy* learning algorithm, is used to improve the estimate of  $Q^\pi$  by using transitions generated by the behavioural policy (the policy derived from

$Q^\pi$ ), which results in setting  $Y = r_t + \gamma Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})$ . *Q*-learning is *off-policy*, as  $Q^\pi$  is instead updated by transitions that were not necessarily generated by the derived policy. Instead, *Q*-learning uses  $Y = r_t + \gamma \max_{\mathbf{a}} Q^\pi(\mathbf{s}_{t+1}, \mathbf{a})$ , which directly approximates  $Q^*$ .

To find  $Q^*$  from an arbitrary  $Q^\pi$ , we use *generalised policy iteration*, where policy iteration consists of *policy evaluation* and *policy improvement*. Policy evaluation improves the estimate of the value function, which can be achieved by minimising TD errors from trajectories experienced by following the policy. As the estimate improves, the policy can naturally be improved by choosing actions greedily based on the updated value function. Instead of performing these steps separately to convergence (as in policy iteration), generalised policy iteration allows for interleaving the steps, such that progress can be made more rapidly.

#### B. Sampling

Instead of bootstrapping value functions using dynamic programming methods, Monte Carlo methods estimate the expected return (2) from a state by averaging the return from multiple rollouts of a policy. Because of this, pure Monte Carlo methods can also be applied in non-Markovian environments. On the other hand, they can only be used in episodic MDPs, as a rollout has to terminate for the return to be calculated. It is possible to get the best of both methods by combining TD learning and Monte Carlo policy evaluation, as in done in the TD( $\lambda$ ) algorithm [135]. Similarly to the discount factor, the  $\lambda$  in TD( $\lambda$ ) is used to interpolate between Monte Carlo evaluation and bootstrapping. As demonstrated in Figure 3, this results in an entire spectrum of RL methods based around the amount of sampling utilised.

Another major value-function based method relies on learning the *advantage* function  $A^\pi(\mathbf{s}, \mathbf{a})$  [6, 43]. Unlike producing absolute state-action values, as with  $Q^\pi$ ,  $A^\pi$  instead represents relative state-action values. Learning relative values is akin to removing a baseline or average level of a signal; more intuitively, it is easier to learn that one action has better consequences than another, than it is to learn the actual return from taking the action.  $A^\pi$  represents a relative advantage of actions through the simple relationship  $A^\pi = Q^\pi - V^\pi$ , and is also closely related to the baseline method of variance reduction within gradient-based policy search methods [164]. The idea of advantage updates has been utilised in many recent DRL algorithms [157, 40, 85, 123].

#### C. Policy Search

Policy search methods do not need to maintain a value function model, but directly search for an optimal policy  $\pi^*$ . Typically, a parameterised policy  $\pi_\theta$  is chosen, whose parameters are updated to maximise the expected return  $\mathbb{E}[R|\theta]$  using either gradient-based or gradient-free optimisation [26]. Neural networks that encode policies have been successfully trained using both gradient-free [37, 23, 64] and gradient-based [164, 163, 46, 79, 122, 123, 74] methods. Gradient-free optimisation can effectively cover low-dimensional parameter spaces, but despite some successes in applying them to large networks [64], gradient-based training remains the method of choice for most DRL algorithms, being more sample-efficient

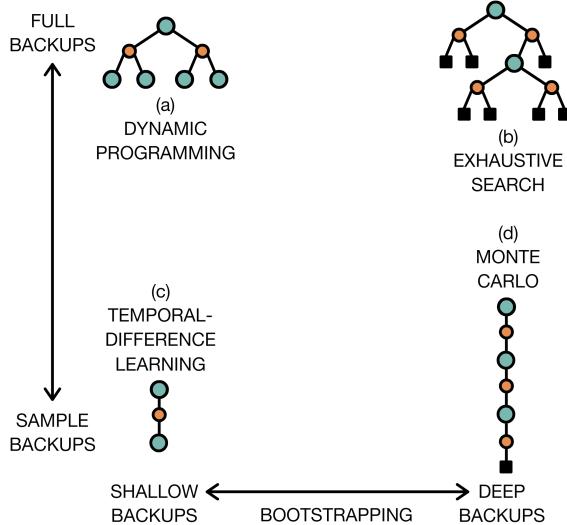


Fig. 3. Two dimensions of RL algorithms, based on the *backups* used to learn or construct a policy. At the extremes of these dimensions are (a) dynamic programming, (b) exhaustive search, (c) one-step TD learning and (d) pure Monte Carlo approaches. Bootstrapping extends from (c) 1-step TD learning to  $n$ -step TD learning methods [135], with (d) pure Monte Carlo approaches not relying on bootstrapping at all. Another possible dimension of variation is choosing to (c, d) sample actions versus (a, b) taking the expectation over all choices. Recreated from [135].

when policies possess a large number of parameters.

When constructing the policy directly, it is common to output parameters for a probability distribution; for continuous actions, this could be the mean and standard deviations of Gaussian distributions, whilst for discrete actions this could be the individual probabilities of a multinomial distribution. The result is a stochastic policy from which we can directly sample actions. With gradient-free methods, finding better policies requires a heuristic search across a predefined class of models. Methods such as evolution strategies essentially perform hill-climbing in a subspace of policies [116], whilst more complex methods, such as compressed network search, impose additional inductive biases [64]. Perhaps the greatest advantage of gradient-free policy search is that they can also optimise non-differentiable policies.

**Policy Gradients:** Gradients can provide a strong learning signal as to how to improve a parameterised policy. However, to compute the expected return (1) we need to average over plausible trajectories induced by the current policy parameterisation. This averaging requires either deterministic approximations (e.g., linearisation) or stochastic approximations via sampling [26]. Deterministic approximations can only be applied in a model-based setting where a model of the underlying transition dynamics is available. In the more common model-free RL setting, a Monte Carlo estimate of the expected return is determined. For gradient-based learning, this Monte Carlo approximation poses a challenge since gradients cannot pass through these samples of a stochastic function. Therefore, we turn to an estimator of the gradient, known in RL as the REINFORCE rule [164], elsewhere known as the score function [34] or likelihood-ratio estimator [36]. The latter name is telling as using the estimator is similar to the practice of optimising

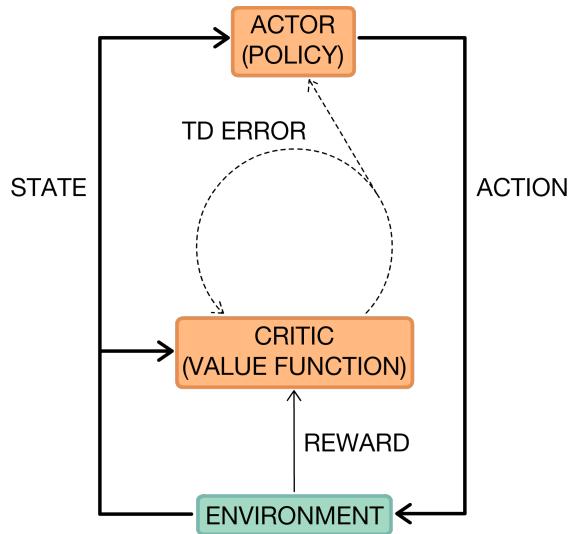


Fig. 4. Actor-critic set-up. The actor (policy) receives a state from the environment and chooses an action to perform. At the same time, the critic (value function) receives the state and reward resulting from the previous interaction. The critic uses the TD error calculated from this information to update itself and the actor. Recreated from [135].

the log-likelihood in supervised learning. Intuitively, gradient ascent using the estimator increases the log probability of the sampled action, weighted by the return. More formally, the REINFORCE rule can be used to compute the gradient of an expectation over a function  $f$  of a random variable  $X$  with respect to parameters  $\theta$ :

$$\nabla_{\theta} \mathbb{E}_X[f(X; \theta)] = \mathbb{E}_X[f(X; \theta) \nabla_{\theta} \log p(X)]. \quad (7)$$

As this computation relies on the empirical return of a trajectory, the resulting gradients possess a high variance. By introducing unbiased estimates that are less noisy it is possible to reduce the variance. The general methodology for performing this is to subtract a baseline, which means weighting updates by an advantage rather than the pure return. The simplest baseline is the average return taken over several episodes [164], but many more options are available [123].

**Actor-critic Methods:** It is possible to combine value functions with an explicit representation of the policy, resulting in actor-critic methods, as shown in Figure 4. The “actor” (policy) learns by using feedback from the “critic” (value function). In doing so, these methods trade off variance reduction of policy gradients with bias introduction from value function methods [63, 123].

Actor-critic methods use the value function as a baseline for policy gradients, such that the only fundamental difference between actor-critic methods and other baseline methods are that actor-critic methods utilise a *learned* value function. For this reason, we will later discuss actor-critic methods as a subset of policy gradient methods.

#### D. Planning and Learning

Given a model of the environment, it is possible to use dynamic programming over all possible actions (Figure 3

(a)), sample trajectories for heuristic search (as was done by AlphaGo [128]), or even perform an exhaustive search (Figure 3 (b)). Sutton and Barto [135] define *planning* as any method which utilises a model to produce or improve a policy. This includes *distribution models*, which include  $\mathcal{T}$  and  $\mathcal{R}$ , and *sample models*, from which only samples of transitions can be drawn.

In RL, we focus on learning without access to the underlying model of the environment. However, interactions with the environment could be used to learn value functions, policies, and also a model. Model-free RL methods learn directly from interactions with the environment, but model-based RL methods can simulate transitions using the learned model, resulting in increased sample efficiency. This is particularly important in domains where each interaction with the environment is expensive. However, learning a model introduces extra complexities, and there is always the danger of suffering from model errors, which in turn affects the learned policy; a common but partial solution in this latter scenario is to use model predictive control, where planning is repeated after small sequences of actions in the real environment [16]. Although deep neural networks can potentially produce very complex and rich models [95, 132, 32], sometimes simpler, more data-efficient methods are preferable [40]. These considerations also play a role in actor-critic methods with learned value functions [63, 123].

#### E. The Rise of DRL

Many of the successes in DRL have been based on scaling up prior work in RL to high-dimensional problems. This is due to the learning of low-dimensional feature representations and the powerful function approximation properties of neural networks. By means of representation learning, DRL can deal efficiently with the curse of dimensionality, unlike tabular and traditional non-parametric methods [15]. For instance, convolutional neural networks (CNNs) can be used as components of RL agents, allowing them to learn directly from raw, high-dimensional visual inputs. In general, DRL is based on training deep neural networks to approximate the optimal policy  $\pi^*$ , and/or the optimal value functions  $V^*$ ,  $Q^*$  and  $A^*$ .

Although there have been DRL successes with gradient-free methods [37, 23, 64], the vast majority of current works rely on gradients and hence the backpropagation algorithm [162, 111]. The primary motivation is that when available, gradients provide a strong learning signal. In reality, these gradients are estimated based on approximations, through sampling or otherwise, and as such we have to craft algorithms with useful inductive biases in order for them to be tractable.

The other benefit of backpropagation is to view the optimisation of the expected return as the optimisation of a stochastic function [121, 46]. This function can comprise of several parts—models, policies and value functions—which can be combined in various ways. The individual parts, such as value functions, may not directly optimise the expected return, but can instead embody useful information about the RL domain. For example, using a differentiable model and policy, it is possible to forward propagate and backpropagate through entire rollouts; on the other hand, inaccuracies can accumulate

over long time steps, and it may be pertinent to instead use a value function to summarise the statistics of the rollouts [46]. We have previously mentioned that representation learning and function approximation are key to the success of DRL, but it is also true to say that the field of deep learning has inspired new ways of thinking about RL.

Following our review of RL, we will now partition the next part of the survey into value function and policy search methods in DRL, starting with the well-known deep  $Q$ -network (DQN) [84]. In these sections, we will focus on state-of-the-art techniques, as well as the historical works they are built upon. The focus of the state-of-the-art techniques will be on those for which the state space is conveyed through visual inputs, e.g., images and video. To conclude, we will examine ongoing research areas and open challenges.

#### IV. VALUE FUNCTIONS

The well-known function approximation properties of neural networks led naturally to the use of deep learning to regress functions for use in RL agents. Indeed, one of the earliest success stories in RL is TD-Gammon, a neural network that reached expert-level performance in Backgammon in the early 90s [141]. Using TD methods, the network took in the state of the board to predict the probability of black or white winning. Although this simple idea has been echoed in later work [128], progress in RL research has favoured the explicit use of value functions, which can capture the structure underlying the environment. From early value function methods in DRL, which took simple states as input [109], current methods are now able to tackle visually and conceptually complex environments [84, 122, 85, 96, 167].

##### A. Function Approximation and the DQN

We begin our survey of value-function-based DRL algorithms with the DQN [84], pictured in Figure 5, which achieved scores across a wide range of classic Atari 2600 video games [10] that were comparable to that of a professional video games tester. The inputs to the DQN are four greyscale frames of the game, concatenated over time, which are initially processed by several convolutional layers in order to extract spatiotemporal features, such as the movement of the ball in “Pong” or “Breakout.” The final feature map from the convolutional layers is processed by several fully connected layers, which more implicitly encode the effects of actions. This contrasts with more traditional controllers that use fixed preprocessing steps, which are therefore unable to adapt their processing of the state in response to the learning signal.

A forerunner of the DQN—neural fitted  $Q$  iteration (NFQ)—involved training a neural network to return the  $Q$ -value given a state-action pair [109]. NFQ was later extended to train a network to drive a slot car using raw visual inputs from a camera over the race track, by combining a deep autoencoder to reduce the dimensionality of the inputs with a separate branch to predict  $Q$ -values [69]. Although the previous network could have been trained for both reconstruction and RL tasks simultaneously, it was both more reliable and computationally efficient to train the two parts of the network sequentially.

The DQN [84] is closely related to the model proposed

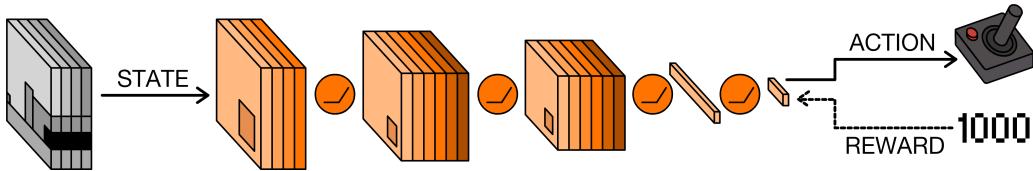


Fig. 5. The deep  $Q$ -network [84]. The network takes the state—a stack of greyscale frames from the video game—and processes it with convolutional and fully connected layers, with ReLU nonlinearities in between each layer. At the final layer, the network outputs a discrete action, which corresponds to one of the possible control inputs for the game. Given the current state and chosen action, the game returns a new score. The DQN uses the reward—the difference between the new score and the previous one—to learn from its decision. More precisely, the reward is used to update its estimate of  $Q$ , and the error between its previous estimate and its new estimate is backpropagated through the network.

by Lange et al. [69], but was the first RL algorithm that was demonstrated to work directly from raw visual inputs and on a wide variety of environments. It was designed such that the final fully connected layer outputs  $Q^\pi(s, \cdot)$  for all action values in a discrete set of actions—in this case, the various directions of the joystick and the fire button. This not only enables the best action,  $\text{argmax}_a Q^\pi(s, a)$ , to be chosen after a single forward pass of the network, but also allows the network to more easily encode action-independent knowledge in the lower, convolutional layers. With merely the goal of maximising its score on a video game, the DQN learns to extract salient visual features, jointly encoding objects, their movements, and, most importantly, their interactions. Using techniques originally developed for explaining the behaviour of CNNs in object recognition tasks, we can also inspect what parts of its view the agent considers important (see Figure 6).

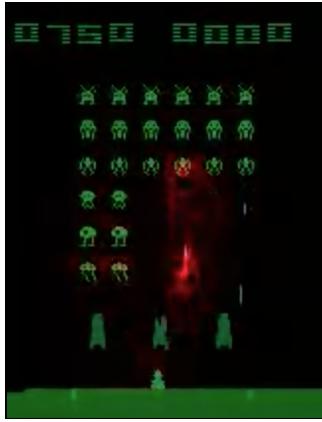


Fig. 6. Saliency map of a trained DQN [84] playing “Space Invaders” [10]. By backpropagating the training signal to the image space, it is possible to see what a neural-network-based agent is attending to. In this frame, the most salient points—shown with the red overlay—are the laser that the agent recently fired, and also the enemy that it anticipates hitting in a few time steps.

The true underlying state of the game is contained within 128 bytes of Atari 2600 RAM. However, the DQN was designed to directly learn from visual inputs ( $210 \times 160$ pixel 8-bit RGB images), which it takes as the state  $s$ . It is impractical to represent  $Q^\pi(s, a)$  exactly as a lookup table: When combined with 18 possible actions, we obtain a  $Q$ -table of size  $|\mathcal{S}| \times |\mathcal{A}| = 18 \times 256^{3 \times 210 \times 160}$ . Even if it were feasible to create such a table, it would be sparsely populated, and information gained from one state-action pair cannot be propagated to other state-action pairs. The strength of the DQN

lies in its ability to compactly represent both high-dimensional observations and the  $Q$ -function using deep neural networks. Without this ability, tackling the discrete Atari domain from raw visual inputs would be impractical.

The DQN addressed the fundamental instability problem of using function approximation in RL [145] by the use of two techniques: experience replay [80] and target networks. Experience replay memory stores transitions of the form  $(s_t, a_t, s_{t+1}, r_{t+1})$  in a cyclic buffer, enabling the RL agent to sample from and train on previously observed data offline. Not only does this massively reduce the amount of interactions needed with the environment, but batches of experience can be sampled, reducing the variance of learning updates. Furthermore, by sampling uniformly from a large memory, the temporal correlations that can adversely affect RL algorithms are broken. Finally, from a practical perspective, batches of data can be efficiently processed in parallel by modern hardware, increasing throughput. Whilst the original DQN algorithm used uniform sampling [84], later work showed that prioritising samples based on TD errors is more effective for learning [118]. We note that although experience replay is typically thought of as a model-free technique, it could actually be considered a simple model [150].

The second stabilising method, introduced by Mnih et al. [84], is the use of a target network that initially contains the weights of the network enacting the policy, but is kept frozen for a large period of time. Rather than having to calculate the TD error based on its own rapidly fluctuating estimates of the  $Q$ -values, the policy network uses the fixed target network. During training, the weights of the target network are updated to match the policy network after a fixed number of steps. Both experience replay and target networks have gone on to be used in subsequent DRL works [40, 79, 158, 89].

### B. $Q$ -Function Modifications

Considering that one of the key components of the DQN is a function approximator for the  $Q$ -function, it can benefit from fundamental advances in RL. van Hasselt [148] showed that the single estimator used in the  $Q$ -learning update rule overestimates the expected return due to the use of the maximum action value as an approximation of the maximum *expected* action value. Double- $Q$  learning provides a better estimate through the use of a double estimator [148]. Whilst double- $Q$  learning requires an additional function to be learned, later work proposed using the already available target network from the DQN algorithm, resulting in significantly better results with only a small change in the update step [149]. A more

radical proposal by Bellemare et al. [12] was to actually learn the full *value distribution*, rather than just the expectation; this provides additional information, such as whether the potential rewards come from a skewed or multimodal distribution. Although the resulting algorithm—based on learning categorical distributions—was used to construct the Categorical DQN, the benefits can potentially be applied to any RL algorithm that utilises learned value functions.

Yet another way to adjust the DQN architecture is to decompose the  $Q$ -function into meaningful functions, such as constructing  $Q^\pi$  by adding together separate layers that compute the state-value function  $V^\pi$  and advantage function  $A^\pi$  [157]. Rather than having to come up with accurate  $Q$ -values for all actions, the duelling DQN [157] benefits from a single baseline for the state in the form of  $V^\pi$ , and easier-to-learn relative values in the form of  $A^\pi$ . The combination of the duelling DQN with prioritised experience replay [118] is one of the state-of-the-art techniques in discrete action settings. Further insight into the properties of  $A^\pi$  by Gu et al. [40] led them to modify the DQN with a convex advantage layer that extended the algorithm to work over sets of continuous actions, creating the normalised advantage function (NAF) algorithm. Benefiting from experience replay, target networks and advantage updates, NAF is one of several state-of-the-art techniques in continuous control problems [40].

Some RL domains, such as recommender systems, have very large discrete action spaces, and hence may be difficult to directly deal with. Dulac-Arnold et al. [30] proposed learning “action embeddings” over the large set of original actions, and then using  $k$ -nearest neighbors to produce “proto-actions” which can be used with traditional RL methods. The idea of using representation learning to create distributed embeddings is a particular strength of DRL, and has been successfully utilised for other purposes [161, 100]. Another related scenario in RL is when many actions need to be made simultaneously, such as specifying the torques in a many-jointed robot, which results in the action space growing exponentially. A naive but reasonable approach is to factorise the policy, treating each action independently [115]. An alternative is to construct an autoregressive policy, where each action in a single timestep is predicted conditionally on the state and previously chosen actions from the same timestep [106, 5, 168]. Metz et al. [81] used this idea in order to construct the sequential DQN, allowing them to discretise a large action space and outperform NAF—which is limited by its quadratic advantage function—in continuous control problems. In a broader context, rather than dealing directly with primitive actions directly, one may choose to invoke “subpolicies” from higher-level policies [136]; this concept, known as hierarchical reinforcement learning (HRL), will be discussed later.

## V. POLICY SEARCH

Policy search methods aim to directly find policies by means of gradient-free or gradient-based methods. Prior to the current surge of interest in DRL, several successful methods in DRL eschewed the commonly used backpropagation algorithm in favour of evolutionary algorithms [37, 23, 64], which are gradient-free policy search algorithms. Evolutionary methods

rely on evaluating the performance of a population of agents. Hence, they are expensive for large populations or agents with many parameters. However, as black-box optimisation methods they can be used to optimise arbitrary, non-differentiable models and naturally allow for more exploration in parameter space. In combination with a compressed representation of neural network weights, evolutionary algorithms can even be used to train large networks; such a technique resulted in the first deep neural network to learn an RL task, straight from high-dimensional visual inputs [64]. Recent work has reignited interest in evolutionary methods for RL as they can potentially be distributed at larger scales than techniques that rely on gradients [116].

### A. Backpropagation through Stochastic Functions

The workhorse of DRL, however, remains backpropagation [162, 111]. The previously discussed REINFORCE rule [164] allows neural networks to learn stochastic policies in a task-dependent manner, such as deciding where to look in an image to track [120], classify [83] or caption objects [166]. In these cases, the stochastic variable would determine the coordinates of a small crop of the image, and hence reduce the amount of computation needed. This usage of RL to make discrete, stochastic decisions over inputs is known in the deep learning literature as *hard attention*, and is one of the more compelling uses of basic policy search methods in recent years, having many applications outside of traditional RL domains. More generally, the ability to backpropagate through stochastic functions, using techniques such as REINFORCE [164] or the “reparameterisation trick” [61, 108], allows neural networks to be treated as stochastic computation graphs that can be optimised over [121], which is a key concept in algorithms such as stochastic value gradients (SVGs) [46].

### B. Compounding Errors

Searching directly for a policy represented by a neural network with very many parameters can be difficult and can suffer from severe local minima. One way around this is to use guided policy search (GPS), which takes a few sequences of actions from another controller (which could be constructed using a separate method, such as optimal control). GPS learns from them by using supervised learning in combination with importance sampling, which corrects for off-policy samples [73]. This approach effectively biases the search towards a good (local) optimum. GPS works in a loop, by optimising policies to match sampled trajectories, and optimising trajectory distributions to match the policy and minimise costs. Initially, GPS was used to train neural networks on simulated continuous RL problems [72], but was later utilised to train a policy for a real robot based on visual inputs [74]. This research by Levine et al. [74] showed that it was possible to train visuomotor policies for a robot “end-to-end”, straight from the RGB pixels of the camera to motor torques, and, hence, is one of the seminal works in DRL.

A more commonly used method is to use a trust region, in which optimisation steps are restricted to lie within a region where the approximation of the true cost function still holds. By preventing updated policies from deviating too wildly from previous policies, the chance of a catastrophically bad

update is lessened, and many algorithms that use trust regions guarantee or practically result in monotonic improvement in policy performance. The idea of constraining each policy gradient update, as measured by the Kullback-Leibler (KL) divergence between the current and proposed policy, has a long history in RL [57, 4, 59, 103]. One of the newer algorithms in this line of work, trust region policy optimisation (TRPO), has been shown to be relatively robust and applicable to domains with high-dimensional inputs [122]. To achieve this, TRPO optimises a *surrogate* objective function—specifically, it optimises an (importance sampled) advantage estimate, constrained using a quadratic approximation of the KL divergence. Whilst TRPO can be used as a pure policy gradient method with a simple baseline, later work by Schulman et al. [123] introduced generalised advantage estimation (GAE), which proposed several, more advanced variance reduction baselines. The combination of TRPO and GAE remains one of the state-of-the-art RL techniques in continuous control. However, the constrained optimisation of TRPO requires calculating second-order gradients, limiting its applicability. In contrast, the newer proximal policy optimisation (PPO) algorithm performs unconstrained optimisation, requiring only first-order gradient information [1, 47, 125]. The two main variants include an adaptive penalty on the KL divergence, and a heuristic clipped objective which is independent of the KL divergence [125]. Being less expensive whilst retaining the performance of TRPO means that PPO (with or without GAE) is gaining popularity for a range of RL tasks [47, 125].

### C. Actor-Critic Methods

Instead of utilising the average of several Monte Carlo returns as the baseline for policy gradient methods, actor-critic approaches have grown in popularity as an effective means of combining the benefits of policy search methods with learned value functions, which are able to learn from full returns and/or TD errors. They can benefit from improvements in both policy gradient methods, such as GAE [123], and value function methods, such as target networks [84]. In the last few years, DRL actor-critic methods have been scaled up from learning simulated physics tasks [46, 79] to real robotic visual navigation tasks [167], directly from image pixels.

One recent development in the context of actor-critic algorithms are deterministic policy gradients (DPGs) [127], which extend the standard policy gradient theorems for stochastic policies [164] to deterministic policies. One of the major advantages of DPGs is that, whilst stochastic policy gradients integrate over both state and action spaces, DPGs only integrate over the state space, requiring fewer samples in problems with large action spaces. In the initial work on DPGs, Silver et al. [127] introduced and demonstrated an off-policy actor-critic algorithm that vastly improved upon a stochastic policy gradient equivalent in high-dimensional continuous control problems. Later work introduced deep DPG (DDPG), which utilised neural networks to operate on high-dimensional, visual state spaces [79]. In the same vein as DPGs, Heess et al. [46] devised a method for calculating gradients to optimise stochastic policies, by “reparameterising” [61, 108] the stochasticity away from the network, thereby

allowing standard gradients to be used (instead of the high-variance REINFORCE estimator [164]). The resulting SVG methods are flexible, and can be used both with (SVG(0) and SVG(1)) and without (SVG( $\infty$ )) value function critics, and with (SVG( $\infty$ ) and SVG(1)) and without (SVG(0)) models. Later work proceeded to integrate DPGs and SVGs with RNNs, allowing them to solve continuous control problems in POMDPs, learning directly from pixels [45].

Value functions introduce a broadly applicable benefit in actor-critic methods—the ability to use off-policy data. On-policy methods can be more stable, whilst off-policy methods can be more data efficient, and hence there have been several attempts to merge the two [158, 94, 41, 39, 42]. Earlier work has either utilised a mix of on-policy and off-policy gradient updates [158, 94, 39], or used the off-policy data to train a value function in order to reduce the variance of on-policy gradient updates [41]. The more recent work by Gu et al. [42] unified these methods under interpolated policy gradients (IPGs), resulting in one of the newest state-of-the-art continuous DRL algorithms, and also providing insights for future research in this area. Together, the ideas behind IPGs and SVGs (of which DPGs can be considered a special case) form algorithmic approaches for improving learning efficiency in DRL.

An orthogonal approach to speeding up learning is to exploit parallel computation. In particular, methods for training networks through asynchronous gradient updates have been developed for use on both single machines [107] and distributed systems [25]. By keeping a canonical set of parameters that are read by and updated in an asynchronous fashion by multiple copies of a single network, computation can be efficiently distributed over both processing cores in a single CPU, and across CPUs in a cluster of machines. Using a distributed system, Nair et al. [91] developed a framework for training multiple DQNs in parallel, achieving both better performance and a reduction in training time. However, the simpler asynchronous advantage actor-critic (A3C) algorithm [85], developed for both single and distributed machine settings, has become one of the most popular DRL techniques in recent times. A3C combines advantage updates with the actor-critic formulation, and relies on asynchronously updated policy and value function networks trained in parallel over several processing threads. The use of multiple agents, situated in their own, independent environments, not only stabilises improvements in the parameters, but conveys an additional benefit in allowing for more exploration to occur. A3C has been used as a standard starting point in many subsequent works, including the work of Zhu et al. [167], who applied it to robotic navigation in the real world through visual inputs. For simplicity, the underlying algorithm may be used with just one agent, termed advantage actor-critic (A2C) [156]. Alternatively, segments from the trajectories of multiple agents can be collected and processed together in a batch, with batch processing more efficiently enabled by GPUs; this synchronous version also goes by the name of A2C [125].

There have been several major advancements on the original A3C algorithm that reflect various motivations in the field of DRL. The first is actor-critic with experience replay [158, 39],

which adds Retrace( $\lambda$ ) off-policy bias correction [88] to a  $Q$ -value-based A3C, allowing it to use experience replay in order to improve sample complexity. Others have attempted to bridge the gap between value and policy-based RL, utilising theoretical advancements to improve upon the original A3C [89, 94, 124]. Finally, there is a growing trend towards exploiting auxiliary tasks to improve the representations learned by DRL agents, and, hence, improve both the learning speed and final performance of these agents [77, 54, 82].

## VI. CURRENT RESEARCH AND CHALLENGES

To conclude, we will highlight some current areas of research in DRL, and the challenges that still remain. Previously, we have focused mainly on model-free methods, but we will now examine a few model-based DRL algorithms in more detail. Model-based RL algorithms play an important role in making RL data-efficient and in trading off exploration and exploitation. After tackling exploration strategies, we shall then address HRL, which imposes an inductive bias on the final policy by explicitly factorising it into several levels. When available, trajectories from other controllers can be used to bootstrap the learning process, leading us to imitation learning and inverse RL (IRL). For the final topic specific to RL, we will look at multi-agent systems, which have their own special considerations. We then bring to attention two broader areas—the use of RNNs, and transfer learning—in the context of DRL. We then examine the issue of evaluating RL, and current benchmarks for DRL.

### A. Model-based RL

The key idea behind model-based RL is to learn a transition model that allows for simulation of the environment without interacting with the environment directly. Model-based RL does not assume specific prior knowledge. However, in practice, we can incorporate prior knowledge (e.g., physics-based models [58]) to speed up learning. Model learning plays an important role in reducing the amount of required interactions with the (real) environment, which may be limited in practice. For example, it is unrealistic to perform millions of experiments with a robot in a reasonable amount of time and without significant hardware wear and tear. There are various approaches to learn predictive models of dynamical systems using pixel information. Based on the deep dynamical model [154], where high-dimensional observations are embedded into a lower-dimensional space using autoencoders, several model-based DRL algorithms have been proposed for learning models and policies from pixel information [95, 160, 155]. If a sufficiently accurate model of the environment can be learned, then even simple controllers can be used to control a robot directly from camera images [32]. Learned models can also be used to guide exploration purely based on simulation of the environment, with deep models allowing these techniques to be scaled up to high-dimensional visual domains [132].

A compelling insight on the benefits of neural-network-based models is that they can overcome some of the problems incurred by planning with imperfect models; in effect, by *embedding* the activations and predictions (outputs) of these models into a vector, a DRL agent can not only obtain more information than just the final result of any model rollouts, but

it can also learn to downplay this information if it believes that the model is inaccurate [161]. This can be more efficient, though less principled, than Bayesian methods for propagating uncertainty [52]. Another way to make use of the flexibility of neural-network-based models is to let them decide when to plan, that is, given a finite amount of computation, whether it is worth modelling one long trajectory, several short trajectories, anything in-between, or simply to take an action in the real environment [100].

Although deep neural networks can make reasonable predictions in simulated environments over hundreds of timesteps [21], they typically require many samples to tune the large amount of parameters they contain. Training these models often requires more samples (interaction with the environment) than simpler models. For this reason, Gu et al. [40] train locally linear models for use with the NAF algorithm—the continuous equivalent of the DQN [84]—to improve the algorithm’s sample complexity in the robotic domain where samples are expensive. In order to spur the adoption of deep models in model-based DRL, it is necessary to find strategies that can be used in order to improve their data efficiency [90].

A less common but potentially useful paradigm exists between model-free and model-based methods—the successor representation (SR) [24]. Rather than picking actions directly or performing planning with models, learning  $\mathcal{T}$  is replaced with learning expected (discounted) future occupancies (SRs), which can be linearly combined with  $\mathcal{R}$  in order to calculate the optimal action; this decomposition makes SRs more robust than model-free methods when the reward structure changes (but still fallible when  $\mathcal{T}$  changes). Work extending SRs to deep neural networks has demonstrated its usefulness in multi-task settings, whilst within a complex visual environment [66].

### B. Exploration vs. Exploitation

One of the greatest difficulties in RL is the fundamental dilemma of *exploration versus exploitation*: When should the agent try out (perceived) non-optimal actions in order to explore the environment (and potentially improve the model), and when should it exploit the optimal action in order to make useful progress? Off-policy algorithms, such as the DQN [84], typically use the simple  $\epsilon$ -greedy exploration policy, which chooses a random action with probability  $\epsilon \in [0, 1]$ , and the optimal action otherwise. By decreasing  $\epsilon$  over time, the agent progresses towards exploitation. Although adding independent noise for exploration is usable in continuous control problems, more sophisticated strategies inject noise that is correlated over time (e.g., from stochastic processes) in order to better preserve momentum [79].

The observation that temporal correlation is important led Osband et al. [97] to propose the bootstrapped DQN, which maintains several  $Q$ -value “heads” that learn different values through a combination of different weight initialisations and bootstrapped sampling from experience replay memory. At the beginning of each training episode, a different head is chosen, leading to temporally-extended exploration. Usunier et al. [147] later proposed a similar method that performed exploration in policy space by adding noise to a single output head, using zero-order gradient estimates to allow backpropo-

gation through the policy.

One of the main principled exploration strategies is the *upper confidence bound* (UCB) algorithm, based on the principle of “optimism in the face of uncertainty” [67]. The idea behind UCB is to pick actions that maximise  $\mathbb{E}[R] + \kappa\sigma[R]$ , where  $\sigma[R]$  is the standard deviation of the return and  $\kappa > 0$ . UCB therefore encourages exploration in regions with high uncertainty and moderate expected return. Whilst easily achievable in small tabular cases, the use of powerful density models [11], or conversely, hashing [139], has allowed this algorithm to scale to high-dimensional visual domains with DRL. UCB is only one technique for trading off exploration and exploitation in the context of Bayesian optimisation [126]; future work in DRL may benefit from investigating other successful techniques that are used in Bayesian optimisation.

UCB can also be considered one way of implementing *intrinsic motivation*, which is a general concept that advocates decreasing uncertainty/making progress in learning about the environment [119]. There have been several DRL algorithms that try to implement intrinsic motivation via minimising model prediction error [132, 101] or maximising information gain [86, 52].

### C. Hierarchical RL

In the same way that deep learning relies on hierarchies of features, HRL relies on hierarchies of policies. Early work in this area introduced *options*, in which, apart from *primitive actions* (single-timestep actions), policies could also run other policies (multi-timestep “actions”) [136]. This approach allows top-level policies to focus on higher-level *goals*, whilst *subpolicies* are responsible for fine control. Several works in DRL have attempted HRL by using one top-level policy that chooses between subpolicies, where the division of states or goals in to subpolicies is achieved either manually [2, 143, 65] or automatically [3, 151, 152]. One way to help construct subpolicies is to focus on discovering and reaching goals, which are specific states in the environment; they may often be locations, which an agent should navigate to. Whether utilised with HRL or not, the discovery and generalisation of goals is also an important area of ongoing research [117, 66, 152].

### D. Imitation Learning and Inverse RL

One may ask why, if given a sequence of “optimal” actions from expert demonstrations, it is not possible to use supervised learning in a straightforward manner—a case of “learning from demonstration”. This is indeed possible, and is known as *behavioural cloning* in traditional RL literature. Taking advantage of the stronger signals available in supervised learning problems, behavioural cloning enjoyed success in earlier neural network research, with the most notable success being ALVINN, one of the earliest autonomous cars [104]. However, behavioural cloning cannot adapt to new situations, and small deviations from the demonstration during the execution of the learned policy can compound and lead to scenarios where the policy is unable to recover. A more generalisable solution is to use provided trajectories to guide the learning of suitable state-action pairs, but fine-tune the agent using RL [49]. Alternatively, if the expert is still available to query during training, the agent can use active learning to gather extra data

when it is unsure, allowing it to learn from states away from the optimal trajectories [110]. This has been applied to a deep learning setting, where a CNN trained in a visual navigation task with active learning significantly improved upon a pure imitation learning baseline [53].

The goal of IRL is to estimate an unknown reward function from observed trajectories that characterise a desired solution [92]; IRL can be used in combination with RL to improve upon demonstrated behaviour. Using the power of deep neural networks, it is now possible to learn complex, nonlinear reward functions for IRL [165]. Ho and Ermon [51] showed that policies are uniquely characterised by their *occupancies* (visited state and action distributions) allowing IRL to be reduced to the problem of measure matching. With this insight, they were able to use generative adversarial training [38] to facilitate reward function learning in a more flexible manner, resulting in the generative adversarial imitation learning (GAIL) algorithm. GAIL was later extended to allow IRL to be applied even when receiving expert trajectories from a different visual viewpoint to that of the RL agent [131]. In complementary work, Baram et al. [7] exploit gradient information that was not used in GAIL to learn models within the IRL process.

### E. Multi-agent RL

Usually, RL considers a single learning agent in a stationary environment. In contrast, multi-agent RL (MARL) considers multiple agents learning through RL, and often the non-stationarity introduced by other agents changing their behaviours as they learn [18]. In DRL, the focus has been on enabling (differentiable) communication between agents, which allows them to co-operate. Several approaches have been proposed for this purpose, including passing messages to agents sequentially [33], using a bidirectional channel (providing ordering with less signal loss) [102], and an all-to-all channel [134]. The addition of communication channels is a natural strategy to apply to MARL in complex scenarios and does not preclude the usual practice of modelling co-operative or competing agents as applied elsewhere in the MARL literature [18]. Other DRL works of note in MARL investigate the effects of learning and sequential decision making in game theory [48, 71].

### F. Memory and Attention

As one of the earliest works in DRL the DQN spawned many extensions. One of the first extensions was converting the DQN into an RNN, which allows the network to better deal with POMDPs by integrating information over long time periods. Like recursive filters, recurrent connections provide an efficient means of acting conditionally on temporally distant prior observations. By using recurrent connections between its hidden units, the deep recurrent *Q*-network (DRQN) introduced by Hausknecht and Stone [44] was able to successfully infer the velocity of the ball in the game “Pong,” even when frames of the game were randomly blanked out. Further improvements were gained by introducing *attention*—a technique where additional connections are added from the recurrent units to lower layers—to the DRQN, resulting in the deep attention recurrent *Q*-network (DARQN) [130]. Attention gives a network the ability to choose which part of its next

input to focus on, and allowed the DARQN to beat both the DQN and DRQN on games, which require longer-term planning. However, the DQN outperformed the DRQN and DARQN on games requiring quick reactions, where  $Q$ -values can fluctuate more rapidly.

Taking recurrent processing further, it is possible to add a differentiable memory to the DQN, which allows it to more flexibly process information in its “working memory” [96]. In traditional RNNs, recurrent units are responsible for both performing calculations and storing information. Differentiable memories add large matrices that are purely used for storing information, and can be accessed using differentiable read and write operations, analogously to computer memory. With their key-value-based memory  $Q$ -network (MQN), Oh et al. [96] constructed an agent that could solve a simple maze built in Minecraft, where the correct goal in each episode was indicated by a coloured block shown near the start of the maze. The MQN, and especially its more sophisticated variants, significantly outperformed both DQN and DRQN baselines, highlighting the importance of using decoupled memory storage. More recent work, where the memory was given a 2D structure in order to resemble a spatial map, hints at future research where more specialised memory structures will be developed to address specific problems, such as 2D or 3D navigation [98]. Alternatively, differentiable memories can be used as approximate hash tables, allowing DRL algorithms to store and retrieve successful experiences to facilitate rapid learning [105].

Note that RNNs are not restricted to value-function-based methods but have also been successfully applied to policy search [163] and actor-critic methods [45, 85].

#### G. Transfer Learning

Even though DRL algorithms can process high-dimensional inputs, it is rarely feasible to train RL agents directly on visual inputs in the real world, due to the large number of samples required. To speed up learning in DRL, it is possible to exploit previously acquired knowledge from related tasks, which comes in several guises: transfer learning, multitask learning [20] and curriculum learning [14] to name a few. There is much interest in transferring learning from one task to another, particularly from training in physics simulators with visual renderers and fine-tuning the models in the real world. This can be achieved in a naive fashion, directly using the same network in both the simulated and real phases [167], or with more sophisticated training procedures that directly try to mitigate the problem of neural networks “catastrophically forgetting” old knowledge by adding extra layers when transferring domain [114, 115]. Other approaches involve directly learning an alignment between simulated and real visuals [146], or even between two different camera viewpoints [131].

A different form of transfer can be utilised to help RL in the form of multitask training [77, 54, 82]. Especially with neural networks, supervised and unsupervised learning tasks can help train features that can be used by RL agents, making optimising the RL objective easier to achieve. For example, the “unsupervised reinforcement and auxiliary learning” A3C-based agent is additionally trained with “pixel control” (maxi-

mally changing pixel inputs), plus reward prediction and value function learning from experience replay [54]. Meanwhile, the A3C-based agent of Mirowski et al. [82] was additionally trained to construct a depth map given RGB inputs, which helps it in its task of learning to navigate a 3D environment. In an ablation study, Mirowski et al. [82] showed the predicting depth was more useful than receiving depth as an extra input, lending further support to the idea that gradients induced by auxiliary tasks can be extremely effective at boosting DRL.

Transfer learning can also be used to construct more data- and parameter-efficient policies. In the student-teacher paradigm in machine learning, one can first train a more powerful “teacher” model, and then use it to guide the training of a less powerful “student” model. Whilst originally applied to supervised learning, the neural network knowledge transfer technique known as *distillation* [50] has been utilised to both transfer policies learned by large DQNs to smaller DQNs, and transfer policies learned by several DQNs trained on separate games to one single DQN [99, 113]. Together, the combination of multitask and transfer learning can improve the sample efficiency and robustness of current DRL algorithms [140]. These are important topics if we wish to construct agents that can accomplish a wide range of tasks, since naively training on multiple RL objectives at once may be infeasible.

#### H. Benchmarks

One of the challenges in any field in machine learning is developing a standardised way to evaluate new techniques. Although much early work focused on simple, custom MDPs, there shortly emerged control problems that could be used as standard benchmarks for testing new algorithms, such as the Cartpole [8] and Mountain Car [87] domains.

However, these problems were limited to relatively small state spaces, and therefore failed to capture the complexities that would be encountered in most realistic scenarios. Arguably the initial driver of DRL, the ALE provided an interface to Atari 2600 video games, with code to access over 50 games provided with the initial release [10]. As video games can vary greatly, but still present interesting and challenging objectives for humans, they provide an excellent testbed for RL agents. As the first algorithm to successfully play a range of these games directly from their visuals, the DQN [84] has secured its place as a milestone in the development of RL algorithms. This success story has started a trend of using video games as standardised RL testbeds, with several interesting options now available. ViZDoom provides an interface to the Doom first-person shooter [60], and echoing the popularity of e-sports competitions, ViZDoom competitions are now held at the yearly IEEE Conference on Computational Intelligence in Games. Facebook’s TorchCraft [137] and DeepMind’s StarCraft II Learning Environment [153] respectively provide interfaces to the StarCraft and StarCraft II real-time strategy games, presenting challenges in both micromanagement and long-term planning. In an aim to provide more flexible environments, DeepMind Lab was developed on top of the Quake III Arena first-person shooter engine [9], and Microsoft’s Project Malmo exposed an interface to the Minecraft sandbox game [55]. Both environments provide customisable platforms

for RL agents in 3D environments.

Most DRL approaches focus on discrete actions, but some solutions have also been developed for continuous control problems. Many DRL papers in continuous control [122, 46, 79, 85, 7, 131] have used the MuJoCo physics engine to obtain relatively realistic dynamics for multi-joint continuous control problems [144], and there has now been some effort to standardise these problems [28].

To help with standardisation and reproducibility, most of the aforementioned RL domains and more have been made available in the OpenAI Gym, a library and online service that allows people to easily interface with and publicly share the results of RL algorithms on these domains [17].

## VII. CONCLUSION: BEYOND PATTERN RECOGNITION

Despite the successes of DRL, many problems need to be addressed before these techniques can be applied to a wide range of complex real-world problems [68]. Recent work with (non-deep) generative causal models demonstrated superior generalisation over standard DRL algorithms [85, 114] in some benchmarks [10], achieved by reasoning about causes and effects in the environment [58]. For example, the schema networks of Kanksy et al. [58] trained on the game “Breakout” immediately adapted to a variant where a small wall was placed in front of the target blocks, whilst progressive (A3C) networks [114] failed to match the performance of the schema networks even after training on the new domain. Although DRL has already been combined with AI techniques, such as search [128] and planning [138], a deeper integration with other traditional AI approaches promises benefits such as better sample complexity, generalisation and interpretability [35]. In time, we also hope that our theoretical understanding of the properties of neural networks (particularly within DRL) will improve, as it currently lags far behind practice.

To conclude, it is worth revisiting the overarching goal of all of this research: the creation of general-purpose AI systems that can interact with and learn from the world around them. Interaction with the environment is simultaneously the advantage and disadvantage of RL. Whilst there are many challenges in seeking to understand our complex and ever-changing world, RL allows us to choose how we explore it. In effect, RL endows agents with the ability to perform experiments to better understand their surroundings, enabling them to learn even high-level causal relationships. The availability of high-quality visual renderers and physics engines now enables us to take steps in this direction, with works that try to learn intuitive models of physics in visual environments [27]. Challenges remain before this will be possible in the real world, but steady progress is being made in agents that learn the fundamental principles of the world through observation and action. Perhaps, then, we are not too far away from AI systems that learn and act in more human-like ways in increasingly complex environments.

## ACKNOWLEDGMENTS

The authors would like to thank the reviewers and broader community for their feedback on this survey; in particular, we would like to thank Nicolas Heess for clarifications on several points. Kai Arulkumaran would like to acknowledge

PhD funding from the Department of Bioengineering, Imperial College London. This research has been partially funded by a Google Faculty Research Award to Marc Deisenroth.

## REFERENCES

- [1] Pieter Abbeel and John Schulman. Deep Reinforcement Learning through Policy Optimization, 2016. Tutorial at NIPS 2016.
- [2] Kai Arulkumaran, Nat Dilokthanakul, Murray Shanahan, and Anil Anthony Bharath. Classifying Options for Deep Reinforcement Learning. In *IJCAI Workshop on Deep Reinforcement Learning: Frontiers and Challenges*, 2016.
- [3] Pierre-Luc Bacon, Jean Harb, and Doina Precup. The Option-Critic Architecture. In *AAAI*, 2017.
- [4] J Andrew Bagnell and Jeff Schneider. Covariant Policy Search. In *IJCAI*, 2003.
- [5] Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron Courville, and Yoshua Bengio. An Actor-Critic Algorithm for Sequence Prediction. In *ICLR*, 2017.
- [6] Leemon C Baird III. Advantage Updating. Technical report, DTIC, 1993.
- [7] Nir Baram, Oron Anschel, and Shie Mannor. Model-Based Adversarial Imitation Learning. In *NIPS Workshop on Deep Reinforcement Learning*, 2016.
- [8] Andrew G Barto, Richard S Sutton, and Charles W Anderson. Neuronlike Adaptive Elements That Can Solve Difficult Learning Control Problems. *IEEE Trans. on Systems, Man, and Cybernetics*, (5):834–846, 1983.
- [9] Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler, Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. DeepMind Lab. *arXiv:1612.03801*, 2016.
- [10] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The Arcade Learning Environment: An Evaluation Platform for General Agents. In *IJCAI*, 2015.
- [11] Marc G Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi Munos. Unifying Count-Based Exploration and Intrinsic Motivation. In *NIPS*, 2016.
- [12] Marc G Bellemare, Will Dabney, and Rémi Munos. A Distributional Perspective on Reinforcement Learning. In *ICML*, 2017.
- [13] Richard Bellman. On the Theory of Dynamic Programming. *PNAS*, 38(8):716–719, 1952.
- [14] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum Learning. In *ICML*, 2009.
- [15] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation Learning: A Review and New Perspectives. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013.
- [16] Dimitri P Bertsekas. Dynamic Programming and Suboptimal Control: A Survey from ADP to MPC. *European Journal of Control*, 11(4-5): 310–334, 2005.
- [17] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym. *arXiv:1606.01540*, 2016.
- [18] Lucian Busoniu, Robert Babuska, and Bart De Schutter. A Comprehensive Survey of Multiagent Reinforcement Learning. *IEEE Trans. on Systems, Man, And Cybernetics*, 2008.
- [19] Murray Campbell, A Joseph Hoane, and Feng-hsiung Hsu. Deep Blue. *Artificial Intelligence*, 134(1-2):57–83, 2002.
- [20] Rich Caruana. Multitask Learning. *Machine Learning*, 28(1):41–75, 1997.
- [21] Silvia Chiappa, Sébastien Racaniere, Daan Wierstra, and Shakir Mohamed. Recurrent Environment Simulators. In *ICLR*, 2017.
- [22] Paul Christiano, Zain Shah, Igor Mordatch, Jonas Schneider, Trevor Blackwell, Joshua Tobin, Pieter Abbeel, and Wojciech Zaremba. Transfer from Simulation to Real World through Learning Deep Inverse Dynamics Model. *arXiv:1610.03518*, 2016.
- [23] Giuseppe Cuccu, Matthew Luciw, Jürgen Schmidhuber, and Faustino Gomez. Intrinsically Motivated Neuroevolution for Vision-Based Reinforcement Learning. In *ICDL*, volume 2, 2011.
- [24] Peter Dayan. Improving Generalization for Temporal Difference Learning: The Successor Representation. *Neural Computation*, 5(4): 613–624, 1993.
- [25] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Andrew Senior, Paul Tucker, Ke Yang, Quoc V Le, et al. Large Scale Distributed Deep Networks. In *NIPS*, 2012.
- [26] Marc P Deisenroth, Gerhard Neumann, and Jan Peters. A Survey on Policy Search for Robotics. *Foundations and Trends® in Robotics*, 2 (1–2), 2013.

- [27] Misha Denil, Pulkit Agrawal, Tejas D Kulkarni, Tom Erez, Peter Battaglia, and Nando de Freitas. Learning to Perform Physics Experiments via Deep Reinforcement Learning. In *ICLR*, 2017.
- [28] Yan Duan, Xi Chen, Rein Houthooft, John Schulman, and Pieter Abbeel. Benchmarking Deep Reinforcement Learning for Continuous Control. In *ICML*, 2016.
- [29] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RL<sup>2</sup>: Fast Reinforcement Learning via Slow Reinforcement Learning. In *NIPS Workshop on Deep Reinforcement Learning*, 2016.
- [30] Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt, Peter Sunehag, Timothy Lillicrap, Jonathan Hunt, Timothy Mann, Theophane Weber, Thomas Degris, and Ben Coppin. Deep Reinforcement Learning in Large Discrete Action Spaces. *arXiv:1512.07679*, 2015.
- [31] David Ferrucci, Eric Brown, Jennifer Chu-Carroll, James Fan, David Gondek, Aditya A Kalyanpur, Adam Lally, J William Murdock, Eric Nyberg, John Prager, et al. Building Watson: An Overview of the DeepQA Project. *AI Magazine*, 31(3):59–79, 2010.
- [32] Chelsea Finn, Xin Yu Tan, Yan Duan, Trevor Darrell, Sergey Levine, and Pieter Abbeel. Deep Spatial Autoencoders for Visuomotor Learning. In *ICRA*, 2016.
- [33] Jakob Foerster, Yannis M Assael, Nando de Freitas, and Shimon Whiteson. Learning to Communicate with Deep Multi-Agent Reinforcement Learning. In *NIPS*, 2016.
- [34] Michael C Fu. Gradient Estimation. *Handbooks in Operations Research and Management Science*, 13:575–616, 2006.
- [35] Marta Garnelo, Kai Arulkumaran, and Murray Shanahan. Towards Deep Symbolic Reinforcement Learning. In *NIPS Workshop on Deep Reinforcement Learning*, 2016.
- [36] Peter W Glynn. Likelihood Ratio Gradient Estimation for Stochastic Systems. *Communications of the ACM*, 33(10):75–84, 1990.
- [37] Faustino Gomez and Jürgen Schmidhuber. Evolving Modular Fast-Weight Networks for Control. In *ICANN*, 2005.
- [38] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. In *NIPS*, 2014.
- [39] Audrunas Gruslys, Mohammad Gheshlaghi Azar, Marc G Bellemare, and Rémi Munos. The Reactor: A Sample-Efficient Actor-Critic Architecture. *arXiv:1704.04651*, 2017.
- [40] Shixiang Gu, Timothy Lillicrap, Ilya Sutskever, and Sergey Levine. Continuous Deep Q-Learning with Model-Based Acceleration. In *ICLR*, 2016.
- [41] Shixiang Gu, Timothy Lillicrap, Zoubin Ghahramani, Richard E Turner, and Sergey Levine. Q-Prop: Sample-Efficient Policy Gradient with an Off-Policy Critic. In *ICLR*, 2017.
- [42] Shixiang Gu, Timothy Lillicrap, Zoubin Ghahramani, Richard E Turner, Bernhard Schölkopf, and Sergey Levine. Interpolated Policy Gradient: Merging On-Policy and Off-Policy Gradient Estimation for Deep Reinforcement Learning. In *NIPS*, 2017.
- [43] Mance E Harmon and Leemon C Baird III. Multi-Player Residual Advantage Learning with General Function Approximation. Technical report, DTIC, 1996.
- [44] Matthew Hausknecht and Peter Stone. Deep Recurrent Q-Learning for Partially Observable MDPs. In *AAAI Fall Symposium Series*, 2015.
- [45] Nicolas Heess, Jonathan J Hunt, Timothy P Lillicrap, and David Silver. Memory-Based Control with Recurrent Neural Networks. In *NIPS Workshop on Deep Reinforcement Learning*, 2015.
- [46] Nicolas Heess, Gregory Wayne, David Silver, Tim Lillicrap, Tom Erez, and Yuval Tassa. Learning Continuous Control Policies by Stochastic Value Gradients. In *NIPS*, 2015.
- [47] Nicolas Heess, Srinivasan Sriram, Jay Lemmon, Josh Merel, Greg Wayne, Yuval Tassa, Tom Erez, Ziyu Wang, Ali Eslami, Martin Riedmiller, et al. Emergence of Locomotion Behaviours in Rich Environments. *arXiv:1707.02286*, 2017.
- [48] Johannes Heinrich and David Silver. Deep Reinforcement Learning from Self-Play in Imperfect-Information Games. 2016.
- [49] Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Andrew Sendonaris, Gabriel Dulac-Arnold, Ian Osband, John Agapiou, et al. Learning from Demonstrations for Real World Reinforcement Learning. *arXiv:1704.03732*, 2017.
- [50] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the Knowledge in a Neural Network. 2014.
- [51] Jonathan Ho and Stefano Ermon. Generative Adversarial Imitation Learning. In *NIPS*, 2016.
- [52] Rein Houthooft, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. VIME: Variational Information Maximizing Exploration. In *NIPS*, 2016.
- [53] Ahmed Hussein, Mohamed Medhat Gaber, and Eyad Elyan. Deep Active Learning for Autonomous Navigation. In *EANN*, 2016.
- [54] Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David Silver, and Koray Kavukcuoglu. Reinforcement Learning with Unsupervised Auxiliary Tasks. In *ICLR*, 2017.
- [55] Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. The Malmo Platform for Artificial Intelligence Experimentation. In *IJCAI*, 2016.
- [56] Leslie P Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and Acting in Partially Observable Stochastic Domains. *Artificial Intelligence*, 101(1):99–134, 1998.
- [57] Sham M Kakade. A Natural Policy Gradient. In *NIPS*, 2002.
- [58] Ken Kansky, Tom Silver, David A Mély, Mohamed Eldawy, Miguel Lázaro-Gredilla, Xinghua Lou, Nimrod Dorfman, Szymon Sidor, Scott Phoenix, and Dileep George. Schema Networks: Zero-Shot Transfer with a Generative Causal Model of Intuitive Physics. In *ICML*, 2017.
- [59] Hilbert J Kappen. Path Integrals and Symmetry Breaking for Optimal Control Theory. *JSTAT*, 2005(11):P11011, 2005.
- [60] Michał Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczek, and Wojciech Jaskowski. ViZDoom: A Doom-Based AI Research Platform for Visual Reinforcement Learning. In *CIG*, 2016.
- [61] Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. In *ICLR*, 2014.
- [62] Nate Kohl and Peter Stone. Policy Gradient Reinforcement Learning for Fast Quadrupedal Locomotion. In *ICRA*, volume 3, 2004.
- [63] Vijay R Konda and John N Tsitsiklis. On Actor-Critic Algorithms. *SICON*, 42(4):1143–1166, 2003.
- [64] Jan Koutník, Giuseppe Cuccu, Jürgen Schmidhuber, and Faustino Gomez. Evolving Large-Scale Neural Networks for Vision-Based Reinforcement Learning. In *GECCO*, 2013.
- [65] Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation. In *NIPS*, 2016.
- [66] Tejas D Kulkarni, Ardavan Saeedi, Simanta Gautam, and Samuel J Gershman. Deep Successor Reinforcement Learning. In *NIPS Workshop on Deep Reinforcement Learning*, 2016.
- [67] Tze Leung Lai and Herbert Robbins. Asymptotically Efficient Adaptive Allocation Rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985.
- [68] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building Machines That Learn and Think Like People. *The Behavioral and Brain Sciences*, page 1, 2016.
- [69] Sascha Lange, Martin Riedmiller, and Arne Voigtlander. Autonomous Reinforcement Learning on Raw Visual Input Data in a Real World Application. In *IJCNN*, 2012.
- [70] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep Learning. *Nature*, 521(7553):436–444, 2015.
- [71] Joel Z Leibo, Vinicius Zambaldi, Marc Lanctot, Janusz Marecki, and Thore Graepel. Multi-Agent Reinforcement Learning in Sequential Social Dilemmas. In *AAMAS*, 2017.
- [72] Sergey Levine and Pieter Abbeel. Learning Neural Network Policies with Guided Policy Search under Unknown Dynamics. In *NIPS*, 2014.
- [73] Sergey Levine and Vladlen Koltun. Guided Policy Search. In *ICLR*, 2013.
- [74] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-End Training of Deep Visuomotor Policies. *JMLR*, 17(39):1–40, 2016.
- [75] Sergey Levine, Peter Pastor, Alex Krizhevsky, and Deirdre Quillen. Learning Hand-Eye Coordination for Robotic Grasping with Deep Learning and Large-Scale Data Collection. In *ISER*, 2016.
- [76] Ke Li and Jitendra Malik. Learning to Optimize. 2017.
- [77] Xiuju Li, Lihong Li, Jianfeng Gao, Xiaodong He, Jianshu Chen, Li Deng, and Ji He. Recurrent Reinforcement Learning: A Hybrid Approach. *arXiv:1509.03044*, 2015.
- [78] Yuxi Li. Deep Reinforcement Learning: An Overview. *arXiv:1701.07274*, 2017.
- [79] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous Control with Deep Reinforcement Learning. In *ICLR*, 2016.
- [80] Long-Ji Lin. Self-Improving Reactive Agents Based on Reinforcement Learning, Planning and Teaching. *Machine Learning*, 8(3–4):293–321, 1992.
- [81] Luke Metz, Julian Ibarz, Navdeep Jaithly, and James Davidson. Discrete Sequential Prediction of Continuous Actions for Deep RL. *arXiv:1705.05035*, 2017.
- [82] Piotr Mirowski, Razvan Pascanu, Fabio Viola, Hubert Soyer, Andy Ballard, Andrea Banino, Misha Denil, Ross Goroshin, Laurent Sifre, Koray

- Kavukcuoglu, et al. Learning to Navigate in Complex Environments. In *ICLR*, 2017.
- [83] Volodymyr Mnih, Nicolas Heess, Alex Graves, and Koray Kavukcuoglu. Recurrent Models of Visual Attention. In *NIPS*, 2014.
- [84] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-Level Control through Deep Reinforcement Learning. *Nature*, 518(7540):529–533, 2015.
- [85] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy P Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous Methods for Deep Reinforcement Learning. In *ICLR*, 2016.
- [86] Shakir Mohamed and Danilo Jimenez Rezende. Variational Information Maximisation for Intrinsically Motivated Reinforcement Learning. In *NIPS*, 2015.
- [87] Andrew William Moore. Efficient Memory-Based Learning for Robot Control. Technical report, University of Cambridge, Computer Laboratory, 1990.
- [88] Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc G Bellemare. Safe and Efficient Off-Policy Reinforcement Learning. In *NIPS*, 2016.
- [89] Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and Dale Schuurmans. Bridging the Gap Between Value and Policy Based Reinforcement Learning. *arXiv:1702.08892*, 2017.
- [90] Anusha Nagabandi, Gregory Kahn, Ronald S Fearing, and Sergey Levine. Neural Network Dynamics for Model-Based Deep Reinforcement Learning with Model-Free Fine-Tuning. *arXiv:1708.02596*, 2017.
- [91] Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro De Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, et al. Massively Parallel Methods for Deep Reinforcement Learning. In *ICML Workshop on Deep Learning*, 2015.
- [92] Andrew Y Ng and Stuart J Russell. Algorithms for Inverse Reinforcement Learning. In *ICML*, 2000.
- [93] Andrew Y Ng, Adam Coates, Mark Diel, Varun Ganapathi, Jamie Schulte, Ben Tse, Eric Berger, and Eric Liang. Autonomous Inverted Helicopter Flight via Reinforcement Learning. *Experimental Robotics*, pages 363–372, 2006.
- [94] Brendan O'Donoghue, Rémi Munos, Koray Kavukcuoglu, and Volodymyr Mnih. PGQ: Combining Policy Gradient and Q-Learning. In *ICLR*, 2017.
- [95] Junhyuk Oh, Xiaoxiao Guo, Honglak Lee, Richard L Lewis, and Satinder Singh. Action-Conditional Video Prediction using Deep Networks in Atari Games. In *NIPS*, 2015.
- [96] Junhyuk Oh, Valliappa Chockalingam, Satinder Singh, and Honglak Lee. Control of Memory, Active Perception, and Action in Minecraft. In *ICLR*, 2016.
- [97] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep Exploration via Bootstrapped DQN. In *NIPS*, 2016.
- [98] Emilio Parisotto and Ruslan Salakhutdinov. Neural Map: Structured Memory for Deep Reinforcement Learning. *arXiv:1702.08360*, 2017.
- [99] Emilio Parisotto, Jimmy L Ba, and Ruslan Salakhutdinov. Actor-Mimic: Deep Multitask and Transfer Reinforcement Learning. In *ICLR*, 2016.
- [100] Razvan Pascanu, Yujia Li, Oriol Vinyals, Nicolas Heess, Lars Buesing, Sébastien Racanière, David Reichert, Théophane Weber, Daan Wierstra, and Peter Battaglia. Learning Model-Based Planning from Scratch. *arXiv:1707.06170*, 2017.
- [101] Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-Driven Exploration by Self-supervised Prediction. In *ICML*, 2017.
- [102] Peng Peng, Ying Wen, Yaodong Yang, Quan Yuan, Zhenkun Tang, Haitao Long, and Jun Wang. Multiagent Bidirectionally-Coordinated Nets: Emergence of Human-level Coordination in Learning to Play StarCraft Combat Games. *arXiv:1703.10069*, 2017.
- [103] Jan Peters, Katharina Mülling, and Yasemin Altun. Relative Entropy Policy Search. In *AAAI*, 2010.
- [104] Dean A Pomerleau. ALVINN, an Autonomous Land Vehicle in a Neural Network. Technical report, Carnegie Mellon University, Computer Science Department, 1989.
- [105] Alexander Pritzel, Benigno Uria, Sriram Srinivasan, Adrià Puigdomènech, Oriol Vinyals, Denis Hassabis, Daan Wierstra, and Charles Blundell. Neural Episodic Control. In *ICML*, 2017.
- [106] Marc'Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. Sequence Level Training with Recurrent Neural Networks. In *ICLR*, 2016.
- [107] Benjamin Recht, Christopher Re, Stephen Wright, and Feng Niu. Hogwild: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent. In *NIPS*, 2011.
- [108] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic Backpropagation and Approximate Inference in Deep Generative Models. In *ICML*, 2014.
- [109] Martin Riedmiller. Neural Fitted Q Iteration—First Experiences with a Data Efficient Neural Reinforcement Learning Method. In *ECML*, 2005.
- [110] Stéphane Ross, Geoffrey J Gordon, and Drew Bagnell. A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning. In *AISTATS*, 2011.
- [111] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning Representations by Back-Propagating Errors. *Cognitive Modeling*, 5(3):1, 1988.
- [112] Gavin A Rummery and Mahesan Niranjan. *On-line Q-learning using Connectionist Systems*. University of Cambridge, Department of Engineering, 1994.
- [113] Andrei A Rusu, Sergio Gomez Colmenarejo, Caglar Gulcehre, Guillaume Desjardins, James Kirkpatrick, Razvan Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and Raia Hadsell. Policy Distillation. In *ICLR*, 2016.
- [114] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive Neural Networks. *arXiv:1606.04671*, 2016.
- [115] Andrei A Rusu, Matej Vecerik, Thomas Rothörl, Nicolas Heess, Razvan Pascanu, and Raia Hadsell. Sim-to-Real Robot Learning from Pixels with Progressive Nets. In *CoRL*, 2017.
- [116] Tim Salimans, Jonathan Ho, Xi Chen, and Ilya Sutskever. Evolution Strategies as a Scalable Alternative to Reinforcement Learning. *arXiv:1703.03864*, 2017.
- [117] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal Value Function Approximators. In *ICML*, 2015.
- [118] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized Experience Replay. In *ICLR*, 2016.
- [119] Jürgen Schmidhuber. A Possibility for Implementing Curiosity and Boredom in Model-Building Neural Controllers. In *SAB*, 1991.
- [120] Jürgen Schmidhuber and Rudolf Huber. Learning to Generate Artificial Fovea Trajectories for Target Detection. *IJNS*, 2(01n02):125–134, 1991.
- [121] John Schulman, Nicolas Heess, Theophane Weber, and Pieter Abbeel. Gradient Estimation using Stochastic Computation Graphs. In *NIPS*, 2015.
- [122] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust Region Policy Optimization. In *ICML*, 2015.
- [123] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-Dimensional Continuous Control using Generalized Advantage Estimation. In *ICLR*, 2016.
- [124] John Schulman, Pieter Abbeel, and Xi Chen. Equivalence Between Policy Gradients and Soft Q-Learning. *arXiv:1704.06440*, 2017.
- [125] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms. *arXiv:1707.06347*, 2017.
- [126] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando de Freitas. Taking the Human out of the Loop: A Review of Bayesian Optimization. *Proc. of the IEEE*, 104(1):148–175, 2016.
- [127] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic Policy Gradient Algorithms. In *ICML*, 2014.
- [128] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the Game of Go with Deep Neural Networks and Tree Search. *Nature*, 529(7587):484–489, 2016.
- [129] Satinder Singh, Diane Litman, Michael Kearns, and Marilyn Walker. Optimizing Dialogue Management with Reinforcement Learning: Experiments with the NJFun System. *JAIR*, 16:105–133, 2002.
- [130] Ivan Sorokin, Alexey Seleznev, Mikhail Pavlov, Aleksandr Fedorov, and Anastasiia Ignateva. Deep Attention Recurrent Q-Network. In *NIPS Workshop on Deep Reinforcement Learning*, 2015.
- [131] Bradley C Stadie, Pieter Abbeel, and Ilya Sutskever. Third Person Imitation Learning. In *ICLR*, 2017.
- [132] Bradley C Stadie, Sergey Levine, and Pieter Abbeel. Incentivizing Exploration in Reinforcement Learning with Deep Predictive Models. In *NIPS Workshop on Deep Reinforcement Learning*, 2015.
- [133] Alexander L Strehl, Lihong Li, Eric Wiewiora, John Langford, and

- Michael L Littman. PAC Model-Free Reinforcement Learning. In *ICML*, 2006.
- [134] Sainbayar Sukhbaatar, Arthur Szlam, and Rob Fergus. Learning Multiagent Communication with Backpropagation. In *NIPS*, 2016.
- [135] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- [136] Richard S Sutton, Doina Precup, and Satinder Singh. Between MDPs and Semi-MDPs: A Framework for Temporal Abstraction in Reinforcement Learning. *Artificial Intelligence*, 112(1–2):181–211, 1999.
- [137] Gabriel Synnaeve, Nantas Nardelli, Alex Auvolat, Soumith Chintala, Timothée Lacroix, Zeming Lin, Florian Richoux, and Nicolas Usunier. TorchCraft: A Library for Machine Learning Research on Real-Time Strategy Games. *arXiv:1611.00625*, 2016.
- [138] Aviv Tamar, Yi Wu, Garrett Thomas, Sergey Levine, and Pieter Abbeel. Value Iteration Networks. In *NIPS*, 2016.
- [139] Haoran Tang, Rein Houthooft, Davis Foote, Adam Stooke, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. #Exploration: A Study of Count-Based Exploration for Deep Reinforcement Learning. In *NIPS*, 2017.
- [140] Yee Whye Teh, Victor Bapst, Wojciech Marian Czarnecki, John Quan, James Kirkpatrick, Raia Hadsell, Nicolas Heess, and Razvan Pascanu. Distral: Robust Multitask Reinforcement Learning. In *NIPS*, 2017.
- [141] Gerald Tesauro. Temporal Difference Learning and TD-Gammon. *Communications of the ACM*, 38(3):58–68, 1995.
- [142] Gerald Tesauro, Rajarshi Das, Hoi Chan, Jeffrey Kephart, David Levine, Freeman Rawson, and Charles Lefurgy. Managing Power Consumption and Performance of Computing Systems using Reinforcement Learning. In *NIPS*, 2008.
- [143] Chen Tessler, Shahar Givony, Tom Zahavy, Daniel J Mankowitz, and Shie Mannor. A Deep Hierarchical Approach to Lifelong Learning in Minecraft. In *AAAI*, 2017.
- [144] Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A Physics Engine for Model-Based Control. In *IROS*, 2012.
- [145] John N Tsitsiklis and Benjamin Van Roy. Analysis of Temporal-Difference Learning with Function Approximation. In *NIPS*, 1997.
- [146] Eric Tzeng, Coline Devin, Judy Hoffman, Chelsea Finn, Xingchao Peng, Sergey Levine, Kate Saenko, and Trevor Darrell. Towards Adapting Deep Visuomotor Representations from Simulated to Real Environments. In *WAFR*, 2016.
- [147] Nicolas Usunier, Gabriel Synnaeve, Zeming Lin, and Soumith Chintala. Episodic Exploration for Deep Deterministic Policies: An Application to StarCraft Micromanagement Tasks. In *ICLR*, 2017.
- [148] Hado van Hasselt. Double Q-Learning. In *NIPS*, 2010.
- [149] Hado van Hasselt, Arthur Guez, and David Silver. Deep Reinforcement Learning with Double Q-Learning. In *AAAI*, 2016.
- [150] Harm Vanseijen and Rich Sutton. A Deeper Look at Planning as Learning from Replay. In *ICML*, 2015.
- [151] Alexander Vezhnevets, Volodymyr Mnih, Simon Osindero, Alex Graves, Oriol Vinyals, John Agapiou, and Koray Kavukcuoglu. Strategic Attentive Writer for Learning Macro-Actions. In *NIPS*, 2016.
- [152] Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. FeUdal Networks for Hierarchical Reinforcement Learning. In *ICML*, 2017.
- [153] Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schriftwieser, et al. StarCraft II: A New Challenge for Reinforcement Learning. *arXiv:1708.04782*, 2017.
- [154] Niklas Wahlström, Thomas B Schön, and Marc P Deisenroth. Learning Deep Dynamical Models from Image Pixels. *IFAC SYSD*, 48(28), 2015.
- [155] Niklas Wahlström, Thomas B Schön, and Marc P Deisenroth. From Pixels to Torques: Policy Learning with Deep Dynamical Models. In *ICML Workshop on Deep Learning*, 2015.
- [156] Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Rémi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to Reinforcement Learn. In *CogSci*, 2017.
- [157] Ziyu Wang, Nando de Freitas, and Marc Lanctot. Dueling Network Architectures for Deep Reinforcement Learning. In *ICLR*, 2016.
- [158] Ziyu Wang, Victor Bapst, Nicolas Heess, Volodymyr Mnih, Rémi Munos, Koray Kavukcuoglu, and Nando de Freitas. Sample Efficient Actor-Critic with Experience Replay. In *ICLR*, 2017.
- [159] Christopher JCH Watkins and Peter Dayan. Q-Learning. *Machine Learning*, 8(3–4):279–292, 1992.
- [160] Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to Control: A Locally Linear Latent Dynamics Model for Control from Raw Images. In *NIPS*, 2015.
- [161] Théophane Weber, Sébastien Racanière, David P Reichert, Lars Buesing, Arthur Guez, Danilo Jimenez Rezende, Adrià Puigdomènec Badia, Oriol Vinyals, Nicolas Heess, Yujia Li, et al. Imagination-Augmented Agents for Deep Reinforcement Learning. In *NIPS*, 2017.
- [162] Paul John Werbos. Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences. Technical report, Harvard University, Applied Mathematics, 1974.
- [163] Daan Wierstra, Alexander Förster, Jan Peters, and Jürgen Schmidhuber. Recurrent Policy Gradients. *Logic Journal of the IGPL*, 18(5):620–634, 2010.
- [164] Ronald J Williams. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning. *Machine Learning*, 8(3–4): 229–256, 1992.
- [165] Markus Wulfmeier, Peter Ondruska, and Ingmar Posner. Maximum Entropy Deep Inverse Reinforcement Learning. In *NIPS Workshop on Deep Reinforcement Learning*, 2015.
- [166] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron C Courville, Ruslan Salakhutdinov, Richard S Zemel, and Yoshua Bengio. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In *ICML*, volume 14, 2015.
- [167] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. Target-Driven Visual Navigation in Indoor Scenes using Deep Reinforcement Learning. In *ICRA*, 2017.
- [168] Barret Zoph and Quoc V Le. Neural Architecture Search with Reinforcement Learning. In *ICLR*, 2017.

**Kai Arulkumaran** (ka709@imperial.ac.uk) is a Ph.D. candidate in the Department of Bioengineering at Imperial College London. He received a B.A. in Computer Science at the University of Cambridge in 2012, and an M.Sc. in Biomedical Engineering at Imperial College London in 2014. He was a Research Intern in Twitter Magic Pony and Microsoft Research in 2017. His research focus is deep reinforcement learning and transfer learning for visuomotor control.

**Marc Peter Deisenroth** (m.deisenroth@imperial.ac.uk) is a Lecturer in Statistical Machine Learning in the Department of Computing at Imperial College London and with PROWLER.io. He received an M.Eng. in Computer Science at the University of Karlsruhe in 2006 and a Ph.D. in Machine Learning at the Karlsruhe Institute of Technology in 2009. He has been awarded an Imperial College Research Fellowship in 2014 and received Best Paper Awards at ICRA 2014 and ICCAS 2016. He is a recipient of a Google Faculty Research Award and a Microsoft Ph.D. Scholarship. His research is centred around data-efficient machine learning for autonomous decision making.

**Miles Brundage** (miles.brundage@philosophy.ox.ac.uk) is a Ph.D. candidate in Human and Social Dimensions of Science and Technology at Arizona State University, and a Research Fellow at the University of Oxford’s Future of Humanity Institute. He received a B.A. in Political Science at George Washington University in 2010. His research focuses on governance issues related to artificial intelligence.

**Anil Anthony Bharath** (a.bharath@imperial.ac.uk) is a Reader in the Department of Bioengineering at Imperial College London and a Fellow of the Institution of Engineering and Technology. He received a B.Eng. in Electronic and Electrical Engineering from University College London in 1988, and a Ph.D. in Signal Processing from Imperial College London in 1993. He was an academic visitor in the Signal Processing Group at the University of Cambridge in 2006. He is a co-founder of Cortexica Vision Systems. His research interests are in deep architectures for visual inference.