# Exploration in Deep Reinforcement Learning: A Comprehensive Survey

Tianpei Yang, Hongyao Tang, Chenjia Bai, Jinyi Liu, Jianye Hao, Zhaopeng Meng, Peng Liu and Zhen Wang

*Abstract*—Deep Reinforcement Learning (DRL) and Deep Multi-agent Reinforcement Learning (MARL) have achieved significant successes across a wide range of domains, including game AI, autonomous vehicles, robotics and so on. However, DRL and deep MARL agents are widely known to be sample-inefficient that millions of interactions are usually needed even for relatively simple problem settings, thus preventing the wide application and deployment in real-industry scenarios. One bottleneck challenge behind is the well-known exploration problem, i.e., how to efficiently explore the environment and collect informative experiences that could benefit policy learning towards the optimal ones. This problem becomes more challenging in complex environments with sparse rewards, noisy distractions, long horizons, and non-stationary co-learners. In this paper, we conduct a comprehensive survey on existing exploration methods for both single-agent and multi-agent RL. We start the survey by identifying several key challenges to efficient exploration. Then we provide a systematic survey of existing approaches by classifying them into two major categories: uncertainty-oriented exploration and intrinsic motivation-oriented exploration. Beyond the above two main branches, we also include other notable exploration methods with different ideas and techniques. In addition to algorithmic analysis, we provide a comprehensive and unified empirical comparison of different exploration methods for DRL on a set of commonly used benchmarks. According to our algorithmic and empirical investigation, we finally summarize the open problems of exploration in DRL and deep MARL and point out a few future directions.

*Index Terms*—Deep Reinforcement Learning, Multi-Agent Systems, Exploration, Uncertainty, Intrinsic Motivation.

## I. INTRODUCTION

In recent years, Deep Reinforcement Learning (DRL) and deep Multi-agent Reinforcement Learning (MARL) have achieved excellent results across a wide range of domains, including Go [1], [2], Atari [3], StarCraft [4], [5], Robotics [6], and so on. This reveals the tremendous potential of DRL, and it is considered to be a promising solution to real-world sequential decision-making problems. Despite the successes in many domains, there is still a long way for DRL and deep MARL agents to be widely applied and deployed in real-world problems, because they are widely known to be sample-inefficient and millions of interactions are usually needed even for some relatively simple problem settings. For example, Agent57 [7] is the first DRL algorithm that is able to outperform

Tianpei Yang, HongyaoTang, JinyiLiu, Jianye Hao and Zhaopeng Meng are with the College of Intelligence and Computing, Tianjin University, Tianjin 300350, China. Chenjia Bai and Peng Liu are with the School of Computer Science and Technology, Harbin Institute of Technology, Harbin 150001, China. Zhen Wang is with Northwestern Polytechnical University, Xi'an 710060, China. (*Tianpei Yang, Hongyao Tang, Chenjia Bai and Jinyi Liu equally contribute to this work.*) (*Corresponding author: Jianye Hao.*)
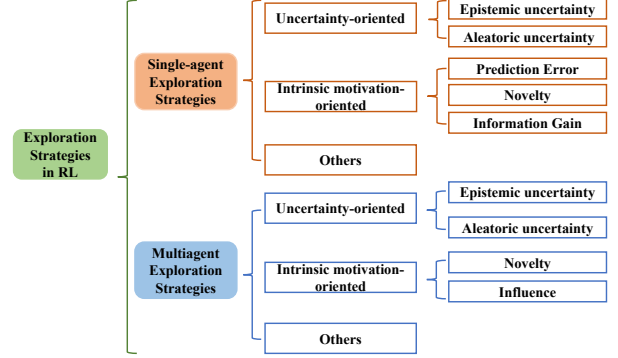


Fig. 1: Illustration of our taxonomy of the current literature on methods for DRL and deep MARL.

the average human player on all 57 Atari games; but the number of interactions it needs is several orders of magnitude larger than that of the average human player. The issue of sample-inefficiency naturally becomes more severe in multi-agent settings since the state-action space grows exponentially with the number of agents involved. One bottleneck challenge behind is the exploration problem, i.e., how to efficiently explore the unknown environments and collect informative experiences that could benefit the policy learning most towards optimal ones. This can be even more challenging in complex environments with sparse rewards, noisy distractions, long horizons, and non-stationary co-learners. Therefore, how to efficiently explore the environments is a significant and fundamental problem for DRL and deep MARL.

In recent years, many progresses have been made on exploration from different perspectives. However, a comprehensive survey on exploration in DRL and deep MARL is currently missing. A few prior papers contain the investigation on exploration methods for DRL [8], [9], [10], [11] and MARL [12], [13], [14]. However, these works are on general topics of DRL and deep MARL, thus their surveys on exploration methods are incomplete and lack of in-depth analysis. Aubret et al. [15] conduct a survey on intrinsic motivation in DRL. They investigate how this idea is developed to learn complicated and generalizable behaviours for problems like exploration, hierarchical learning, skill discovery and so on. Since exploration is not their focus, the ability of intrinsic motivation in addressing different exploration problems is not well studied. Besides, several works study the exploration methods only for multi-arm bandits [16], while they are often incompatible with deep neural networks directly and have the issue of scalability in complex problems. Due to the lack of a comprehensive

survey on exploration, the advantages and limitations of existing methods are seldom compared and studied in a unified scheme. This prevents further understandings of the exploration problem in RL community. To this end, we aim at providing an in-depth analysis of existing exploration methods for DRL and deep MARL in a unified view from both algorithmic and experimental perspectives.

In this paper, we propose a comprehensive survey for exploration methods in DRL and deep MARL. We focus on the model-free setting where the agent learns its policy without access to the environment model since the setting is more popular in the existing works on exploration. The general principles of exploration studied in the model-free setting are also shared by model-based methods in spite of different ways to realize. We start our survey by identifying five key challenges to achieve efficient exploration in DRL and deep MARL. The abilities in addressing these challenges also serve as criteria when we analyze and compare the existing exploration methods. The overall taxonomy of this survey is given in Fig. 1. For both exploration in DRL and deep MARL, we classify the existing exploration methods into two major categories based on their core ideas and algorithmic characteristics. The first category is uncertainty-oriented exploration whose core idea originates from the principle of *Optimism in the Face of Uncertainty* (OFU). The essence of this category is to leverage the quantification of epistemic and aleatoric uncertainty as a general means to measure the sufficiency of learning and the intrinsic stochasticity, based on which efficient exploration can be derived by following the OFU principle. The second category is the intrinsic motivation-oriented exploration. In developmental psychology, intrinsic motivation is considered as the primary driver in the early stages of human development [17], [18], e.g., children often employ less goal-oriented exploration but use curiosity to gain knowledge about the world [19]. Taking such an inspiration, the methods in this category heuristically make use of various reward-agnostic information as intrinsic motivation of exploration. Note that to some degree, methods in one of the above two categories may have underlying connections to some methods in the other one, usually from specifically intuitive or theoretical perspectives. In our taxonomy, we classify each method by sticking to its origination in motivation and algorithm as we describe above, as well as referring to the conventional literature it anchors. Beyond the above two main streams, we also include a few other advanced exploration methods, which show potential in solving hard-exploration tasks.

In addition to analysis and comparison from the algorithmic perspective, we provide a unified empirical evaluation of representative DRL exploration methods among several typical exploration environments, in terms of cumulative rewards and sample efficiency. The benchmarks demonstrate the successes and failures of the compared methods, showing the efficacy of corresponding algorithmic techniques. Uncertainty-oriented exploration methods show general improvements on exploration and learning in most environments. Nevertheless, it is nontrivial to estimate the uncertainty with a high quality in complex environments, which is significant to the effectiveness of uncertainty-oriented exploration methods. By contrast, intrinsic motivation can significantly boost exploration in environments with sparse, delayed rewards but may also cause deterioration in conventional environments due to the deviation of learning objectives. At present, the research on exploration in deep MARL is at an early stage. Most current methods for deep MARL exploration share the same major ideas of exploration and extend the techniques in single-agent exploration. In addition to the issues in single-agent setting, these methods face difficulties in addressing the larger joint-action space, the inconsistency of individual exploration behaviors, etc. Besides, the lack of common benchmarks prevents a unified empirical comparison. Finally, we conclude that both the algorithmic and empirical analysis above and highlight several significant remaining challenges of exploration in DRL and deep MARL, followed by some potential directions for future study.

We summarize our main contributions as follows:

- We give a comprehensive survey on exploration in DRL and deep MARL with a novel taxonomy for the first time.
- We analyze the strengths and weaknesses of representative articles on exploration for DRL and deep MARL, along with their abilities in addressing different challenges.
- We provide a unified empirical comparison of representative exploration methods in DRL among several typical exploration benchmarks.
- We highlight existing exploration challenges, open problems, and future directions for DRL and deep MARL.

The following of this paper is organized as follows. Sec. II describes preliminaries on basic RL algorithms and exploration methods. Then we introduce the challenges for exploration in Sec. III. According to our taxonomy, we present the exploration methods for DRL and deep MARL in Sec. IV and Sec. V, respectively. In Sec. VI, we provide a unified empirical analysis of different exploration methods in commonly adopted exploration environments; moreover, we discuss several open problems in this field and several promising directions for future research. Finally, Sec. VII gives the conclusion.

## II. PRELIMINARIES

### A. Markov Decision Process and Markov Game

**Markov Decision Process (MDP).** An MDP is usually defined as $\langle S, A, T, R, \rho_0, \gamma \rangle$, with a set of states $S$, a set of actions $A$, a stochastic transition function $T : S \times A \to P(S)$, which represents the probability distribution over possible next states, given the current state and action, a reward function $R : S \times A \to \mathbb{R}$, an initial state distribution $\rho_0 : S \to \mathbb{R}_{\in[0,1]}$, and a discounted factor $\gamma \in [0, 1)$. An agent interacts with the environment by performing its policy $\pi : S \to P(A)$. The agent's objective is to maximize the expected cumulative discounted reward: $J(\pi) = \mathbb{E}_{\rho_0, \pi, T} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right]$, where $s_0 \sim \rho_0(s_0)$, $a_t \sim \pi(s_t)$, $s_{t+1} \sim T(\cdot|s_t, a_t)$ and $r_t = R(s_t, a_t)$.

**Markov Game (MG).** MG is a multi-agent extension of MDP, which is defined as $\langle S, N, \{A^i\}_{i=1}^N, T, \{R^i\}_{i=1}^N, \rho_0, \gamma, Z, \{O^i\}_{i=1}^N \rangle$, additionally with action sets for each of $N$ agents, $A^1, ..., A^N$, a state transition function, $T : S \times A^1 \times ... \times A^N \to P(S)$, a reward function for each agent $R^i : S \times A^1 \times ... \times A^N \to \mathbb{R}$. For partially observable Markov games, each agent $i$ receives a local observation $o^i : Z(S, i) \to O^i$ and interacts

with environment with its policy $\pi^i : O^i \rightarrow P(A^i)$. The goal of each agent is to learn a policy that maximizes its expected discounted return, i.e., $J^i(\pi^i) = \mathbb{E}_{\rho_0, \pi^1, ..., \pi^N, T} \left[ \sum_{t=0}^{\infty} \gamma^t r_t^i \right]$, where $r_t^i = R^i(s_t, a_t^1, ..., a_t^N)$.

### B. Reinforcement Learning

In this section, we briefly introduce basic concepts and representative methods of RL and MARL, based on which most exploration approaches we review in this paper are designed.

RL is a learning paradigm of learning from interactions with the environment [20]. One central notion of RL is value function, which defines the expected return for a state or state-action pair to obtain by following a policy. Formally, *value function* $v^\pi(s)$ and *action-value function* $q^\pi(s, a)$ are defined as, $v^\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s \right]$ and $q^\pi(s, a) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s, a_0 = a \right]$. Based on the above definition, most RL algorithms can be divided into two categories: value-based methods and policy-based methods. Value-based methods usually learn the value functions from which the policies are derived implicitly. In contrast, policy-based methods maintain explicit policies and optimize to maximize the RL objective $J(\pi)$. In this survey, we focus on model-free RL algorithms.

**Value-based Methods.** Deep $Q$-Network (DQN) [21] is the most representative algorithm in DRL that derived from $Q$-learning. The $Q$-function $Q(s, a; \theta)$ parameterized with $\theta$ is learned by minimizing Temporal Difference (TD) loss:

$$\mathcal{L}^{\text{DQN}}(\theta) = \mathbb{E}_{(s_t, a_t, r_t, s'_{t+1}) \sim D} \left[ y - Q(s_t, a_t; \theta) \right]^2, \quad (1)$$

where $y = r_t + \gamma \max_{a'} Q(s_{t+1}, a'; \theta^-)$ is the target value, $D$ is the replay buffer and $\theta^-$ denotes the parameters of the target network. Further, a variety of variants are proposed to improve the learning performance of DQN, from different perspectives such as addressing the approximation error [22], [23], modeling value functions in distributions [24], [25], and other advanced structures and training techniques [26], [27], [3].

**Policy Gradients Methods.** Policy-based methods optimize a parameterized policy $\pi_\phi$ by performing gradient ascent on $J(\pi_\phi) = J(\phi)$ with regard to policy parameters $\phi$. According to the *Policy Gradients Theorem* [20], $\phi$ is updated as below:

$$\nabla_\phi J(\phi) = \mathbb{E}_{\pi_\phi} \left[ \nabla_\phi \log \pi_\phi(a_t | s_t) q^{\pi_\phi}(s_t, a_t) \right]. \quad (2)$$

One typical policy gradient algorithm is REINFORCE [28] that uses the complete return $G_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i}$ as estimates $\hat{Q}(s_t, a_t)$ of $q^{\pi_\phi}(s, a)$, i.e., Monte Carlo value estimation. To further ensure an effective policy updates, Proximal Policy Optimization (PPO) [29] proposes a modified surrogate objective:

$$\mathcal{L}^{\text{PPO}}(\phi) = \mathbb{E}_{\pi_{\phi^-}} \big[ \min \big( \rho_t \hat{A}(s_t, a_t), \\ \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) \hat{A}(s_t, a_t) \big) \big], \quad (3)$$

where $\hat{A}(s_t, a_t)$ is the estimation of the advantage of old policy $\pi_{\phi^-}$, $A^{\pi_{\phi^-}}(s_t, a_t) = q^{\pi_{\phi^-}}(s_t, a_t) - v^{\pi_{\phi^-}}(s_t)$, and $\rho_t = \frac{\pi_\phi(a_t, s_t)}{\pi_{\phi^-}(a_t, s_t)}$ is the importance sampling ratio.

**Actor-Critic Methods.** Actor-Critic methods usually approximate value functions and estimate $q^{\pi_\phi}(s, a)$ in Eq. (2) with

bootstrapping [20], conventionally $\hat{Q}(s_t, a_t) = r_t + \gamma V(s_t)$. Apart from the stochastic policy gradients (Eq. (2)), *Deterministic Policy Gradients* (DPG) [30] allows deterministic policy, formally $\mu_\phi$, to be updated via maximizing an approximated $Q$-function. In DDPG [31], the policy is updated as:

$$\nabla_\phi J(\phi) = \mathbb{E}_{\mu_\phi} \left[ \nabla_\phi \mu_\phi(s) \nabla_a Q^{\mu_\phi}(s, a) |_{a = \mu_\phi(s)} \right]. \quad (4)$$

**MARL Algorithms.** In MARL, agents learn their policies in a shared environment whose dynamics are influenced by all agents. The most straightforward way is to let each agent learn a decentralized policy independently, treating other agents as part of the environment, i.e., *Independent Learning* (IL). The main issue of IL is the non-stationary problem due to independent updates of multiple agents. A few works study how to stabilize the learning process [32], [33]. In contrast, some works assume global information can be accessible or at least sometimes during execution, on which communication protocols can be built [34], [35]. The most popular MARL learning paradigm is *Centralized Training and Decentralized Execution* (CTDE). CTDE allows agents to fully utilize global information (e.g., global states, other agents' actions) during training while only local information (e.g., local observations) is required during execution [36], [37], [38]. CTDE provides good training efficiency and practical execution policies.

One representative CTDE method is MADDPG [38], a multi-agent Actor-Critic for mixed cooperative-competitive problems. Concretely, consider a game with $N$ agents with deterministic policies $\{\mu_i\}_{i=1}^N$ parameterized by $\{\phi_i\}_{i=1}^N$, the deterministic policy gradient for each agent $i$ can be:

$$\nabla_{\phi_i} J(\phi_i) = \mathbb{E}_{x, \vec{a}} \big[ \nabla_{\phi_i} \mu_{\phi_i}(o_i) \\ \nabla_{a_i} Q_i^{\vec{\mu}}(x, a_1, a_2, ..., a_N) |_{a_i = \mu_{\phi_i}(o_i)} \big], \quad (5)$$

where $x = (o_1, o_2, ..., o_N)$ is a common choice of global state information and the vector $\vec{a}$, $\vec{\mu}$ denote the concatenation of corresponding variables. Each agent maintains its own critic $Q_i$ which estimates the joint observation-action value function and uses the critic to update its decentralized policy.

### C. Basic Exploration Techniques

In this section, we briefly review several exploration methods that are commonly used in DRL methods.

$\epsilon$**-Greedy.** The most basic exploration method is $\epsilon$-greedy: with a probability of $1 - \epsilon$, the agent chooses the action greedily (i.e., exploitation); and a random choice is made otherwise (i.e., exploration). Albeit the popularity and simplicity, $\epsilon$-greedy is inefficient in complex problems with large state-action space.

*Boltzmann* **Exploration.** Another exploration method in RL is *Boltzmann* (*softmax*) exploration: agent draws actions from a Boltzmann distribution over its $Q$-values. Formally, the probability of choosing action $a$ is as, $p(a) = \frac{e^{Q(s,a)/\tau}}{\sum_i e^{Q(s,a^i)/\tau}}$, where the temperature parameter $\tau$ controls the degree of the selection strategy towards a purely random strategy, i.e., the higher value the $\tau$, the more randomness selection strategy. The drawback of Boltzmann exploration is that it cannot be directly applied in continuous state-action spaces.

**Upper Confidence Bounds (UCB).** UCB is a classic exploration method originally from Multi-Armed Bandits

(MABs) problems [39]. In contrast to performing unintended and inefficient exploration as for naive random exploration methods (e.g., $\epsilon$-greedy), the methods among UCB family measure the potential of each action by an upper confidence bound of the reward expectation. Formally, the selected action can be calculated as follows,

$$a = \arg\max_a Q(a) + \sqrt{\frac{-\ln k}{2N(a)}} \quad (6)$$

where $Q(a)$ is the reward expectation of action $a$, $N(a)$ is the count of selecting action $a$ and $k$ is a constant.

**Entropy Regularization.** While value-based RL methods add randomness in action selection based on $Q$-values, entropy regularization is widely used to promote exploration for RL algorithms [40] with stochastic policies, by adding the policy entropy $H(\pi(a|s))$ to the objective function as a regularization term, encouraging the policy to take diverse actions. Such a regularization may deviate the original objective function to optimize. One solution is to gradually decay the influence of the entropy regularization during the learning process.

**Noise Perturbation.** As to deterministic policies, noise perturbation is a natural way to induce exploration. Concretely, by adding noise perturbation sampled from a stochastic process $\mathcal{N}$ to the deterministic policy $\pi(s)$, an exploration policy $\pi'(s)$ is constructed, i.e., $\pi'(s) = \pi(s) + \mathcal{N}$. The stochastic process can be selected to suit the environment, e.g., an Ornstein-Uhlenbeck process is preferred in physical control problems [38]; more generally, a standard Gaussian distribution is simply adopted in many works [41]. However, such vanilla noise-based exploration is unintentioned exploration which can be inefficient in complex learning tasks with exploration challenges.

### D. Exploration based on Bayesian Optimization

In this section, we take a brief view with respect to Bayesian Optimization (BO) [42], which could perform efficient exploration to find where the global optimum locates. The objective of BO is to solve the problem: $x^* = \arg\max_{x \in X} f(x)$, where $x$ is what going to be optimized, $X$ is the candidate set, and $f$ is the objective function: $x \to \mathbb{R}$, which is always black-box.

BO strategies treat the objective function $f$ as a random function, thus a Bayesian statistic model for $f$ is the basic component of Bayesian optimization. The Bayesian statistic model is then used to construct an acquisition function $\alpha(x)$, which is tractable and provides the selection metric for each sample $x$. The acquisition function should be designed for more efficient exploration, to seek for a probably better selection that has not been attempted yet, such as GP-UCB [43], Thompson Sampling (TS) [44], etc. In the following, we mainly introduce these acquisition functions, which have been applied in DRL algorithms, and explain how they facilitate the efficient exploration.

**Gaussian Process-Upper Confidence Bound (GP-UCB).** UCB [39], [45] is a commonly used exploration strategy, as shown in Sec. II-C. In BO, modelling the Bayesian statistic model as Gaussian Process (GP), GP-UCB [43] uses the mean $\mu(x)$ and standard deviation $\sigma(x)$ of the GP posterior to determine the upper confidence bound, and makes decisions in

such optimistic estimation. Acquisition function of GP-UCB is $\alpha(x) = \mu_{t-1}(x) + \beta_t^{1/2} \sigma_{t-1}(x)$, where the added bonus $\sigma(x)$ guides optimistic exploration.

**Thompson Sampling** [44], also known as posterior sampling or probability matching, samples a function $f$ from Posterior($f$) each time, and takes actions greedily with respect to such randomly drawn belief $f$ as follows: $\alpha(x) = f(x), s.t. f \sim$ Posterior($f$). TS captures the uncertainty according to the posterior and enables deep exploration. TS has been employed in DRL for efficient exploration, which we discuss in Sec. IV-A.

## III. WHY EXPLORATION IN RL IS HARD

In this section, we identify several typical environmental challenges to efficient exploration for DRL and deep MARL. For each of these challenges, we analyze their causes and characteristics by illustrative examples and also point out some key factors of a potential solution. Most exploration methods are proposed to address one or some of these challenges. Therefore, these challenges serve as significant touchstones for the analysis and evaluation in the following of this paper.

**Large State-action Space.** The difficulty of DRL naturally increases with the growing of the state-action space. For example, real-world robots often have high-dimensional sensory inputs like images or high-frequency radar signals, and have a large number of degrees of freedom for delicate manipulation. Another practical example is the recommendation system with graph-structured data as states and large number of discrete actions. In general, when facing large state-action space, exploration can be inefficient, since it can take an unaffordable budget to well explore the state-action space. Moreover, for more complex cases, the state-action space may have complex underlying structures: states are inaccessible at equal efforts where causal dependencies exist among states; the actions can be combinatorial and hybrid of discrete and continuous components. These practical problems make efficient exploration more challenging.

Most exploration methods compatible with deep neural networks are able to deal with high-dimensional state space. With low-dimensional compact state representations learned, exploration and learning can be conducted efficiently. Towards efficient exploration in complex state space, sophisticated strategies may be needed to quantify the degree of familiarity of states, avoid unnecessary exploration, reach the 'bottleneck' states, and so on. By contrast, the exploration issues in large action space are lack of study currently. Though several specific algorithms that tailored for complex action structures are proposed [46], [47], [48], how to efficiently explore in a large action space is still an open question. A further discussion can be seen in Sec. VI-B.
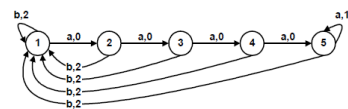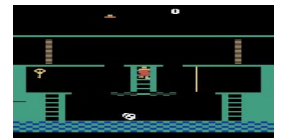


Fig. 2: Chain MDP.  Fig. 3: Montezuma's Revenge.

**Sparse, Delayed Rewards.** Another major challenge is to explore in environments with sparse, delayed rewards. Basic exploration strategies can rarely discover meaningful states or obtain informative feedback in such scenarios. One of the most representative environments is Chain MDP [49], [50], as the instance with five states and two actions (*left* and *right*) illustrated in Fig. 2. The optimal behavior is to always go right to receive a reward of $+10$ when reaching the state 5. Even with finite states and simple state representation, chain MDP is a classic hard exploration environment since the degree of reward sparsity and the difficulty of exploration increase as the chain becomes longer. Montezuma's Revenge is another notorious example with sparse, delayed rewards. Fig. 3 shows the first room of Montezuma's Revenge. A long sequence of specific actions is required to solve this task. An ideal agent needs to climb down the ladder, move left and collect the key where obtains a reward ($+100$); then the agent backtracks and navigates to the door and opening it with the key, resulting in final reward ($+300$). Without effective exploration, the agent would not be able to finish the task in such an environment. Intuitively, to solve such problems, an effective exploration method is to leverage reward-agnostic information as dense signals to explore the environment. In addition, the ability to perform a temporally-extended (deep) exploration is another key point in such circumstances. Temporally-extended exploration means the agent follows the same exploratory policy for a longer horizon (e.g., an entire episode). From the above two perspectives, several recent works [51], [52], [53], [50] have shown promising results, which will be discussed later.
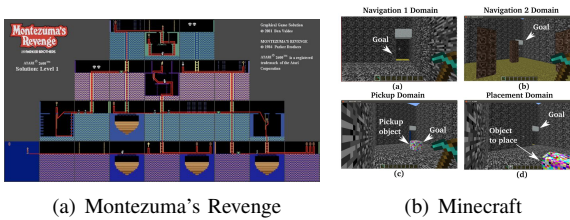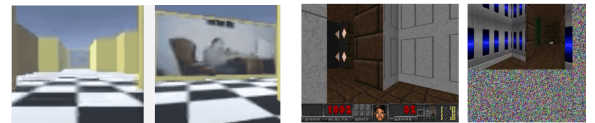


Fig. 4: Examples of environments with long horizon and extremely sparse, delayed rewards. (a) A panorama of 24 rooms in Montezuma's Revenge. (b) A four-domains environment in Minecraft.

Although a little progress has been made for this exploration problem in several representative environments, unfortunately more practical environments with long horizons and extremely sparse, delayed rewards remain far away from being solved. Reconsider the whole game of Montezuma's Revenge and a panorama of 24 different rooms is shown in Fig. 4(a). The agent needs to collect different objects (e.g., keys and treasures) and solve relations among objects (e.g, the correspondence between keys and doors) to trigger a series of critical events (e.g., discover the sword) in the environment. This long-horizon task requires a more complicated policy to accomplish the whole game than the one-room scene mentioned above. Another example is the four-domain environment in Minecraft [54] shown in Fig. 4(b). The agent with first-person view is expected to traverse the four domains with different obstacles and textures, to pick up the object and place it on the target location, then obtain the final reward. At present, few approaches are able to learn effective policies in such problems even with specific heuristics and prior knowledge. This is also regarded as a significant gap between DRL and real-world applications, which will be discussed an an open problem in Sec. VI-B.

**White-noise Problem.** The real-world environments often have high randomness, where usually unpredictable things appear in the observation or action spaces. For example, visual observations of autonomous cars contain predominantly irrelevant information, like continuously changing positions and shapes of clouds. In exploration literature, white-noise is often used to generate high entropy states and inject randomness into the environment. Because the white noise is unpredictable, the agent cannot build an accurate dynamics model to predict the next state. 'Noisy-TV' is a typical kind of white-noise environment. The noise of 'Noisy-TV' includes constantly changing background, irrelevant objects in observations, Gaussian noise in pixel space, and so on. We illustrate the typical 'Noisy-TV' scenarios in Fig. 5. In Fig. 5(a), we show a typical navigation task in Unity maze on the left and a 'Noisy-TV' variant on the right. The TV constantly changes channels when the agent chooses actions. The high entropy of TV becomes an irresistible attraction to the agent. In Fig. 5(b), we show a similar 'Noisy-TV' in VizDoom [55] on the right. The uncontrollable Gaussian noise is added to the observation space, which attracts the agent to stay in the current room and prevents it from passing through more rooms. Exploration methods that measure the novelty through predicting the future become unstable when confronting 'Noisy-TV' or similarly unpredictable stimuli. Existing methods mainly focus on learning state representation [51], [56], [57], [58] to improve robustness when faced with the white-noise problem. How an agent can explore robustly in stochastic environments is an important research branch in DRL. We will discuss this problem further in Sec. VI-B.



(a) A typical Unity maze (left) and a 'Noisy-TV' variant (right)  (b) A typical VizDoom (left) and a noisy variant (right)

Fig. 5: Tasks with Noisy-TV Problem

**Multi-agent Exploration.** Except for the above challenges, explorations in multi-agent settings are more arduous. 1) Exponential increase of the state-action space. The straightforward challenge is that the joint state-action space increases exponentially with the increase in the number of agents, making it much more difficult to explore the environment. Therefore, the crucial point is how to carry out a comprehensive and effective exploration strategy while reducing the cost of exploration. A naive exploration strategy is that agents execute individual exploration based on their local information.

2) Coordinated exploration. Although individual exploration avoids the exponential increase of the joint state-action space, it induces extra difficulties in the exploration measurement due to partial observation and non-stationary problems. The estimation based on local observation is biased and cannot reflect the
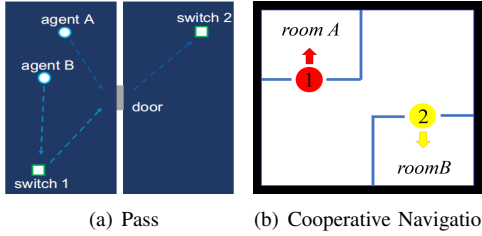
(a) Pass      (b) Cooperative Navigation

Fig. 6: Tasks with multi-agent exploration problem.

global information of the environment. Furthermore, an agent's behaviors are influenced by other coexisting agents, which induces extra stochasticity. Thus, such an approach always fails in tasks which need agents to behave cooperatively. For example, Fig. 6(a) illustrates a maze consisting of two rooms with a door and two switches. This is a typical environment that requires agents to perform coordinated exploration. Two agents are only rewarded when both of them go through the door and reach the right room; while the door will be open only when one of the two switches is covered. Thus, the optimal strategy should be one agent first steps on switch 1 to open the door, then the other agent goes to the right room, and further steps on switch 2 to hold the door open, and lets the remaining agent in. Challenges arise in such an environment where agents should visit some critical states in a cooperative manner, through which the optimal policy can be achieved.

3) Local- and global-exploration balance. To achieve such coordinated exploration strategies, only local information is insufficient, the global information is also necessary. However, one main challenge is the inconsistency between local information and global information, which requires agents to make a balance between local and global perspectives, otherwise, it may lead to inadequate or redundant exploration. Fig. 6(b) shows such inconsistency where a homogeneous agent 1 (red dot) and agent 2 (yellow dot) are required to search the big room with two small rooms, (i.e., room $A$ and room $B$). In this environment, from the global perspective, either agent 1 visits room $A$, agent 2 visits room $B$ or agent 1 visits room $B$, agent 2 visits room $A$ are equally treated, which means the latter does not make an increase of the known global states. However, if agents explore relying more on local information, each agent will still try to visit another room to increase its known local states, which will cause redundant exploration from the global view. Furthermore, agents cannot simply explore relying on global information, since how much each agent contributes to the global exploration is unknown. The trade-off of exploring locally and globally is the key problem that needs to be addressed to facilitate more efficient multi-agent exploration, which will be further discussed in Sec. VI-B.

## IV. EXPLORATION IN SINGLE-AGENT DRL

In this section, we investigate exploration works in single-agent RL that leverage different theories and heuristics to address or alleviate the above challenges. Based on the different key ideas and principles of these methods, we classify them into two major categories (shown in Fig. 1). The first category is uncertainty-oriented exploration, which originates from the

optimistic principle. This kind of work usually explicitly models the epistemic uncertainty via the Bayesian posterior of the value function, or additionally measures the aleatoric uncertainty through distributional value functions. The agent is encouraged to explore the regions with high epistemic uncertainty and avoid visiting the areas with high aleatoric uncertainty. The second category is intrinsic motivation-oriented exploration, which is inspired by intrinsic motivation in psychology [19], [17], [18] that intrinsically rewards exploration activities to encourage such behaviors. By extending such an idea to RL, this kind of work usually uses reward-agnostic information to design intrinsic rewards to induce exploration. In addition, we also conclude other techniques for exploration beyond these two mainstream methods.

### A. Uncertainty-oriented Exploration

We first introduce uncertainty-based exploration methods. In model-free RL, this kind of methods often measure the uncertainty of the value function; and model-based RL usually considers the uncertainty of dynamics. There exist two types of uncertainty-based exploration. 1) *Epistemic uncertainty*, which represents the errors that arise from insufficient and inaccurate knowledge about environment [73], [69], is also named parametric uncertainty. Many strategies provide efficient exploration by following the principle of *Optimism in the Face of Uncertainty* (OFU) [74], [75], where higher epistemic uncertainty is considered as insufficient knowledge of the environment. OFU encourages the agent to visit states and actions with higher epistemic uncertainty to explore the unknown environment. There are some approaches proposed for estimating epistemic uncertainty, such as MC dropout [76], bootstrap sampling [65], ensemble estimators [77], [69], and broadly Bayesian regression [78], [79]. The second kind of uncertainty named 2) *Aleatoric uncertainty* represents the intrinsic randomness of the environment, and can be captured by the return distribution [24], [69], which is also known as intrinsic uncertainty or return uncertainty.

Based on the uncertainty estimation, there exists mainly two ways to utilize the uncertainty. 1) Following the idea of UCB [39], [45], a direct method for exploring states and actions with high uncertainty is **performing optimistic action-selection** by choosing the action to maximize the optimistic value function $Q^+$ in each time step, as

$$a_t = \arg\max_a \ Q^+(s_t, a), \qquad (7)$$

where $Q^+(s_t, a_t) = Q(s_t, a_t) + \text{Uncertainty}(s_t, a_t)$ indicates the ordinary Q-value added by an exploration bonus based on a specific uncertainty measurement, which is usually derived by the posterior estimation of value function or dynamics [80]. The optimistic action selection needs to compute a bonus for each time step to balance the current Q-value and the uncertainty in exploration. 2) Following the idea of Thompson Sampling [44] to perform exploration. The action selection is greedy to a **sampled value function $Q_\theta$ from the $Q$-posterior**, as

$$a_t = \arg\max_a Q_\theta(s_t, a), \quad Q_\theta \sim \text{Posterior}(Q), \qquad (8)$$

TABLE I: Characteristics of uncertainty-oriented exploration algorithms in RL. We use whether the algorithm uses parametric or non-parametric posterior to distinguishes methods. The last two properties correspond to the challenges shown in Sec. III. We use blank, partially, high and check mark to indicate the extent to the method address a specific problem.

| | | Method | Large State Space | Continuous Control | Long-horizon | White-noise |
|---|---|---|---|---|---|---|
| Epistemic Uncertainty | Parametric Posterior | RLSVI[59] [60] | | | high | partially |
| | | Bayesian DQN[61] | | | high | partially |
| | | Successor Uncertainty [62] | ✓ | | high | partially |
| | | Wasserstein DQN [63] | ✓ | | high | partially |
| | | UBE [64] | ✓ | | high | partially |
| | Non-parametric Posterior | Bootstrapped DQN[65] | ✓ | | high | partially |
| | | OAC[66] | ✓ | ✓ | high | partially |
| | | SUNRISE [67] | ✓ | ✓ | high | partially |
| | | OB2I [68] | ✓ | | high | partially |
| Epistemic & Aleatoric Uncertainty | Non-parametric Posterior | DUVN[69] | ✓ | | partially | partially |
| | | IDS[70] | ✓ | | high | high |
| | | DLTV with QR-DQN [71] | ✓ | | | high |
| | | DLTV with NC-QR-DQN[72] | ✓ | | | high |

that is, to first estimate the posterior distribution of Q-function through a parametric or non-parametric posterior, then sample a Q-function $Q_\theta$ from this posterior and use $Q_\theta$ for action-selection when interacting with the environment for a whole episode. Compared with UCB-based exploration, the Thompson sampling method uses the same Q-function in the whole episode rather than performing optimistic action-selection in each time step, thus the agent is enabled to perform *deep* exploration and has advantages in long-horizon exploration tasks. In the following, we first introduce methods based on epistemic uncertainty only, and then present some recent methods that use both kinds of uncertainty. Table I shows the characteristics of uncertainty-oriented exploration methods in terms of whether it addresses the challenges discussed in Sec. III.

*1) Exploration via Epistemic Uncertainty:* Estimating the epistemic uncertainty usually needs to maintain a posterior of the value function. In the following, we categorize existing methods based on which type of posterior is used, including the parametric posterior and the non-parametric posterior.

*Parametric Posterior based Exploration.* Parametric posterior is usually learned by Bayesian regression in linear MDPs, where the transition and reward functions are assumed to be linear to state-action features. Randomized Least-Squares Value Iteration (RLSVI) [60] perform Bayesian regression in linear MDPs so that it is able to sample the value function through Thompson Sampling as Eq. (8). RLSVI is shown to attain a near-optimal worst-case regret bound [81] in linear settings, while such an approach is restricted to linear MDPs and is not applicable to large-scale RL or continuous control tasks. In DRL that uses neural networks as function approximators, Bayesian DQN [61] extends the idea of RLSVI to generalized function approximations. It is a practical Thompson Sampling method that considers the state-action encoder in the Q-network as a fixed feature extractor and performs Bayesian Linear Regression (BLR) in the last layer of the Q-network. BLR constructs the parametric posterior of the value function by approximately considering the value iteration of the last-layer Q-network as a linear MDP problem. Usually, BLR can only be applied for fixed inputs with linear function approximation

and cannot be directly used in DRL. Bayesian DQN solves this problem by approximately considering the feature mapping before the output layer as a fixed feature vector, and then performs BLR based on this feature vector. However, in DRL tasks, the features of high-dimensional state-action pairs are trained in the learning process, which violates the hypothesis of BLR that the features are fixed and may cause unstable performance in large-scale tasks. Based on a similar idea, Successor Uncertainty [62] approximates the posterior through successor features [82], [83]. Because the successor feature contains the discounted representation in an episode, the $Q$-value is linear to the successor feature of the corresponding state-action pair. BLR can be applied to measure the posterior of the value function based on the successor representation. RLSVI, Bayesian DQN and Successor Uncertainties partially address the long-horizon problem since Thomson sampling captures the long-term uncertainties, as shown in Table I.

Another approach that stands in a novel view to use epistemic uncertainty is Wasserstein Q-Learning (WQL) [63], which uses a Bayesian framework to model the epistemic uncertainty, and propagates the uncertainty across state-action pairs. Instead of applying a standard Bayesian update, WQL approximates the posterior distribution based on Wasserstein barycenters [84]. However, the use of the Wasserstein-Temporal-Difference update makes it computationally expensive to apply in complex environments. Uncertainty Bellman Equation and Exploration (UBE) [64] proposes another Bayesian posterior framework that uses an upper bound on the Q-posterior variance. This method is computationally tractable since it estimates the variance of the posterior rather than the whole posterior distribution. However, without the whole posterior distribution, UBE cannot perform Thomson sampling and instead uses the posterior variance as an exploration bonus as Eq. (7).

Nevertheless, the use of function approximation (i.e., neural network) often makes the estimation of the Bayesian posterior difficult. The reason is that most Bayesian methods consider a fixed representation of the state-action pair, and then model the posterior distribution based on this representation. While in large-scale tasks, the representations of high-dimensional state-

action space also need to be learned in the training process. Learning the Q-posterior and representation simultaneously brings instability to Bayesian methods. How to solve this problem is an important direction for future research.

*Non-Parametric Posterior based Exploration.* Beyond parametric posterior, bootstrap-based exploration [85] constructs a non-parametric posterior based on the bootstrapped value functions, which has theoretical guarantees in tabular and linear MDPs. In DRL, Bootstrapped DQN [65] maintains several independent $Q$-estimators and randomly samples one of them at the beginning of each episode, which enables the agent to conduct temporally-extended exploration since the agent considers long-term effects of exploration from the $Q$-function and follows the same exploratory policy in the entire episode. From the perspective of Thompson Sampling, each Q-network has the same probability to reflect the optimal policy, and the uniform distribution over Q-networks induces an instance of the randomized value function. Bootstrapped DQN is similar to RLSVI, but samples value function via bootstrapping instead of Gaussian distribution. Bootstrapped DQN is easy to implement and performs well, thus becoming a common baseline for *deep* exploration. Subsequently, Osband et al. [50] show that via injecting a 'prior' for each bootstrapped $Q$-function, the bootstrapped posterior can further increase the diversity of bootstrapped functions in regions with fewer observations, and thus improve the generalization of uncertainty estimation.

Besides Bootstrapped DQN, there are several works that combine other techniques with bootstrapped sampling to guide more efficient and deep exploration. Bootstrap policy gradient [86], Multi-DDPG [87], MABDDPG [88], and SOUP [89] combine bootstrapped sampling with policy gradient through actor-critic framework. All the exploration strategies explore through Thompson sampling, as in Eq. (8).

There are also several approaches that construct an optimistic value estimation using non-parametric posterior, as in Eq. (7). SUNRISE [67] integrates bootstrapped sampling to provide a bonus for optimistic action selection, and additionally adopts a weighted Bellman backup to prevent instability in error propagation. OAC [66] builds the upper bound of $Q$-value through two bootstrapped networks and explores by choosing optimistic actions based on the upper bound. Considering uncertainty propagation, OB2I [68] performs backward induction of uncertainty to capture the long-term uncertainty in an episode. The bootstrapped-based methods address the long-horizon problem since they sample a Q-function from the bootstrapped posterior and select actions based on this Q-function for a whole episode.

The parametric posterior-based methods can only handle discrete control problems since the update of LSVI and other Bayesian methods require the action space to be countable. However, the non-parametric posterior [50], [65], [68] based methods can be applied in continuous control to choose actions that maximize the sampled value function $Q_\theta(s, a)$ from the posterior. SUNRISE [67] proposes an approximation scheme that first learns an explicit actor and generates $N$ candidate action set from the current policy, and then chooses the action from these $N$ actions to maximize $Q_\theta(s, a)$. OAC [66] approximate the $Q$-function with a linear function $\bar{Q}$ that fits

$Q$ in a small region near the current policy, and then trains an independent actor to maximize $\bar{Q}$.

The above methods perform optimistic action selection or posterior sampling-based solely on epistemic uncertainty. There exist several other methods [69], [90], [91], [71] that consider both the epistemic uncertainty and the aleatoric uncertainty in exploration. Beyond estimating the epistemic uncertainty, additionally preserving the aleatoric uncertainty enables to prevent the agent from exploring areas with high randomness. We introduce these methods in the following.

*2) Exploration under both types Uncertainty:* For an environment with large randomness, the estimated uncertainty may be disturbed by aleatoric uncertainty, e.g., the ensemble estimators may be uncertain about a state-action pair not because that is seldom visited, but due to its large environment randomness. Meanwhile, since the aleatoric uncertainty cannot be reduced during training, to be optimistic about aleatoric uncertainty usually leads the agent to favor actions with higher variances, which hurts the performance. To consider both epistemic uncertainty and aleatoric uncertainty for exploration, Double Uncertain Value Network (DUVN) [69] firstly proposes to use Bayesian dropout [76] to measure the epistemic uncertainty, and return distribution to estimate the aleatoric uncertainty. However, DUVN does not tackle the negative impact of being optimistic to aleatoric uncertainty. Following the principle of OFU, the exploration strategy would be better if it not only is optimistic about epistemic uncertainty, but also tries to avoid the impact of aleatoric uncertainty.

Inspired by Information Directed Sampling (IDS) [70] in bandit settings, Nikolov et al. [90] extend the idea of IDS to general MDPs for efficient exploration by considering both epistemic and aleatoric uncertainties, and try to avoid the impact of aleatoric uncertainty. This method combines distributional RL (i.e., C51 [24]) to measure aleatoric uncertainty and bootstrapped $Q$-values to approximate epistemic uncertainty. Then the behavior policy is designed to balance the instantaneous regret and information gain. To improve the computational efficiency of IDS, Carmel et al. [91] estimate both types of uncertainty on the expected return through two networks. Specifically, aleatoric uncertainty is captured via the learned return distribution using QR-DQN [25], and the epistemic uncertainty is estimated on the Bayesian posterior by sampling parameters from two QR-DQN networks. However, considering the two types of uncertainties needs more computation caused by the use of distributional value functions.

Without explicitly estimating epistemic uncertainty, Decaying Left Truncated Variance (DLTV) [71] studies how to take advantage of the distributions learned by distributional Q-function for exploration directly. In order to gradually reduce the aleatoric uncertainty, DLTV uses the variance of quantiles as bonuses and applies a decay schedule to drop the effect of aleatoric uncertainty along with the training process. NC-QR [72] then improves DLTV through non-crossing quantile regression, which guarantees monotonicity of learned quantile curves by additional constraints.

Finally, we discuss the uncertainty-oriented exploration methods in dealing with the *white-noise problem* shown in Sec. III and Table I. Theoretically, if the posterior of value

function can be solved accurately (e.g., a closed-form solution), the epistemic uncertainty will be accurate and the exploration will be robust to the white-noise. However, in practice, since the Q-functions are usually estimated through parametric or non-parametric methods, exploration based on the posterior can still be affected by the randomness of the environment. Meanwhile, since we do not know the prior of value function exactly, we usually consider the prior is uninformative, which leads to an inaccurate estimation of the true posterior. In addition, for methods that additionally estimate the aleatoric uncertainty, the exploration is more robust since the agent can avoid exploring the noisy states. Nevertheless, estimating the noise is unstable and computationally expensive since it needs to learn the distributional $Q$-functions through a projected Bellman update or quantile regression.

### B. Intrinsic Motivation-oriented Exploration

Intrinsic motivation is originated from humans inherent tendency to interact with the world in an attempt to have an effect, and to feel a sense of accomplishment[92], [93]. It is usually accompanied by positive effects (rewards), thus intrinsic motivation-oriented exploration methods often design intrinsic rewards to create the sense of accomplishment for agents. In this section, we investigate previous works tackling the exploration problem based on intrinsic motivation. These works can be technically classified into three categories: 1) methods that estimate prediction errors of the environmental dynamics; 2) methods that estimate the state novelty; 3) methods based on the information gain. Table II presents the characteristics of all reviewed intrinsic motivation-oriented exploration algorithms in terms of whether it can apply to continuous control problems, and whether it can solve the white-noise and long-horizon problems described in Section III.

*1) Prediction Error:* The first class of works is based on prediction errors, which encourages agents to explore states with higher prediction errors. Specifically, for each state, the intrinsic reward is designed using its prediction error for the next state, which can be measured as the distance between the predicted next state and true one:

$$R(s_t, s_{t+1}) = \text{dist}\left(\phi(s_{t+1}), \hat{f}\left(\phi(s_t), a_t\right)\right) \quad (9)$$

where dist($\cdot$) can be any distance measurement function, $\phi$ is an encoder network that maps the raw state space to a latent space. $\hat{f}$ is a dynamic model that predicts the next latent state with the current latent state and action as input. In this direction, how to learn a suitable encoder $\phi$ is the main challenge.

To address this problem, Dynamic Auto-Encoder (Dynamic-AE) [94] is proposed to compute the distance between the predicted state and the true state in the latent state space, which is learned by an auto-encoder, where $\phi$ is the encoder. However, this approach is unable to handle the white-noise problem (Sec. III) since they design $\phi$ without removing the noisy distractions existing in the environment. Furthermore, it cannot handle the long-horizon problem since it does not consider temporally extended information. A solution to improve the robustness of $\phi$ is Intrinsic Curiosity Module (ICM)[56], which also learns the environment dynamics in a latent space. ICM

TABLE II: Characteristics of reviewed intrinsic motivation-oriented exploration algorithms in RL, where blank, partially, high and ✓ denote the extent that this method addresses this problem.

| | Method | Continuous Control | Long-horizon | White-noise |
|---|---|---|---|---|
| Prediction error | Dynamic-AE[94] | | | |
| | ICM [56] | | | high |
| | Curiosity-Driven[57] | | | |
| | AR4E[95] | | | high |
| | VDM[96] | ✓ | | high |
| | EMI [97] | ✓ | | |
| Novelty | TRPO-AE-hash[98] | ✓ | | partially |
| | A3C+[99] | ✓ | | partially |
| | DQN-PixelCNN[100] | | | partially |
| | $\phi$-EB[101] | ✓ | | partially |
| | VAE+ME[102] | ✓ | | |
| | DQN+SR[103] | | | high |
| | DORA[104] | | | high |
| | A2C+CoEX[52] | ✓ | | |
| | RND[51] | ✓ | | partially |
| | Action balance RND[105] | ✓ | | partially |
| | MADE[106] | ✓ | | |
| | Informed exploration[107] | ✓ | partially | |
| | EX$^2$[108] | | | high |
| | SFC[109] | ✓ | partially | partially |
| | CB[58] | ✓ | partially | high |
| | VSIMR[110] | ✓ | | high |
| | SMM[111] | ✓ | | |
| | DeepCS[112] | | | |
| | Novelty Search[113] | ✓ | | high |
| | ECO[114] | ✓ | partially | high |
| | NGU[53] | ✓ | high | |
| | NovelD[115] | ✓ | | high |
| Information gain | VIME[116] | ✓ | | high |
| | AKL[117] | ✓ | | high |
| | Disagreement[118] | ✓ | | high |
| | MAX[119] | ✓ | | high |

learns $\phi$ by a self-supervised inverse model using states pair $(s_t, s_{t+1})$ to predict the action $a_t$ done between them. Thus, $\phi$ ignores the uncontrollable aspects in the environment, so it can handle white noise problem to some degree. However, one major drawback is that it only considers the influence caused by one-step action, thus it cannot handle the long-horizon problem, for example, the consequences are caused by its action several steps later. Burda et al. [57] propose a detailed analysis of previous works that use different latent spaces to compute the prediction error. They show that using random features is sufficient for many typical RL environments, while its generalization is worse than learned features. Meanwhile, methods based on prediction errors are limited to deterministic environments and do not show much advantage in stochastic environments. Later, VDM [96] learns the stochasticity in the transition model through a variational dynamics model and measures the novelty through evidence lower bound.

Instead of learning a state representation using an encoder $\phi$, there are some works[97], [95] aiming to learn both state and action representations. Similar to ICM, AR4E [95] also learns the state transition model via a self-supervised reverse model. In addition, AR4E contains an action representation module that expands the input low-dimension actions to high-

dimension representations, and inputs the action representation into the dynamics model together with the current state. By increasing the representation power of the dynamics model, the results show improvements over ICM. EMI[97] is another work that both learns the state and action representations $\phi_s(s)$ and $\phi_a(a)$. EMI uses a linear forward model in the representation space which computes the prediction error of the dynamics model as an intrinsic reward. The state and action latent spaces are trained by maximizing the Mutual Information (MI) $I([\phi_s(s); \phi_a(a)]; \phi_s(s'))$ and $I([\phi_s(s); \phi_s(s')]; \phi_a(a))$ with the variational divergence lower bound of MI. EMI outperforms previous works on Atari and is capable of continuous control problems. However, it cannot address the white-noise problem since the state and action representation do not consider to remove the noisy distractions, as well as the long-horizon problem since it ignores the temporally-extended information.

*2) Novelty:* The second category focuses on motivating agents to approach states they have never visited (high novelty) by assigning agents intrinsic rewards as bonuses [120], [121]. Formally, this intrinsic reward is in inverse proportion to the visiting time of states $N(s_t)$: $R(s_t) = 1/N(s_t)$. These methods are called to be count-based. Nevertheless, it is hard to apply these methods to very large or continuous state spaces since an agent is impossible to cover the whole state space. To address this problem, Tang et al. [98] propose TRPO-AE-hash, which uses a hash function, SimHash [122] to discrete state space. However, the benefit of this method is marginal compared with those basic exploration policies. There are other attempts that have been proposed to deal with the large or continuous state space, like A3C+ [99] and DQN-PixelCNN [100], which rely on density models [123], [124]. Density models compute the pseudo-count $\hat{N}(s_t)$ [99], which is defined as the generalised visit-count: $\hat{N}(s_t) = \rho(s_t)(1 - \rho'(s_t))/(\rho'(s_t) - \rho(s_t))$, where $\rho(s_t)$ is the density model which produces the probability of observing $s_t$, and $\rho'(s_t)$ is the probability of observing $s_t$ after a new occurrence of $s_t$. Although these methods that are based on the pseudo-count work well in environments with sparse, delayed rewards, extra computational complexity is caused by estimating the density model[100].

In order to reduce the computational costs, $\phi$-EB [101] models the density on a latent space rather than on the raw state space. Their results on Montezuma's revenge show significant advantages considering the decrease in computational costs. Other indirect count-based methods are proposed, e.g., DQN+SR [103] uses the norm of the successor representation [83] as the intrinsic reward. DORA [104] proposes a generalization of counters, called $E$-values, that can be used to evaluate the propagating exploratory value over trajectories. This approach is applicable to the continuous space, and it shows superior performance compared with density models. However, DORA lacks experiments comparing it with previous approaches [99], [100], which is not compelling. Choi et al. [52] propose Attentive Dynamics Model (ADM) to discover the contingent region for state representation and exploration purposes. The intrinsic reward is defined as the visit count of the state representations consisting of the contingent region. Experiments on Montezuma's Revenge show ADM successfully extracts the location of the character and achieves a high score

of 11618 combined with PPO [29].

RND [51] estimates the state novelty by distilling a fixed random network (target network) into another network (predictor network). For each state, the target network produces random features of the state. The predictor network reconstructs the output of the target network for each state and is trained by minimizing the prediction error. The intrinsic reward is set as the prediction error. However, RND cannot handle the long-horizon problem because it only counts the visit times of each state and ignores the temporally-extended information. Furthermore, random features may be insufficient to represent the environment. Later, Song et al. [105] propose the action balance exploration strategy, which is based on RND and concentrates on finding unknown states. This approach aims to balance the frequency of each action selection to prevent an agent from paying too much attention to individual actions, thus encouraging the agent to visit unknown states. The action balance exploration achieves the highest score on Montezuma's revenge. However, like RND, it cannot handle the long-horizon problem since it also ignores the temporally-extended information. Recently, Zhang et al. [106] propose a new exploration method, called MADE, which maximizes the deviation of the occupancy of the policy from explored regions. This term is added as a regularizer to the original RL objective, which results in an intrinsic reward which can be incorporated to improve the performance of RL algorithms.

Apart from count-based methods, another way to estimate the state novelty is to measure the distance between the current state $s_t$ and states usually visited: $R(s_t) = \mathbb{E}_{s' \in \mathcal{B}}[\text{dist}(s_t; s')]$, where $\text{dist}(\cdot)$ is a distance measurement function and $\mathcal{B}$ is a distribution over recently visited states. Informed exploration [107] adopts a forward model to predict the action that leads the agent to the least often visited state in the last $d$ time steps. Then, an informed exploration strategy is proposed following the $\epsilon$-greedy strategy, which has a probability of $\epsilon$ to select the exploratory action that leads to the least often visited state in the last $d$ time step, rather than random actions. Later, to improve the efficiency of training the dynamics model, EX$^2$ [108] learns a classifier to differentiate each visited state from each other: when the classifier is not able to discriminate the current state against those in a buffer, this means the state has not been visited enough and the agent will be given a bonus for visiting this state, and vice versa. Successor Feature Control (SFC) [109] is a new kind of intrinsic reward, which takes statistics over trajectories into account, differing from previous works that use local information only to evaluate the intrinsic motivation. SFC can find the environmental bottlenecks where the dynamics change a lot and encourage the agent to go through these bottlenecks using intrinsic rewards.

CB [58] learns compressive state representations by maximizing the MI between states and latent variables via variational information bottleneck[125]. The intrinsic reward is defined as the KL-divergence between a fixed Gaussian prior and the posterior distribution of latent variables. CB ignores the task-irrelevant information when learning the state representations, thus it addresses the white-noise problem to a high degree. However, it requires the extrinsic reward to handle the stochasticity. Similarly, VSIMR [110] also adopts a KL-divergence

intrinsic reward, but uses a variational auto-encoder (VAE) to learn the latent space. State Marginal Matching (SMM) [111] is another method using a KL-divergence intrinsic reward. It computes the KL-divergence between the state distribution derived by the policy and a target uniform distribution which in fact, is equal to maximizing the state entropy. Lastly, Tao et al. [113] propose a new kind of intrinsic rewards based on the distance between each state and its nearest neighboring states in the low dimensional feature space to solve tasks with sparse rewards. However, using low-dimensional features may lose information that is crucial for the exploration of the entire state space, which restricts its generalization to complex scenarios.

All the above methods can be classified as the inter-episode novelty, where the novelty of each state is considered from the perspective of across episodes. In contrast, the intra-episode novelty resets the state novelty at the beginning of each episode. This kind of novelty encourages the agent to visit more different states within an episode. Stanton and Clune [112] firstly distinguish the cross-training novelty (inter-episode novelty) and intra-life novelty (intra-episode novelty). They introduce Deep Curiosity Search (DeepCS) to improve the intra-episode exploration by encouraging agents to visit more different states within an episode. The intrinsic reward is binary, which is set as 1 for unexplored states and 0 otherwise. The episodic curiosity module (ECO) [114] uses an episodic memory to form such a novelty bonus. To compute the state's intra-episode novelty, ECO compares each state with states in the memory. If from the current state, the agent takes too many steps to reach those states contained in the episodic memory, the agent is rewarded a bonus. Therefore, it computes the probability that the current state and each state in the episodic memory are reachable within $k$ steps as an intrinsic reward. Besides, ECO stores states with larger intrinsic rewards than a threshold, which means the stored states are more unreachable than other states, like bottleneck states. Thus, in fact, ECO encourages the agent to visit bottleneck states more. However, the intrinsic reward calculation requires the agent to compare the current state with each state in the episodic memory, which restricts its scalability to a very large state space since it requires more effort to find all bottleneck states.

Further, *Never Give Up* (NGU) [53] proposes a new intrinsic reward mechanism that combines both episodic novelty and life-long novelty. An episodic intrinsic reward is constructed by using k-nearest neighbors over the visited states stored in the episodic memory; the life-long novelty is driven by RND [51] error and multiplicatively modulates the episodic similarity signal. Such an episodic novelty encourages the agent to visit as much as possible states within the current episode no matter how often the state is visited previously, while the life-long novelty serves as a global measure with regard to the whole learning process. This kind of novelty shows generalization among complex tasks and allows both intra- and cross-episode temporally-extended exploration. More recently, Zhang et al. [115] propose a new criterion called NovelD which assigns intrinsic rewards to states at the boundary between already explored and unexplored regions. To avoid the agent to exploit the intrinsic reward by going back and forth between novel states and previous states, NovelD considers the episodic

restriction that the agent is only rewarded for its first visit to the novel state in an episode. Their results show NovelD outperforms SOTA in Mini-Grid.

*3) Information Gain:* The last class of work aims to lead the agents towards unknown areas, as well as to prevent agents from paying much attention to stochastic areas. This is achieved by using the information gain as an intrinsic reward. This reward is computed based on the decrease in the uncertainty about environment dynamics[126], which can also be assimilated to Bayesian surprise [127] or the learning progress [126], [128]. If the environment is deterministic, the transitions are predictable, so the uncertainty of dynamics can be decreased. On the contrary, when faced with stochastic environments, the agent is hardly capable to predict dynamics accurately, thus cannot reduce the uncertainty, which implicitly restricts this kind of methods. Specifically, the intrinsic reward based on information gain is defined as:

$$R(s_t, s_{t+k}) = \text{Uncertainty}_{t+k}(\theta) - \text{Uncertainty}_t(\theta) \quad (10)$$

where $\theta$ denotes the parameter of a dynamics model, and Uncertainty refers to the model uncertainty, which can be estimated in different ways as described in Sec. IV-A.

Variational Information Maximizing Exploration (VIME) [116] is an exploration strategy that encourages the agent to take actions that maximize the information gain about its belief of environment dynamics. VIME measures the information gain using variational inference, where the dynamics are approximated using a Bayesian neural network (BNN)[129]. Thus, the reward is computed as the uncertainty reduction on weights of BNN. However, using a BNN as the dynamics model makes the VIME hard to apply to complex scenarios due to the high computation costs. Later, Achiam and Sastry [117] propose a more efficient way than VIME to learn the dynamics model, by replacing BNNs with a neural network followed by fully-factored Gaussian distributions. They design two kinds of rewards: the first one (NLL) is the cross-entropy, which approximates the KL-divergence of the true transition probabilities and the learned model. The second reward (AKL) is designed as the learning progress, which is the improvement of the prediction between the current time step $t$ and after $k$ improvements at $t + k$. This method is simpler compared with VIME, however, the benefit in performance is also marginal.

Pathak et al. [118] train an ensemble of dynamics models and use the mean of outputs as the final prediction. The intrinsic reward is designed as the variance over the ensemble of network output. The variance is high when dynamics models are not learned well, and low when the training process continues and all models will converge to the mean value finally, thus ignoring the noisy distractions since noises in the environment are task-irrelevant and do not affect the convergence. A similar idea is MAX[119], while it uses the JS-divergence instead of the variance over the outputs of dynamics models. These methods handle the white-noise problem to some degree (Table II) since the ensemble technique ignores the stochasticity. However, the main issue is the high computational complexity since it requires to train an ensemble of models.

## C. Other Advanced Methods for Exploration

The previous two sections introduce two main streams of exploration in DRL, i.e., uncertainty-oriented exploration and intrinsic motivation-oriented exploration. Next, we discuss several other branches of approaches that pursue efficient exploration from other perspectives which cannot be classified into the previous two main streams exactly. These methods provide different insights of how to achieve a general and effective exploration in DRL.

*1) Distributed Exploration:* One straightforward idea to improve exploration is distributed exploration: using heterogeneous actors with different exploration behaviors to discover the environment in a diverse manner. One representative work is Ape-X [130], in which a bunch of DQN workers perform $\epsilon$-greedy with different values of $\epsilon$ among independent environment instances in parallel. The independent randomness of different instances and the different exploration degree (i.e., $\epsilon$) of distributed workers allow efficient exploration of the environment regarding the wall time. Prioritized Experience Replay (PER) [27] is then applied to improve the learning efficiency from diverse experiences collected by different workers. Later after Ape-X, R2D2 [131] is proposed to further develop Ape-X architecture by integrating recurrent state information, making the learning more efficient.

Moreover, aim at solving hard-exploration games, distributed exploration is adopted to enhance the advanced exploration methods. Taking R2D2 as the base architecture, *Never Give Up* (NGU) [53] performs efficient exploration with the help of distributed agents, which are of different exploration degrees with regard to both episodic novelty and life-long novelty. Furthermore, Agent57 [7] improves NGU by adopting a separate parametrization of Q-functions and a meta-controller for adaptive selection of exploration policies. The former makes use of independent networks to adapt to the scale and variance of intrinsic and extrinsic reward, thus making the learning more stable. The later addresses the limitation of NGU where all exploration policies are trained in equivalence through learning to select exploration policies with the consideration of their contribution to the overall learning process. Both NGU and Agent57 achieve significant improvement in hard-exploration tasks of the Atari-57 suite while maintaining good performance across the remaining ones, which indicates the strong ability of advanced distributed methods in dealing with large state-action space and sparse reward.

*2) Exploration with Parametric Noise:* Another perspective to encourage exploration is to inject noise directly in the parameter space. Compared with naively adding noise in the outputs of the policy, parametric noise is used to perturb the policy parameters for more diverse and consistent exploration. Plappert et al. [132] propose to inject spherical Gaussian noise directly in policy parameters, i.e., the weights and biases of policy networks or value networks. The noise scale is adaptively adjusted according to a heuristic mechanism which depends on a distance measure in action space (e.g., KL divergence) between the perturbed and non-perturbed policy. Such an exploration technique is combined with both on-policy and off-policy algorithms (e.g., DQN, DDPG, and TRPO),

and empirically shows the benefits. Another similar work is NoisyNet [133] which also injects noise in network parameters for effective exploration. One thing that differs the most is that the noise scale (variance) in NoisyNet is learnable and updated from the RL loss function along with the other network parameters, in contrast to the heuristic adaptive mechanism in [132]. Parametric noise is an effective branch of methods to improve the exploration and learning performance in usual environments; however, it may not do a great favor in dealing with the exploration challenges like sparse, delayed rewards.

Beyond the above two branches, we introduce several other remarkable works with different exploration ideas. Arguably, Go-Explore [134], [135] may be the most powerful methods to solve Montezuma's Revenge and Pitfall, the most notorious hard-exploration problems in 57 Atari games. The recipe of Go-Explore is *return-then-explore*: policy first arrives at interesting states and then explores with higher efficiency. The interesting states are selected heuristically (e.g., count subscore) from the state archive (similar to the episodic memory); the arrival of the selected states can be achieved through resetting the simulation or performing a goal-conditioned policy which is trained additionally. After the arrival, random exploration is conducted and the archive is updated with newly visited states. Finally, a robustification phase is carried out to learn a robust policy from high-performing trajectories via Learning from Demonstration (LfD) [136]. Go-Explore demonstrates its amazing ability in solving large-state space, sparse reward, and extremely long-horizon problems. However, the superiority of Go-Explore comes much from sophisticated and specific designs, and it may not be a general method for other hard-exploration problems. Another different work is Potentialized Experience Replay (PotER) [137]. The key idea of PotER is the introduction of Artificial Potential Field (APF). Through defining a potential energy function for each state in experience replay, including an attractive potential energy that encourages the agent to be far away from the initial state and a repulsive one that prevents from going towards obstacles, i.e., death states. This allows the agent to learn from both superior and inferior experiences using intrinsic potential signals. Although PotER also relies on task-specific designs, potential-based exploration is seldom studied in DRL, thus is worthwhile further study.

There are also many other remotely related works that demonstrate the efficacy in exploration, even though some of them are not proposed to be intended for efficient exploration. One representative and popular research domain is Skill Discovery [138], [139], [140], [141]. The skill discovery methods aim to learn skills in an unsupervised fashion by exploring the environment without any extrinsic rewards. The basic principle of skill discovery is *Empowerment* [142], which describes what the agent can be done while learning how to do it. Intuitively, the empowerment principle encourages diverse skills regarding the coverage of state spaceresulting in efficient exploration in state space.

## V. EXPLORATION IN DEEP MARL

After investigating the exploration methods for single-agent DRL, in this section, we move to multi-agent exploration

methods. At present, the study on exploration for deep MARL is roughly at the preliminary stage. Most of these exploration methods extend the ideas in the single-agent setting and propose different mechanisms by integrating the characteristics of deep MARL. Recall the challenges faced by MARL exploration (Sec. III), the dimensionality of state-action space increases rapidly with the growth of the number of agents. As a consequence, this scales up the difficulty of quantifying the uncertainty and computing various forms of intrinsic motivation. Beyond the large exploration space, the multi-agent interaction is another critical thing, which can be considered from several aspects: 1) multi-agent environments usually involve partial observations and non-stationary dynamics, thus cooperative and coordinated behaviors among agents are nontrivial to achieve; 2) multiple agents jointly influence the environmental dynamics and often share reward signals, raising a challenge in reasoning and inferring the effects of joint exploration consequences, and coordinating the individual exploration behaviors; 3) the agents are expected to perform coordinated exploration, where local (individual) and global (joint) exploration and exploitation are to be dealt with properly; 4) multi-agent interactions induce the mutual influence among agents, providing richer reward-agnostic information which can be utilized for exploration. Next, we investigate these methods by following the similar taxonomy (shown in Fig. 1) adopted in single-agent cases, i.e., uncertainty-oriented exploration, intrinsic motivation-oriented exploration and other advanced exploration techniques.

### A. Uncertainty-oriented Exploration

In the multi-agent domain, estimating the uncertainty is difficult since the joint state-action space is significantly larger than single-agent domain. Meanwhile, quantifying the uncertainty has special difficulties according to the partial observation and non-stationary problems. 1) Each agent only draws a local observation from the joint state space. In such a case, the uncertainty measurement is a kind of 'local' uncertainty. Exploration based on the local uncertainty is unreliable since the estimation is biased and cannot reflect the global information of the environment. 2) The multi-agent domain is non-stationary since an agent's behaviors are influenced by other coexisting agents. This problem also leads to a 'local' uncertainty measurement since an agent cannot obtain other agents' policies [143]. Both problems increase the randomness of the environment and make the uncertainty estimation difficult in such stochastic environments. Meanwhile, the agent should balance the local and global uncertainty to explore the novel states concerning the local information, and also avoid duplicate exploration by considering the other cooperative agents' uncertainty.

The epistemic uncertainty-based approach can be directed extended to the multi-agent problem. Following the OFU principle, Zhu et al.[144] propose Multi-agent Safe Q-Learning that uses the epistemic uncertainty for exploration. This method assumes each agent observes the joint state space of all agents. For uncertainty estimation, they use the Gaussian process to represent the posterior of the $Q$-value function as Eq. (7), and obtain an upper bound of the $Q$-value function as $Q^+$, in which

the variance of the Gaussian process portrays the epistemic uncertainty. Then the agent follows a Boltzmann policy to explore according to $Q^+$. To overcome the non-stationary problems, they further constrain the actions to ensure the low risk through the joint unsafety measurement.

There exist several methods that use both the epistemic uncertainty and aleatoric uncertainty in multi-agent exploration, where the use of aleatoric uncertainty addresses the stochasticity in the value function. Martin et al. [145] measure the aleatoric uncertainty following the similar idea of distributional value estimation in single-agent RL [24], and perform exploration based on the distributional value function as well as the epistemic uncertainty. Specifically, they extend several single-agent exploration methods [74], [146], [147] to zero-sum stochastic games, and find the most effective approaches among them are Thomson sampling and Bayes-UCB based methods. The Thompson sampling-based method samples from the posterior of value function, and the Bayes-UCB based method extends the Bayes-UCB [147] to zero-sum stochastic form game. It samples several payoff matrices from the posterior distribution and chooses the action with the highest mean payoff quantile. Both strategies maintain a posterior distribution to measure the epistemic uncertainty, thus they can perform *deep* exploration. In addition, QR-MIX [148] combines QMIX [37] with IQN [149] through modeling the joint state-action values as a distribution. Based on the return likelihoods, LH-IRQN [150] introduces a decentralized quantile estimator to distinguish the non-stationary samples. These methods consider aleatoric uncertainty only for better value estimation, nevertheless, it is more desirable to use the aleatoric uncertainty to improve the robustness in exploration through explicitly avoid the agent to explore states with high aleatoric uncertainty.

### B. Intrinsic motivation-oriented Exploration

Intrinsic motivation is widely used as the basis of exploration bonuses to encourage agents to explore unseen regions. Following the great success in single-agent RL, a number of works [151], [152] tend to apply intrinsic motivation in multi-agent domain. However, the difficulty of measuring the intrinsic motivation increases exponentially with the number of agents increasing. Furthermore, assigning intrinsic rewards to agents has special difficulties according to the partial observation and non-stationary problems. 1) Each agent only draws a local observation from the joint state space. In such a case, the intrinsic motivation (such as novelty) measurement is a kind of 'local' intrinsic motivation, which is unreliable since the estimation is biased and cannot reflect the global information of the environment. 2) The multi-agent domain is non-stationary since an agent's behaviors are influenced by other coexisting agents. This problem also leads to a 'local' intrinsic motivation measurement since an agent usually does not know the policy of other agents [143]. 3) Agents are expected to perform coordinated exploration, where local (individual) and global (joint) intrinsic motivation are to be dealt with properly and consistently. Meanwhile, multi-agent interactions induce the mutual influence among agents, providing the reward-agnostic information which can be utilized. In this section,

we investigate previous methods that aim to apply intrinsic motivation in multi-agent domains. Some works have addressed some of the above challenges to a certain degree, such as coordinated exploration[153], [154]. Nevertheless, there still exist open questions, such that how to decrease the difficulties in estimating the intrinsic motivation in large-scale multi-agent systems, and how to balance the inconsistency between local and global intrinsic motivation.

Some works assign agents extra bonuses based on novelty to encourage exploration. For example, Wendelin Böhmer et al. [151] introduce an intrinsically rewarded centralized agent to interact with the environment and store experiences in a shared replay buffer while decentralized agents update their local policies using the experience from this buffer. They find that although only the centralized agent is rewarded, decentralized agents can still benefit from this and improve exploration efficiency. Instead of simply designing intrinsic rewards according to global states' novelty, Iqbal and Sha [152] define several types of intrinsic rewards by combining decentralized curiosity of each agent. One type is selected dynamically during each episode which is controlled by a meta-policy. However, these types of intrinsic rewards are domain-specific and cannot be extended to other scenarios.

Instead of designing intrinsic rewards based on state novelty, Learning Individual Intrinsic Reward (LIIR) is proposed [153] which learns the individual intrinsic reward and uses it to update an agent's policy with the objective of maximizing the team reward. LIIR extends a similar idea in single-agent domains[155], [156] that learns an extra proxy critic for each agent, with the input of intrinsic rewards and extrinsic rewards. The intrinsic rewards are learned by building the connection between the parameters of intrinsic rewards and the centralized critic based on the chain rule, thus achieving the objective of maximizing the team reward. Jaques et al. [157] define the intrinsic reward function from another perspective called "social influence", which measures the influence of one agent's actions on others' behavior. By maximizing this function, agents are encouraged to take actions with the most strong influence on the policies of other agents, those joint actions lead agents to behave cooperatively. Instead of using intrinsic rewards as in [157], Wang et al. integrate such influence as a regularizer into the learning objective[154]. They measure the influence of one agent on other agents' transition function (EITI) and rewarding structure (EDTI), and encourage agents to visit critical states in the state-action space (Fig. 6(a)), through where agents can transit to potentially important unexplored regions.

Chitnis et al. [158] tackle the coordinated exploration problem from a different view by considering that the environment dynamics caused by joint actions are different from that caused by individually sequential actions. Therefore, this method incentivizes agents to take joint actions whose effects cannot be achieved via a composition of the predicted effect into individual actions executed sequentially. However, they need manually modify the environment to disable some of the agents, which is unrealistic for real-world scenarios. To summarize, the research of the intrinsic motivation-oriented exploration in MARL is mainly extended from single-agent domains, like state novelty estimation and so on. Meanwhile, there exist some works that leverage the mutual influence among agents to guide coordinated exploration. However, with the inconsistency between local and global information, how to obtain a robust and accurate intrinsic motivation estimation that balances the local and global information to derive coordinated exploration is a promising direction that needs to be studied.

### C. Others Methods for Multi-agent Exploration

Different from the works introduced previously which are derived from uncertainty estimation or intrinsic motivation, there are a few notable multi-agent exploration works which cannot be classified into the former two categories. In the following, we introduce these works in chronological order.

In the literature of multi-agent repeated matrix games and multi-agent MAB problems, the previous works propose exploration methods from different perspectives [159], [160], [161], [162]. These works are studied in environments with relatively small state-action space, which is far from the scales often considered in deep MARL. Towards complex multi-agent environments with large state-action space, Dimakopoulou and Van Roy [163], [164] identify three essential properties for efficient coordinated exploration: adaptivity, commitment, and diversity. They present the failures of straightforward extensions of single-agent posterior sampling approaches in satisfying the above properties. For a practical method to meet these properties, they propose a Seed Sampling by incorporating randomized value networks for scalability.

From different angles, Chen [165] proposes a new framework to address the coordinated exploration problem under the paradigm of CTDE. The key idea of the framework is to train a centralized policy first and then derive decentralized policies via policy distillation. Efficient exploration and learning are conducted by optimizing a global maximum entropy RL objective with global information. Decentralized policies distill cooperative behaviors from the centralized policy with the favor of agent-to-agent communication protocols. Another notable exploration method that follows CTDE paradigm is MAVEN [166]. MAVEN utilizes a hierarchical policy to control a shared latent variable as a signal of a coordinated exploration mode, in which the value-based agents condition their policies. Through maximizing the mutual information between the latent variable and the induced episodic trajectory, MAVEN achieves diverse and temporally-extended exploration. Both the above two methods leverage the centralized training mechanism to enable coordinated exploration and they demonstrate the potential of such exploration improvements in benefiting deep MARL.

## VI. DISCUSSION

### A. Empirical Analysis

For a unified empirical evaluation of different exploration methods, we summarize the experimental results of some representative methods on three representative benchmarks: Montezuma's Revenge, the overall Atari suit, and Vizdoom. Each of the three benchmarks has different characteristics and evaluation focus on different exploration challenges.

Recall the introduction in Sec. III, Montezuma's Revenge is notorious due to its sparse, delayed rewards and long horizon.

TABLE III: A benchmark of experimental results of exploration methods in DRL. The results are from those reported in their original papers.

| Benchmark scenarios | Method | Basic Algorithm | Convergence Time (frames) | | | Convergence Return | | |
|---|---|---|---|---|---|---|---|---|
| Montezuma's Revenge | Successor Uncertainty [62] | DQN [21] | 200M | | | 0 | | |
| | Bootstrapped DQN [65] | Double DQN [22] | 20M | | | 100 | | |
| | Randomized Prior Functions [50] | Bootstrapped DQN [65] | 200M | | | 2500 | | |
| | Uncertainty Bellman Equation [64] | DQN [21] | 500M | | | $\sim$2750 | | |
| | IDS [90] | Bootstrapped DQN [65] & C51 [24] | 200M | | | 0 | | |
| | DLTV [71] | QR-DQN [25] | 40M | | | 187.5 | | |
| | EMI [97] | TRPO [167] | 50M | | | 387 | | |
| | A2C+CoEX[52] | A2C [40] | 400M | | | 6635 | | |
| | RND [51] | PPO [29] | 1.6B | | | 8152 | | |
| | Action balance RND [105] | RND [51] | 20M | | | 4864 | | |
| | PotER [137] | SIL [168] | 50M | | | 6439 | | |
| | NGU [53] | R2D2 [131] | 35B | | | 10400 | | |
| | Agent57 [7] | NGU [53] | 35B | | | 9300 | | |
| | Go-Explore (no domain knowledge) [134] | PPO [29] + Backward Algorithm [169] | 5.55B (Hypothetically[1]over 35B) | | | 43763 | | |
| | Go-Explore (domain knowledge) [134] | PPO [29] + Backward Algorithm [169] | 5.2B (Hypothetically over 150B) | | | 666474 | | |
| | Go-Explore (best)[2] [134] | PPO [29] + Backward Algorithm [169] | 5.2B (Hypothetically over 150B) | | | 18003200 | | |
| Overall Atari Suit | Bootstrapped DQN [65] | Double DQN [22] | 200M (55 games) | | | 553%(mean), 139%(median) | | |
| | Uncertainty Bellman Equation [64] (1-step) | DQN [21] | 200M (57 games) | | | 776%(mean), 95%(median) | | |
| | Uncertainty Bellman Equation [64] (n-step) | DQN [21] | 200M (57 games) | | | 440%(mean), 126%(median) | | |
| | A3C+ [99] | A3C[40] | 200M (57 games) | | | 273%(mean), 81%(median) | | |
| | Noisy-Net [133] | DQN [21] | 200M (55 games) | | | 389%(mean), 123%(median) | | |
| | Noisy-Net [133] | Dueling DQN [26] | 200M (55 games) | | | 608%(mean), 172%(median) | | |
| | IDS [90] | Bootstrapped DQN [65] | 200M (55 games) | | | 651%(mean), 172%(median) | | |
| | IDS [90] | Bootstrapped DQN [65] & C51 [24] | 200M (55 games) | | | 1058%(mean), 253%(median) | | |
| | Randomized Prior Functions [50] | Bootstrapped DQN [65] | 200M (55 games) | | | 444%(mean), 124%(median) | | |
| | Randomized Prior Functions [50] | Bootstrapped DQN [65]+Dueling [26] | 200M (55 games) | | | 608%(mean), 172%(median) | | |
| | NGU [53] | R2D2 [131] | 35B (57 games) | | | 3421%(mean), 1354%(median) | | |
| | Agent57 [7] | NGU [53] | 35B (57 games) | | | 4766%(mean), 1933%(median) | | |
| Vizdoom (MyWayHome) | | | Dense | Sparse | Very sparse | Dense | Sparse | Very sparse |
| | ICM [56] | A3C [40] | 300M | 500M | 700M | 1.0 | 1.0 | 0.8 |
| | AR4E [95] | ICM [56] | 300M | 400M | 600M | 1.0 | 1.0 | 1.0 |
| | SFC [109] | Ape-X DQN [130] | 250M | - | - | 1.0 | - | - |
| | EX$^2$[108] | TRPO [167] | 200M | - | - | 0.8 | - | - |
| | ECO[114] | PPO [29] | 100M | 100M | 100M | 1.0 | 1.0 | 1.0 |

[1] The hypothetical frame denotes the number of game frames that would have been played if Go-Explore [134] replayed trajectories instead of resetting the emulator state as done in their original paper.

[2] Removing the maximum limit of 400,000 game frames imposed by default in OpenAI Gym, the best single run of Go-Explore with domain knowledge achieved a score of 18,003,200 and solved 1441 levels during 6,198,985 game frames, corresponding to 28.7 hours of game play (at 60 game frames per second, Atari's original speed) before losing all its lives.

It requires agents to have a strong exploration ability to obtain positive feedback; while to traverse multiple rooms and achieve a high score further needs human-level memory and control of events in the environment. On the contrary, the overall Atari suite focuses on a more general evaluation of exploration methods in improving the learning performance of RL agents. Vizdoom is another representative task with multiple reward configurations (from dense to very sparse). Distinct from the previous two tasks, Vizdoom is a navigation (and shooting) game with the first-person view. This simulates a learning environment with severe partial observability and underlying spatial structure, which is more similar to real-world ones faced by humans. We collect the reported results on the three tasks from the original papers. Thus, the exploration methods without such results are not included. We do not provide a benchmark for the exploration methods in MARL because there is hardly any environment commonly used. This is also considered as a desideratum for the study on MARL exploration.

The results are shown in Table III, in terms of the exploration method, the basic RL algorithm, the convergence time and return. The table provides a quick glimpse to the performance comparison: 1) For Montezuma's Revenge, Go-Explore [134], NGU [53] and RND [51] outperforms human-level performance. Go-Explore achieves the best results while it needs the environment to be deterministic and re-settable in its exploration stage, which can be unrealistic in many practical problems. 2)

For the overall Atari suit with 200M training frames, IDS [90] achieves the best results and outperforms its basic methods (i.e., Bootstrapped DQN) significantly, showing the success of the sophisticated utilization of uncertainty in improving exploration and learning generally. Armed with a more advanced distributed architecture like R2D2 [131], Agent57 [7] and NGU [53] achieve extremely high scores in 35B training frames. 3) For Vizdoom (MyWayHome), ECO [114] outperforms other exploration methods especially in achieving higher sample efficiency across all reward settings. This demonstrates the effectiveness of intrinsic motivation for another time and also reveals the potential of episodic memory in dealing with partial observability and spatial structure of navigation tasks.

We conclude the results of the exploration methods as follows:

- For uncertainty-oriented methods, the results show that they achieve better results on the overall Atari suit, demonstrating the effectiveness in general cases. Meanwhile, their performance in Montezuma's Revenge is generally lower than intrinsic motivation-oriented methods. This is because the uncertainty quantification is mainly based on well-learned value functions, which are hard to learn if the extrinsic rewards are almost absent. In principle, the effects of uncertainty-oriented exploration methods heavily rely on the quality of uncertainty estimates. As a consequence, more accurate and robust uncertainty estimates

are expected to achieve among complex environments.

- The intrinsic motivation-oriented methods usually focus on hard-exploration tasks like Montezuma's Revenge. For example, RND [51] outperforms human-level performance by using intrinsic rewards to constantly look for novel states to help the agent pass more room. However, according to a recent empirical study [170], although the existing methods greatly improve the performance in several hard-exploration tasks, they may have no positive effect on or even hinder the learning performance in other tasks. This can be attributed to the introduction of intrinsic motivation, which often alters the original learning objective and may deviate from the optimal policies. This raises a requirement for improving the versatility of intrinsic motivation-oriented methods. A preliminary success is achieved by NGU [53] that learns a series of $Q$ functions corresponding to different coefficients of the intrinsic reward, among which the exploitative $Q$ function is ready for execution. Such a technique is further improved by separate parameterization in Agent57 [7].

- Other advanced exploration methods also pursue efficient exploration from different perspectives. 1) Distributed training greatly improves the exploration and learning performance generally. Distributed training is one of the key components of Agent57 [7], the first DRL agent that achieves superhuman performance on all 57 Atari games. 2) $Q$-network with parametric noise brings stable performance improvement compared to a deterministic $Q$-network. For an instance, Noisy-Net [133] significantly outperforms the intrinsic motivation-oriented methods evaluated by the overall Atari suit. 3) Another notable concept is potential-based exploration. As a representative, PotER [137] achieves good performance in Montezuma's Revenge with significantly higher sample efficiency than other methods, making potential-based exploration promising in addressing hard-exploration environments.

### B. Open Problems

Although encouraging progress has been achieved, efficient exploration remains a challenging problem for DRL and deep MARL. Moreover, we discuss several open problems which are fundamental yet not well addressed by existing methods, and point out a few potential solutions and directions.

**Exploration in Large State-action Space.** The difficulty of exploration escalates as the growth of scale and complexity of state-action space. To deal with large state space, exploration methods often need high-capacity neural network to measure the uncertainty and novelty. For representative uncertainty-oriented methods, Bootstrapped DQN [65], OAC [50] and IDS [70] use Bayesian network to approximate the posterior of $Q$-function, which is computationally intensive. Theoretically, accurately estimating the $Q$-posterior in large state space requires infinite bootstrapped $Q$-networks, which is obviously infeasible in practice, and in stead an ensemble of 10 $Q$-networks is often used. Intrinsic motivation-oriented methods often need additional auxiliary models such as forward dynamics [118] and density model [100] to measure the novelty of states or

transitions. For intrinsic motivation-oriented methods, learning accurate density estimation and effective auxiliary models such as forward dynamics in large state space is also nontrivial within a practical budget. The consequent quality of uncertainty estimation and intrinsic guidance achieved thus in turn affects the exploration performance. Another major limitation of existing works is the incapability of learning and exploring in large and complex action spaces. Most methods consider a relatively small discrete action space or low-dimensional continuous action space. Nevertheless, in many real-world scenarios the action space can consist of a large number of discrete actions or has a complex underlying structure such as a hybrid of discrete and continuous action spaces. Conventional DRL algorithms have scalability issues and even are infeasible to be applied. A few recent works attempt to deal with large and complex action spaces in different ways. For example, [47] proposes to learn a compact action representation of large discrete action spaces and convert the original policy learning into a low-dimensional space. Another idea is to factorize the large action space [171], e.g., into a hierarchy of small action space with different abstraction levels. Besides, for structured action space like discrete-continuous hybrid action space, several works are proposed with sophisticated algorithms [46], [149], [172], [48]. Despite the attempts made in aforementioned works, how efficient exploration can be realized with a large and complex action space remains unclear. Overall, it remains an open problem to RL community on how to develop efficient methods for uncertainty and novelty estimation in large state space and large, complex action space.

Since the main challenge is the large-scale and complex structure of state-action spaces, a natural solution is to construct an abstract and well-behaved space as a surrogate, among which exploration can be conducted efficiently. Thus, one promising way is to leverage the representation learning of states and actions. The potential of representation learning in RL has been demonstrated by some recent works in improving policy performance in environments with image states [173] and hybrid actions [48], as well as generalization across multiple tasks [174]. The representations in these works are learned by following specific criteria, e.g., reconstruction, instance-discriminative contrast and dynamics prediction. However, how an exploration-oriented state and action representations can be obtained is unclear yet. To our knowledge, some efforts have been made in this direction. For example, DB [175] learns dynamics-relevant representations for exploration through the information bottleneck. SSR [103] makes use of successor representation upon which count-based exploration is then performed. The central problem is what exploration-favorable information should be retained by the representation to learn. To fulfill exploration-oriented state representation, we consider that the information of both the environment to explore and the agent's current knowledge are pivotal. One feasible approach of leveraging useful environment information is learning state abstraction based on the topology of the state space with actionable connectivity [176], thus unnecessary exploration of redundant states can be avoided. Taking into account agent's current knowledge, a further abstraction of state space can be achieved for more efficient exploration. A potential way is

establishing an equivalence relation based on the familiarity of states, following which the boundary of exploration can be characterized. In this manner, we expect highly targeted and efficient exploration can be realized. For action representation, one key point may be the utilization of action semantics, i.e., how the action affects the environment especially on the critical states. The similarity of action semantics between actions enables effective generalization of learned knowledge, e.g., the value estimate of an action can be generalized to other actions that have similar impacts on the environment. At the same time, the distinction of action semantics can be made use of to select potential actions to seek for novel information.

**Exploration in Long-horizon Environments Extremely Sparse, Delayed Rewards.** For exploration in environments with sparse, delayed rewards, some promising results have been achieved by a few exploration methods from the perspective of intrinsic motivation [51], [105] or uncertainty guidance [50], [64]. However, most current methods also reveal their incapability when dealing with sparser or extremely sparse rewards, typically in an environment with a long horizon. As a typical example, the whole game of Montezuma's Revenge has not been solved by DRL agents except for Agent57 [7] (although partially) and Go-Explore [134], [135]. However, Agent57 achieves over 9.3k scores by taking advantage of a two-scale novelty mechanism and an advanced distributed architecture base (i.e., R2D2 [131]); Go-Explore [134], [135] achieves superhuman performance through imitating superior trajectory experiences collected in a return-and-explore fashion with the access to controlling the simulation. These methods are far from sample-efficient and highly customed, thus lack of generality. For real-world navigation scenarios, the practical methods often need to combine prior knowledge and intrinsic motivation to perform exploration in a long horizon. In such environments, the prior knowledge usually includes representative landmarks [177] and topological graphs [178] which are designated by human. Apparently, current exploration methods for navigation rely heavily on prior knowledge. However, this is often expensive and nontrivial to obtain in general.

Overall, learning in such an environment with extremely sparse, delayed rewards is an unsolved problem at present, which is of significance to developing RL towards practical applications. Intuitively, beyond sparse and delayed rewards, solving such problems involves higher-level requirements on long-time memorization of environmental context and versatile control of complex environmental semantics. These aspects can be the desiderata of effective exploration methods in the future study. To take a step towards this direction, long-time memorization of the environmental context may be realized by more sophisticated models, e.g., Transformer [179] and Episodic Memory, in an implicit or explicit manner. Another promising yet challenging solution is establishing a universal approach to extract the hierarchical structure of different environments. Learning an ensemble of sub-policies (i.e., skills) is a possible way to fulfilling the versatile control of the environment, based on which temporal abstracted exploration can be performed in higher levels. In addition, incorporating general-form prior knowledge to reduce unnecessary exploration is also a promising perspective.

**Exploration with White-noise Problem.** The stochasticity inherent in dynamics or manually injected in environments usually distracts agents in exploration. Several works like RND [51], ICM [56], and EMI [97] all focus on solving state-space noise by constructing a compact feature representation to discard task-irrelevant features. Count-based exploration handles stochasticity in state-space through attention [52] and state-space VAE [110]. The state representation discards the task-irrelevant information in exploration, which is promising to overcome the white-noise problem. However, most methods require a dynamics model or state-encoder, which increases the computational cost. Although CB[58] does not learn a state-encoder, it needs environmental rewards to remove the noisy information thus cannot work in extremely sparse reward tasks. Other methods to solve this problem include Bayesian disagreement [118] and active world model learning [180]. They use the Bayesian uncertainty and information gain-based model that are insensitive to white-noise to overcome the state-space noise. Although they are promising to handle the state space noise, the action-space noise has not been rigorously discussed in the community. The sticky Atari [181] injects action-space noise in discrete action space, while it is hard to design and represent the realistic noise in real-world applications. How to construct a more realistic distraction is a worth exploring direction in the future. One possible direction in solving the white noise problem is using adversarial training to learn a robust policy. Specifically, we can find the adversarial examples by solving a minmax problem, and then reduce the sensibility of the model to the adversaries. Another promising direction is using direct regularization. For example, enforcing Lipschitz constraints for policy networks to make the predictions provably robust to small perturbations produced by noise; using bisimulation constraints to reduce the effects of noise by learning the dynamics-relevant representations.

**Convergence.** For uncertainty-oriented exploration, optimism and Thomson sampling-based methods need the uncertainty to converge to zero in the training process to make the policy converge. Theoretically, the epistemic uncertainty enables to converge to zero in tabular [182] and linear MDPs [183], [184] according to the theoretical results. In general, MDPs, as the agent learns more about the environment, the uncertainty that encourages exploration gradually decreases to zero, then the confidence set of the MDP posterior will contain the true MDP with a high probability [185], [186]. However, due to the curse of dimensionality, the practical uncertainty estimation usually relies on function approximation like neural networks, which makes the estimation errors hard to converge. Meanwhile, the bootstrapped sampling [65], linear final layer approximation [61] or variational inference [133], [187] only provide the approximated Bayesian posterior rather than the accurate posterior, which may be hard to represent the true confidence of the value function with the changing experiences in exploration. How to build efficient uncertainty measurements and posterior approximation is still an open problem, both in theory and in practice. The recently proposed Single Model Uncertainty [188] is a kind of uncertainty measurement through a single model, which may provide more stable rewards in exploration. Combining out-of-distribution detection methods

[189] with RL exploration is also a promising direction.

For intrinsic motivation-oriented exploration, a regular operation is to add an intrinsic reward to the extrinsic reward, which is a kind of reward shaping mechanism. However, these reward transformation methods are usually heuristically designed while not theoretically justified. Unlike potential-based reward [190] that does not change the optimal policy, the reward transformations may hinder the performance and converge to the suboptimal solution. There are several efforts that propose novel strategies to combine extrinsic and intrinsic policy without reward transformation. For example, scheduled intrinsic drive [109] and MuleX [191] learn intrinsic and extrinsic policies independently, and then use a high-level scheduler or random heuristic to decide which one to use in each time step. This method improves the robustness of the combined policy while also does not make the theoretical guarantee of optimality. The two-objective optimization [192] calculates the gradient of intrinsic and extrinsic reward independently, and then combines the gradient directions by following Fig. 7. However, this technique does not work when the extrinsic rewards are almost absent. Recently, meta-learned rewards [155], [156] are proposed to learn an optimal intrinsic reward by following the gradient of extrinsic reward, which guarantees the optimality of intrinsic rewards. However, how to combine this method with existing bonus-based exploration remains an open question. Another promising research direction is designing intrinsic rewards following the potential-based reward shaping principle to ensure the optimality of the combined policy.
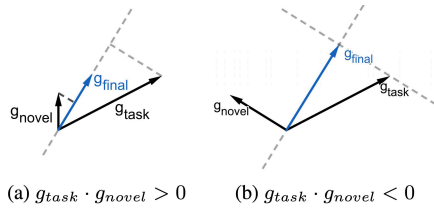


Fig. 7: The policy gradient of extrinsic rewards $g_{task}$ and intrinsic rewards $g_{novel}$ are calculated independently. (a) When $g_{task} \cdot g_{novel} > 0$, the final gradient is the angular bisector of the two gradients because they have similar directions. (b) When $g_{task} \cdot g_{novel} < 0$, the final gradient will rough follow the extrinsic gradient because the two gradients have large disagreement.

**Multi-agent Exploration.** The research of multi-agent exploration is still at an early stage and does not address most of the mentioned challenges well, such as partial observation, non-stationary, high-dimensional state-action space, and coordinated exploration. Since the joint state-action space grows exponentially as the increase in the number of agents, the difficulty in measuring the uncertainty or intrinsic motivation escalates significantly. Referring to the discussion of the large state-action space, the scalability problem in multi-agent exploration can also be alleviated by representation learning. Beyond the representation approaches in single-agent domain, one special and significant information that can be taken advantage of is to build a graph structure of multi-agent interactions and relations. It is of great potential while remains open to incorporate such graph structure information in state and action representation learning for effective abstraction and reduction of the joint space. Furthermore, with the inconsistency

between local and global information (detailed in Section III), how to obtain a robust and accurate uncertainty (or intrinsic motivation) estimation that balances the local and global information to derive a coordinated exploration is worthwhile to be studied. One promising solution is to design exploration strategies by taking credit assignment into consideration, this technique has been successfully applied to achieve multi-agent coordination[36], which may help to assign a reasonable intrinsic reward for each agent to solve this problem. Another problem is that, there is still no well-accepted benchmark yet, most of the previous works have designed specific test beds to verify the effectiveness of their proposed corresponding solutions, such as influence-based exploration strategies[154], [157] which are tested on environments that need strong cooperative and coordinated behaviors (shown in Fig 6(a)); novelty-based strategies [152], [151] are tested on environments where novel states are well correlated with improved rewards. Therefore, how to construct a well-accepted benchmark and derive a general exploration strategy that is suited for the benchmark remains an open problem in multi-agent exploration.

## VII. CONCLUSION

In this paper, we conduct a comprehensive survey on exploration in DRL and deep MARL. We identity the major challenges of exploration for both single-agent and multi-agent RL. We investigate previous efforts on addressing these challenges by a taxonomy consists of two major categories: uncertainty-oriented exploration and intrinsic motivation-oriented exploration. Besides, some other advanced exploration methods with distinct ideas are also concluded. In addition to the study on algorithmic characteristics, we provide a unified comparison of the most current exploration methods for DRL on three representative benchmarks. From the empirical results, we shed some light on the specialties and limitations of different categories of exploration methods. Finally, we summarize several open questions. We discuss in depth what is expected to be solved and provide a few potential solutions that are worthwhile being further studied in the future.

At present, exploration in environments with large state-action space, long-horizon, and complex semantics is still very challenging with current advances. Moreover, exploration in deep MARL is much less studied than that in the single-agent setting. Towards a broader perspective to address exploration problems in DRL and deep MARL, we highlight a few suggestions and insights below: 1) First, current exploration methods are evaluated mainly in terms of cumulative rewards and sample efficiency in only several well-known hard-exploration environments or manually designed environments. The lack of multidimensional criteria and standard experimental benchmarks inhibit the calibration and evaluation of different methods. In this survey, we identify several challenges to exploration and use them for qualitative analysis of different methods from the algorithmic aspect. New environments specific to these challenges and corresponding quantitative evaluation criteria are of necessity. 2) Second, the essential connections between different exploration methods are to be further revealed. For example, although intrinsic motivation-oriented exploration

methods usually come from strong heuristics, some of them are closely related to uncertainty-oriented exploration methods, as the underlying connection between RND [51] and Random Prior Fitting [193]. The study on such essential connections can promote the unification and classification of both the theory and methodology of RL exploration. 3) Third, exploration among large action space, exploration in long-horizon environments and convergence analysis are relatively lack of study. The progress in solving these problems is significant to both the practical application of RL algorithms and RL theoretical study. 4) Lastly, multi-agent exploration can be even more challenging due to complex multi-agent interactions. At present, many multi-agent exploration methods are proposed through integrating the exploration methods in single-agent settings and MARL algorithms. Coordinated exploration with decentralized execution and exploration under non-stationarity may be the key problems to address.

## APPENDIX A
### MODEL-BASED EXPLORATION

In the main text, we mainly conduct a survey on exploration methods in a model-free domain. In this section, we briefly introduce exploration methods for model-based RL. Model-based RL uses an environment model to generate simulated experience or planning. In exploration, RL algorithms can directly use this environment model for uncertainty estimation or intrinsic motivation.

For uncertainty-based exploration, model-based RL is based on optimism for planning and exploration [194]. For example, Model-assisted RL [195] uses an ensemble model to measure the dynamics uncertainty, and makes use of artificial data only in cases of high uncertainty, which encourages the agent to learn difficult parts of the environment. Planning to explore [196] also uses ensemble dynamics for uncertainty estimation, and seeks out future uncertainty by integrating uncertainty to Dreamer-based planning [197]. Noise-Augmented RL [198] uses statistical bootstrap to generalize the optimistic posterior sampling [199] to DRL. Hallucinated UCRL [200] measures the epistemic uncertainty and reduces optimistic exploration to exploitation by enlarging the control space.

For intrinsic motivation-based exploration, model-based RL usually uses the information gain of the dynamics to incentivize exploration. Ready Policy One [201] uses Gaussian distribution to build the environment model and measures the information gain through the Gaussian process. Then they optimize the policies for both reward and model uncertainty reduction, which encourages the agent to explore high-uncertainty areas. Dynamics Bottleneck [175] uses variational methods to learn the latent space of dynamics and uses the information gain measured by the difference between the prior and posterior distributions of the latent variables.

## APPENDIX B
### HINDSIGHT-BASED EXPLORATION

In this section, we discuss a specific exploration problem in goal-conditional RL settings. For example, the agent needs to reach the specific goal in mazes, but only receives a reward if it successfully reaches the desired goal. It is almost impossible to reach the goal by chance, even in the simplest environment.

Hindsight experience replay (HER) [202] shows promising results in solving such problems by using hindsight goals to make the agent learn from failures. Specifically, the HER agent samples the already-achieved goals from a failed experience as the hindsight goals. The hindsight goals are used to substitute the original goals and recompute the reward functions. Because the hindsight goals lie close to the sampled transitions, the agent frequently receives the hindsight rewards and accelerates learning. On the basis of HER, hindsight policy gradient [203] extends HER to on-policy RL by using importance sampling. RIG [204] uses HER to handle image-based observation by learning a latent space from a variational autoencoder. Dynamic HER [205] solves dynamic goals in real robotics tasks by assembling successful experiences from two failures. Competitive experience replay [206] introduces self-play [207] between two agents and generates a curriculum for exploration. Entropy-regularized HER [208] develops a maximum entropy framework to optimize the prioritized multi-goal RL. Curriculum-guided HER [209] selects the replay experiences according to the proximity to true goals and the curiosity of exploration. Hindsight goals generation [210] creates valuable goals that are easy to achieve and valuable in the long term. Directed Exploration [211] uses the prediction error of dynamics to choose high uncertain goals and learns a goal-conditional policy to reach them based on HER.

## APPENDIX C
### EXPLORATION VIA SKILL DISCOVERY

Recent research converts the exploration problem to an unsupervised skill-discovery problem. That is, through learning skills in reward-free exploration, the agent collects information about the environment and learns skills that are useful for downstream tasks. We discuss this kind of learning method in the following.

The unsupervised skill-discovery methods aim to learn skills by exploring the environment without any extrinsic rewards. The basic principle of skill discovery is empowerment [142], which describes what the agent can be done while learning how to do it. The skill is a latent variable $z \sim p(z)$ sampled from a skill space, which can be continuous or discrete. The skill discovery methods seeks to find a skill-conditional policy $\pi(a|s, z)$ to maximize the mutual information between $S$ and $Z$. Formally, the mutual information of $S$ and $Z$ can be solved in reverse and forward forms as follows,

$$I(S; Z) = H(Z) - H(Z|S). \quad \text{(Reverse)} \quad (11)$$

$$I(S; Z) = H(S) - H(S|Z). \quad \text{(Forward)} \quad (12)$$

Most existing skill discovery methods follow the reverse form defined in Eq. (11), where $I(S; Z) = \mathbb{E}_{s,z \sim p(s,z)}[\log p(z|s)] - \mathbb{E}_{z \sim p(z)}[\log p(z)]$. A variational lower bound $I(S; Z) \geq \mathbb{E}_{s,z \sim p(s,z)}[\log q_\phi(z|s)] - \mathbb{E}_{z \sim p(z)}[\log p(z)]$ is derived by using $q_\phi(z|s)$ to approximate the posterior $p(z|s)$ of skill, where $q_\phi(z|s)$ is trained by maximizing the likelihood of $(s, z)$ collected by the agent. The mutual information $I(S; Z)$ is maximized by using the variational lower bound as the reward

function when training the policy. Following this principle, variational intrinsic control [138] considers $s$ as the last state in a trajectory and the skill is sampled from a prior distribution. DIAYN [139] trains the policy to maximize $I(S; Z)$ while minimizing $I(A; Z|S)$ additionally, which pushes the skills away from each other to learn distinguishable skills. VALOR [212] uses a trajectory $\tau$ to calculate the variational distribution $q_\phi(z|\tau)$ conditioned on the trajectory rather than $q_\phi(z|s)$ calculated by the individual states, which encourages the agent to learn dynamical modes. Skew-Fit [213] uses states sampled from the replay buffer as goals $G$ and the agent learns to maximize $I(S; G)$, where $G$ serves as the skill $Z$ in Eq. (11).

There are several recent works that use the forward form of $I(S, Z)$ defined in Eq. (12) to maximize the mutual information, where $I(S; Z) = \mathbb{E}_{s,z \sim p(s,z)}[\log p(s|z)] - \mathbb{E}_{s \sim p(s)}[\log p(s)]$. Similar to the reverse form, a variational distribution $q_\phi(s|z)$ that fits the $(s, z)$ tuple collected by the agent is used to form a variational lower bound as $I(S, Z) \geq \mathbb{E}_{s,z \sim p(s,z)}[\log q_\phi(s|z)] - \mathbb{E}_{s \sim p(s)}[\log p(s)]$. DADS [140] maximizes the mutual information between skill and the next state $I(S', Z|S)$ that conditions on the current state. Then the variational distribution $q_\phi(s|z)$ becomes $q_\phi(s'|s, z)$, which represents a skill-level dynamics model and enables the use of model-based planning algorithms for downstream tasks. EDL [141] finds that existing skill discovery methods tend to visit known states rather than discovering new ones when maximizing the mutual information, which leads to a poor coverage of the state space. Based on the analyses, EDL uses State Marginal Matching (SMM) [111] to perform maximum entropy exploration in the state space before skill discovery, which results in a fixed $p(s)$ that is agnostic to the skill. EDL succeeds at discovering state-covering skills in environments where previous methods failed.

## APPENDIX D
## SELF-IMITATION LEARNING

In addition to the various methods we have discussed, Self-Imitation Learning (SIL) is another branch of methods that can facilitate RL in the face of learning and exploration difficulties. The main idea of SIL is to imitate the agent's own experiences collected during the historical learning process. A SIL agent tries to make the best use of the superior (e.g., successful) experiences encountered by occasion.

Oh et al. [168] firstly propose the concept of SIL. They propose additional off-policy AC losses to ascend the policy towards non-negative advantages and make the value function approximate the corresponding returns. They demonstrate the effectiveness of SIL when combined with A2C and PPO in 49 Atari games and the Key-Door-Treasure domain. Furthermore, they shed light on the theoretical connection between SIL and lower-bound Soft Q-Learning [214]. From a different angle, Gangwani et al. [215] derive a SIL algorithm by casting the divergence minimization problem of the state-action visitation distributions of the policy to optimize and (historical) good policies, to a policy gradient objective. In addition, Gangwani et al. [215] utilize Stein Variation Policy Gradient (SVPG) [216]

for the purpose of optimizing an ensemble of diverse SI policies. Later, SIL is developed from transition level to trajectory level with the help of trajectory [217] or goal conditioned policies [218]. Diverse Trajectory-conditioned SIL (DTSIL) [217] is proposed to maintain a buffer of good and diverse trajectories based on their similarity which are taken as input of a trajectory-conditioned policy during SIL. Episodic SIL (ESIL) [218] is proposed to also perform trajectory-level SIL with a trajectory-based episodic memory; to train a goal-conditioned policy, HER [202] is utilized to alter the original (failed) goals among which good trajectories to self-imitate are then determined and filtered.

Recently, a few works improve original SIL algorithm [168] in different ways. Tang [219] proposes a family of SIL algorithms through generalizing the original return-based lower-bound Q-learning [168] to $n$-step TD lower bound, balancing the trade-off between fixed point bias and contraction rate. Chen et al. [220] propose SIL with Constant Reward (SILCR) to get rid of the need of immediate environment rewards (but make use of episodic reward instead). Beyond policy-based SIL, Ferret et al. propose Self-Imitation Advantage Learning (SAIL) [221] which extends SIL to off-policy value-based RL algorithms through modifying Bellman optimality operator that connects to Advantage Learning [99]. Ning et al. [222] propose Co-Imitation Learning (CoIL) which learns two different agents via letting each of them alternately explore the environment and selectively imitate the heterogeneous good experiences from each other.

Overall, SIL is different from the two main categories of methods considered in our main text, i.e., intrinsic motivation and uncertainty. From some angle, SIL methods can be viewed as passive exploration methods which reinforce the exploration the informative reward signals in the experience buffer.

## REFERENCES

[1] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, pp. 484–489, 2016.

[2] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton *et al.*, "Mastering the game of go without human knowledge," *nature*, vol. 550, no. 7676, pp. 354–359, 2017.

[3] M. Hessel, J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. G. Azar, and D. Silver, "Rainbow: Combining improvements in deep reinforcement learning," in *AAAI*, 2018.

[4] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev *et al.*, "Grandmaster level in starcraft ii using multi-agent reinforcement learning," *Nature*, vol. 575, no. 7782, pp. 350–354, 2019.

[5] L. Han, J. Xiong, P. Sun, X. Sun, M. Fang, Q. Guo, Q. Chen, T. Shi, H. Yu, and Z. Zhang, "Tstarbot-x: An open-sourced and comprehensive study for efficient league training in starcraft ii full game," *arXiv preprint arXiv:2011.13729*, 2020.

[6] M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray *et al.*, "Learning dexterous in-hand manipulation," *International Journal of Robotics Research*, vol. 39, no. 1, pp. 3–20, 2020.

[7] A. P. Badia, B. Piot, S. Kapturowski, P. Sprechmann, A. Vitvitskyi, D. Guo, and C. Blundell, "Agent57: Outperforming the atari human benchmark," *arXiv preprint arXiv:2003.13350*, 2020.

[8] H.-n. Wang, N. Liu, Y.-y. Zhang, D.-w. Feng, F. Huang, D.-s. Li, and Y.-m. Zhang, "Deep reinforcement learning: a survey," *FRONT INFORM TECH EL*, pp. 1–19, 2020.

[9] Y. Li, "Deep reinforcement learning: An overview," *arXiv preprint arXiv:1701.07274*, 2017.

[10] S. S. Mousavi, M. Schukat, and E. Howley, "Deep reinforcement learning: an overview," in *SAI Intelligent Systems Conference*, 2016.

[11] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep reinforcement learning: A brief survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.

[12] L. Busoniu, R. Babuska, and B. De Schutter, "A comprehensive survey of multiagent reinforcement learning," *TSMC*, vol. 38, no. 2, pp. 156–172, 2008.

[13] P. Hernandez-Leal, B. Kartal, and M. E. Taylor, "A survey and critique of multiagent deep reinforcement learning," *JAAMAS*, vol. 33, no. 6, pp. 750–797, 2019.

[14] K. Zhang, Z. Yang, and T. Başar, "Multi-agent reinforcement learning: A selective overview of theories and algorithms," *arXiv preprint arXiv:1911.10635*, 2019.

[15] A. Aubret, L. Matignon, and S. Hassas, "A survey on intrinsic motivation in reinforcement learning," *arXiv preprint arXiv:1908.06976*, 2019.

[16] T. Lattimore and C. Szepesvári, *Bandit Algorithms*. Cambridge University Press, 2020.

[17] R. M. Ryan and E. L. Deci, "Intrinsic and extrinsic motivations: Classic definitions and new directions," *Contemporary educational psychology*, vol. 25, no. 1, pp. 54–67, 2000.

[18] A. G. Barto, "Intrinsic motivation and reinforcement learning," in *Intrinsically motivated learning in natural and artificial systems*, 2013.

[19] A. Gopnik, A. N. Meltzoff, and P. K. Kuhl, *The scientist in the crib: Minds, brains, and how children learn.* William Morrow & Co, 1999.

[20] R. S. Sutton and A. G. Barto, *Reinforcement learning - an introduction*, ser. Adaptive computation and machine learning, 1998.

[21] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.

[22] H. van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *AAAI*, 2016.

[23] O. Anschel, N. Baram, and N. Shimkin, "Averaged-dqn: Variance reduction and stabilization for deep reinforcement learning," in *ICML*, 2017.

[24] M. G. Bellemare, W. Dabney, and R. Munos, "A distributional perspective on reinforcement learning," in *ICML*, 2017.

[25] W. Dabney, M. Rowland, M. G. Bellemare, and R. Munos, "Distributional reinforcement learning with quantile regression," in *AAAI*, 2018.

[26] Z. Wang, T. Schaul, M. Hessel, H. van Hasselt, M. Lanctot, and N. de Freitas, "Dueling network architectures for deep reinforcement learning," in *ICML*, 2016.

[27] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," in *ICLR*, 2016.

[28] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine Learning*, vol. 8, pp. 229–256, 1992.

[29] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

[30] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. A. Riedmiller, "Deterministic policy gradient algorithms," in *ICML*, 2014.

[31] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," in *ICLR*, 2016.

[32] J. N. Foerster, N. Nardelli, G. Farquhar, T. Afouras, P. H. S. Torr, P. Kohli, and S. Whiteson, "Stabilising experience replay for deep multi-agent reinforcement learning," in *ICML*, 2017.

[33] S. Omidshafiei, J. Pazis, C. Amato, J. P. How, and J. Vian, "Deep decentralized multi-task multi-agent reinforcement learning under partial observability," in *ICML*, 2017.

[34] S. Sukhbaatar, A. Szlam, and R. Fergus, "Learning multiagent communication with backpropagation," in *NeurIPS*, 2016.

[35] J. N. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson, "Learning to communicate with deep multi-agent reinforcement learning," in *NeurIPS*, 2016.

[36] J. N. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," in *AAAI*, 2018.

[37] T. Rashid, M. Samvelyan, C. S. de Witt, G. Farquhar, J. N. Foerster, and S. Whiteson, "QMIX: monotonic value function factorisation for deep multi-agent reinforcement learning," in *ICML*, 2018.

[38] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *NeurIPS*, 2017.

[39] T. L. Lai and H. Robbins, "Asymptotically efficient adaptive allocation rules," *Advances in applied mathematics*, vol. 6, no. 1, pp. 4–22, 1985.

[40] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *ICML*, 2016.

[41] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *ICML*, 2018.

[42] P. I. Frazier, "A tutorial on bayesian optimization," *arXiv preprint arXiv:1807.02811*, 2018.

[43] N. Srinivas, A. Krause, S. M. Kakade, and M. W. Seeger, "Gaussian process optimization in the bandit setting: No regret and experimental design," in *ICML*, 2010.

[44] W. R. Thompson, "On the likelihood that one unknown probability exceeds another in view of the evidence of two samples," *Biometrika*, vol. 25, no. 3/4, pp. 285–294, 1933.

[45] P. Auer, "Using confidence bounds for exploitation-exploration trade-offs," *JMLR*, vol. 3, pp. 397–422, 2002.

[46] M. J. Hausknecht and P. Stone, "Deep reinforcement learning in parameterized action space," in *ICLR*, 2016.

[47] Y. Chandak, G. Theocharous, J. Kostas, S. M. Jordan, and P. S. Thomas, "Learning action representations for reinforcement learning," in *ICML*, 2019.

[48] B. Li, H. Tang, Y. Zheng, J. Hao, P. Li, Z. Wang, Z. Meng, and L. Wang, "Hyar: Addressing discrete-continuous action reinforcement learning via hybrid action representation," *arXiv preprint arXiv:2109.05490*, 2021.

[49] M. J. A. Strens, "A bayesian framework for reinforcement learning," in *ICML*, 2000.

[50] I. Osband, J. Aslanides, and A. Cassirer, "Randomized prior functions for deep reinforcement learning," in *NeurIPS*, 2018.

[51] Y. Burda, H. Edwards, A. J. Storkey, and O. Klimov, "Exploration by random network distillation," in *ICLR*, 2019.

[52] J. Choi, Y. Guo, M. Moczulski, J. Oh, N. Wu, M. Norouzi, and H. Lee, "Contingency-aware exploration in reinforcement learning," in *ICLR*, 2019.

[53] A. P. Badia, P. Sprechmann, A. Vitvitskyi, D. Guo, B. Piot, S. Kapturowski, O. Tieleman, M. Arjovsky, A. Pritzel, A. Bolt, and C. Blundell, "Never give up: Learning directed exploration strategies," *arXiv preprint arXiv:2002.06038*, 2020.

[54] C. Tessler, S. Givony, T. Zahavy, D. J. Mankowitz, and S. Mannor, "A deep hierarchical approach to lifelong learning in minecraft," in *AAAI*, 2017.

[55] M. Kempka, M. Wydmuch, G. Runc, J. Toczek, and W. Jaskowski, "Vizdoom: A doom-based AI research platform for visual reinforcement learning," in *CIG*, 2016.

[56] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, "Curiosity-driven exploration by self-supervised prediction," in *ICML*, 2017.

[57] Y. Burda, H. Edwards, D. Pathak, A. J. Storkey, T. Darrell, and A. A. Efros, "Large-scale study of curiosity-driven learning," in *ICLR*, 2019.

[58] Y. Kim, W. Nam, H. Kim, J. Kim, and G. Kim, "Curiosity-bottleneck: Exploration by distilling task-specific novelty," in *ICML*, 2019.

[59] I. Osband, B. Van Roy, D. J. Russo, and Z. Wen, "Deep exploration via randomized value functions." *JMLR*, vol. 20, no. 124, pp. 1–62, 2019.

[60] I. Osband, B. V. Roy, and Z. Wen, "Generalization and exploration via randomized value functions," in *ICML*, 2016.

[61] K. Azizzadenesheli, E. Brunskill, and A. Anandkumar, "Efficient exploration through bayesian deep q-networks," in *Information Theory and Applications Workshop*, 2018.

[62] D. Janz, J. Hron, P. Mazur, K. Hofmann, J. M. Hernández-Lobato, and S. Tschiatschek, "Successor uncertainties: Exploration and uncertainty in temporal difference learning," in *NeurIPS*, 2019.

[63] A. M. Metelli, A. Likmeta, and M. Restelli, "Propagating uncertainty in reinforcement learning via wasserstein barycenters," in *NeurIPS*, 2019.

[64] B. O'Donoghue, I. Osband, R. Munos, and V. Mnih, "The uncertainty bellman equation and exploration," in *ICML*, 2018.

[65] I. Osband, C. Blundell, A. Pritzel, and B. V. Roy, "Deep exploration via bootstrapped DQN," in *NeurIPS*, 2016.

[66] K. Ciosek, Q. Vuong, R. Loftin, and K. Hofmann, "Better exploration with optimistic actor critic," in *NeurIPS*, 2019.

[67] K. Lee, M. Laskin, A. Srinivas, and P. Abbeel, "SUNRISE: A simple unified framework for ensemble learning in deep reinforcement learning," in *ICML*, 2021.

[68] C. Bai, L. Wang, L. Han, J. Hao, A. Garg, P. Liu, and Z. Wang, "Principled exploration via optimistic bootstrapping and backward induction," in *ICML*, 2021.

[69] T. M. Moerland, J. Broekens, and C. M. Jonker, "Efficient exploration with double uncertain value networks," *arXiv preprint arXiv:1711.10789*, 2017.

[70] J. Kirschner and A. Krause, "Information directed sampling and bandits with heteroscedastic noise," in *CoLT*, 2018.

[71] B. Mavrin, H. Yao, L. Kong, K. Wu, and Y. Yu, "Distributional reinforcement learning for efficient exploration," in *ICML*, 2019.

[72] F. Zhou, J. Wang, and X. Feng, "Non-crossing quantile regression for distributional reinforcement learning," in *NeurIPS*, 2020.

[73] R. Dearden, N. Friedman, and S. J. Russell, "Bayesian q-learning," in *AAAI*, 1998.

[74] P. Auer, N. Cesa-Bianchi, and P. Fischer, "Finite-time analysis of the multiarmed bandit problem," *Machine learning*, vol. 47, no. 2-3, pp. 235–256, 2002.

[75] K. Azizzadenesheli, A. Lazaric, and A. Anandkumar, "Reinforcement learning of pomdps using spectral methods," in *CoLT*, 2016.

[76] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *ICML*, 2016.

[77] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," in *NeurIPS*, 2017.

[78] M. Welling and Y. W. Teh, "Bayesian learning via stochastic gradient langevin dynamics," in *ICML*, 2011.

[79] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight uncertainty in neural networks," *arXiv preprint arXiv:1505.05424*, 2015.

[80] I. Osband, D. Russo, and B. V. Roy, "(more) efficient reinforcement learning via posterior sampling," in *NeurIPS*, 2013.

[81] A. Zanette, D. Brandfonbrener, E. Brunskill, M. Pirotta, and A. Lazaric, "Frequentist regret bounds for randomized least-squares value iteration," in *AISTATS*, 2020.

[82] P. Dayan, "Improving generalization for temporal difference learning: The successor representation," *Neural Computation*, vol. 5, no. 4, pp. 613–624, 1993.

[83] A. Barreto, W. Dabney, R. Munos, J. J. Hunt, T. Schaul, D. Silver, and H. van Hasselt, "Successor features for transfer in reinforcement learning," in *NeurIPS*, 2017.

[84] M. Agueh and G. Carlier, "Barycenters in the wasserstein space," *SIAM J. Math. Anal.*, vol. 43, no. 2, pp. 904–924, 2011.

[85] I. Osband and B. V. Roy, "Bootstrapped thompson sampling and deep exploration," *arXiv preprint arXiv:1507.00300*, 2015.

[86] Y. Zhang and W. Goh, "Bootstrapped policy gradient for difficulty adaptation in intelligent tutoring systems," in *AAMAS*, 2019.

[87] G. Kalweit and J. Boedecker, "Uncertainty-driven imagination for continuous deep reinforcement learning," in *CoRL*, 2017.

[88] Z. Yang, K. E. Merrick, H. A. Abbass, and L. Jin, "Multi-task deep reinforcement learning for continuous action control," in *IJCAI*, 2017.

[89] Z. Zheng, C. Yuan, Z. Lin, Y. Cheng, and H. Wu, "Self-adaptive double bootstrapped DDPG," in *IJCAI*, 2018.

[90] N. Nikolov, J. Kirschner, F. Berkenkamp, and A. Krause, "Information-directed exploration for deep reinforcement learning," in *ICLR*, 2019.

[91] W. R. Clements, B. Robaglia, B. V. Delft, R. B. Slaoui, and S. Toth, "Estimating risk and uncertainty in deep reinforcement learning," *arXiv preprint arXiv:1905.09638*, 2019.

[92] E. L. Deci, "Intrinsic motivation and self-determination," in *Encyclopedia of Applied Psychology*, 2004.

[93] P.-Y. Oudeyer, F. Kaplan, and V. V. Hafner, "Intrinsic motivation systems for autonomous mental development," *IEEE transactions on evolutionary computation*, vol. 11, no. 2, pp. 265–286, 2007.

[94] B. C. Stadie, S. Levine, and P. Abbeel, "Incentivizing exploration in reinforcement learning with deep predictive models," *arXiv preprint arXiv:1507.00814*, 2015.

[95] C. Oh and A. Cavallaro, "Learning action representations for self-supervised visual exploration," in *ICRA*, 2019.

[96] C. Bai, P. Liu, Z. Wang, K. Liu, L. Wang, and Y. Zhao, "Variational dynamic for self-supervised exploration in deep reinforcement learning," *arXiv preprint arXiv:2010.08755*, 2020.

[97] H. Kim, J. Kim, Y. Jeong, S. Levine, and H. O. Song, "EMI: exploration with mutual information," in *ICML*, 2019.

[98] H. Tang, R. Houthooft, D. Foote, A. Stooke, X. Chen, Y. Duan, J. Schulman, F. D. Turck, and P. Abbeel, "#exploration: A study of count-based exploration for deep reinforcement learning," in *NeurIPS*, 2017.

[99] M. G. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos, "Unifying count-based exploration and intrinsic motivation," in *NeurIPS*, 2016.

[100] G. Ostrovski, M. G. Bellemare, A. van den Oord, and R. Munos, "Count-based exploration with neural density models," in *ICML*, 2017.

[101] J. Martin, S. N. Sasikumar, T. Everitt, and M. Hutter, "Count-based exploration in feature space for reinforcement learning," in *IJCAI*, 2017.

[102] G. Vezzani, A. Gupta, L. Natale, and P. Abbeel, "Learning latent state representation for speeding up exploration," *arXiv preprint arXiv:1905.12621*, 2019.

[103] M. C. Machado, M. G. Bellemare, and M. Bowling, "Count-based exploration with the successor representation," in *AAAI*, 2020.

[104] L. Fox, L. Choshen, and Y. Loewenstein, "DORA the explorer: Directed outreaching reinforcement action-selection," in *ICLR*, 2018.

[105] Y. Song, Y. Chen, Y. Hu, and C. Fan, "Exploring unknown states with action balance," in *CIG*, 2020.

[106] T. Zhang, P. Rashidinejad, J. Jiao, Y. Tian, J. Gonzalez, and S. Russell, "Made: Exploration via maximizing deviation from explored regions," *arXiv preprint arXiv:2106.10268*, 2021.

[107] J. Oh, X. Guo, H. Lee, R. L. Lewis, and S. P. Singh, "Action-conditional video prediction using deep networks in atari games," in *NeurIPS*, 2015.

[108] J. Fu, J. D. Co-Reyes, and S. Levine, "EX2: exploration with exemplar models for deep reinforcement learning," in *NeurIPS*, 2017.

[109] J. Zhang, N. Wetzel, N. Dorka, J. Boedecker, and W. Burgard, "Scheduled intrinsic drive: A hierarchical take on intrinsically motivated exploration," *arXiv preprint arXiv:1903.07400*, 2019.

[110] M. Klissarov, R. Islam, K. Khetarpal, and D. Precup, "Variational state encoding as intrinsic motivation in reinforcement learning," in *TARL Workshop on ICLR*, 2019.

[111] L. Lee, B. Eysenbach, E. Parisotto, E. P. Xing, S. Levine, and R. Salakhutdinov, "Efficient exploration via state marginal matching," *arXiv preprint arXiv:1906.05274*, 2019.

[112] C. Stanton and J. Clune, "Deep curiosity search: Intra-life exploration improves performance on challenging deep reinforcement learning problems," *arXiv preprint arXiv:1806.00553*, 2018.

[113] R. Y. Tao, V. François-Lavet, and J. Pineau, "Novelty search in representational space for sample efficient exploration," in *NeurIPS*, 2020.

[114] N. Savinov, A. Raichuk, D. Vincent, R. Marinier, M. Pollefeys, T. P. Lillicrap, and S. Gelly, "Episodic curiosity through reachability," in *ICLR*, 2019.

[115] T. Zhang, H. Xu, X. Wang, Y. Wu, K. Keutzer, J. E. Gonzalez, and Y. Tian, "Noveld: A simple yet effective exploration criterion," in *Advances in Neural Information Processing Systems*, 2021.

[116] R. Houthooft, X. Chen, Y. Duan, J. Schulman, F. D. Turck, and P. Abbeel, "VIME: variational information maximizing exploration," in *NeurIPS*, 2016.

[117] J. Achiam and S. Sastry, "Surprise-based intrinsic motivation for deep reinforcement learning," *arXiv preprint arXiv:1703.01732*, 2017.

[118] D. Pathak, D. Gandhi, and A. Gupta, "Self-supervised exploration via disagreement," in *ICML*, 2019.

[119] P. Shyam, W. Jaskowski, and F. Gomez, "Model-based active exploration," in *ICML*, 2019.

[120] R. I. Brafman and M. Tennenholtz, "R-MAX - A general polynomial time algorithm for near-optimal reinforcement learning," *JMLR*, vol. 3, pp. 213–231, 2002.

[121] M. J. Kearns and S. P. Singh, "Near-optimal reinforcement learning in polynomial time," *Machine Learning*, vol. 49, no. 2-3, pp. 209–232, 2002.

[122] M. Charikar, "Similarity estimation techniques from rounding algorithms," in *STOC*, 2002.

[123] A. van den Oord, N. Kalchbrenner, L. Espeholt, K. Kavukcuoglu, O. Vinyals, and A. Graves, "Conditional image generation with pixelcnn decoders," in *NeurIPS*, 2016.

[124] M. G. Bellemare, J. Veness, and E. Talvitie, "Skip context tree switching," in *ICML*, 2014.

[125] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, "Deep variational information bottleneck," in *ICLR*, 2017.

[126] P. Oudeyer and F. Kaplan, "What is intrinsic motivation? A typology of computational approaches," *Frontiers Neurorobotics*, vol. 1, p. 6, 2007.

[127] L. Itti and P. Baldi, "Bayesian surprise attracts human attention," in *NeurIPS*, 2005.

[128] M. Frank, J. Leitner, M. F. Stollenga, A. Förster, and J. Schmidhuber, "Curiosity driven reinforcement learning for motion planning on humanoids," *Frontiers Neurorobotics*, vol. 7, p. 25, 2013.

[129] A. Graves, "Practical variational inference for neural networks," in *NeurIPS*, 2011.

[130] D. Horgan, J. Quan, D. Budden, G. Barth-Maron, M. Hessel, H. van Hasselt, and D. Silver, "Distributed prioritized experience replay," in *ICLR*, 2018.

[131] S. Kapturowski, G. Ostrovski, J. Quan, R. Munos, and W. Dabney, "Recurrent experience replay in distributed reinforcement learning," in *ICLR*, 2019.

[132] M. Plappert, R. Houthooft, P. Dhariwal, S. Sidor, R. Y. Chen, X. Chen, T. Asfour, P. Abbeel, and M. Andrychowicz, "Parameter space noise for exploration," in *ICLR*, 2018.

[133] M. Fortunato, M. G. Azar, B. Piot, J. Menick, M. Hessel, I. Osband, A. Graves, V. Mnih, R. Munos, D. Hassabis, O. Pietquin, C. Blundell, and S. Legg, "Noisy networks for exploration," in *ICLR*, 2018.

[134] A. Ecoffet, J. Huizinga, J. Lehman, K. O. Stanley, and J. Clune, "Go-explore: a new approach for hard-exploration problems," *arXiv preprint arXiv:1901.10995*, 2019.

[135] E. Adrien, H. Joost, L. Joel, S. K. O, and C. Jeff, "First return, then explore," *Nature*, vol. 590, no. 7847, pp. 580–586, 2021.

[136] T. Hester, M. Vecerík, O. Pietquin, M. Lanctot, T. Schaul, B. Piot, D. Horgan, J. Quan, A. Sendonaris, I. Osband, G. Dulac-Arnold, J. P. Agapiou, J. Z. Leibo, and A. Gruslys, "Deep q-learning from demonstrations," in *AAAI*, 2018.

[137] E. Zhao, S. Deng, Y. Zang, Y. Kang, K. Li, and J. Xing, "Potential driven reinforcement learning for hard exploration tasks," in *IJCAI*, 2020.

[138] K. Gregor, D. J. Rezende, and D. Wierstra, "Variational intrinsic control," *arXiv preprint arXiv:1611.07507*, 2016.

[139] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine, "Diversity is all you need: Learning skills without a reward function," in *ICLR*, 2019.

[140] A. Sharma, S. Gu, S. Levine, V. Kumar, and K. Hausman, "Dynamics-aware unsupervised discovery of skills," in *ICLR*, 2020.

[141] V. Campos, A. Trott, C. Xiong, R. Socher, X. Giro-i Nieto, and J. Torres, "Explore, discover and learn: Unsupervised discovery of state-covering skills," in *ICML*, 2020.

[142] C. Salge, C. Glackin, and D. Polani, "Empowerment–an introduction," in *Guided Self-Organization: Inception*, 2014.

[143] T. Bolander and M. B. Andersen, "Epistemic planning for single and multi-agent systems," *JANCL*, vol. 21, no. 1, pp. 9–34, 2011.

[144] Z. Zhu, E. Biyik, and D. Sadigh, "Multi-agent safe planning with gaussian processes," *arXiv preprint arXiv:2008.04452*, 2020.

[145] C. Martin and T. Sandholm, "Efficient exploration of zero-sum stochastic games," *arXiv preprint arXiv:2002.10524*, 2020.

[146] D. J. Russo, B. Van Roy, A. Kazerouni, I. Osband, and Z. Wen, "A tutorial on thompson sampling," *Foundations and Trends® in Machine Learning*, vol. 11, no. 1, pp. 1–96, 2018.

[147] E. Kaufmann, O. Cappé, and A. Garivier, "On bayesian upper confidence bounds for bandit problems," in *AISTATS*, 2012.

[148] J. Hu, S. A. Harding, H. Wu, and S. Liao, "QR-MIX: distributional value function factorisation for cooperative multi-agent reinforcement learning," *arXiv preprint arXiv:2009.04197*, 2020.

[149] J. Xiong, Q. Wang, Z. Yang, P. Sun, L. Han, Y. Zheng, H. Fu, T. Zhang, J. Liu, and H. Liu, "Parametrized deep q-networks learning: Reinforcement learning with discrete-continuous hybrid action space," *arXiv preprint arXiv:1810.06394*, 2018.

[150] X. Lyu and C. Amato, "Likelihood quantile networks for coordinating multi-agent reinforcement learning," in *AAMAS*, 2020.

[151] W. Böhmer, T. Rashid, and S. Whiteson, "Exploration with unreliable intrinsic reward in multi-agent reinforcement learning," *arXiv preprint arXiv:1906.02138*, 2019.

[152] S. Iqbal and F. Sha, "Coordinated exploration via intrinsic rewards for multi-agent reinforcement learning," *CoRR*, vol. abs/1905.12127, 2019.

[153] Y. Du, L. Han, M. Fang, J. Liu, T. Dai, and D. Tao, "LIIR: learning individual intrinsic reward in multi-agent reinforcement learning," in *NeurIPS*, 2019.

[154] T. Wang, J. Wang, Y. Wu, and C. Zhang, "Influence-based multi-agent exploration," in *ICLR*, 2020.

[155] Z. Zheng, J. Oh, and S. Singh, "On learning intrinsic rewards for policy gradient methods," in *NeurIPS*, 2018.

[156] Z. Zheng, J. Oh, M. Hessel, Z. Xu, M. Kroiss, H. van Hasselt, D. Silver, and S. Singh, "What can learned intrinsic rewards capture?" in *ICML*, 2020.

[157] N. Jaques, A. Lazaridou, E. Hughes, Ç. Gülçehre, P. A. Ortega, D. Strouse, J. Z. Leibo, and N. de Freitas, "Social influence as intrinsic motivation for multi-agent deep reinforcement learning," in *ICML*, 2019.

[158] R. Chitnis, S. Tulsiani, S. Gupta, and A. Gupta, "Intrinsic motivation for encouraging synergistic behavior," in *ICLR*, 2020.

[159] D. Carmel and S. Markovitch, "Exploration strategies for model-based learning in multi-agent systems: Exploration strategies," *JAAMAS*, vol. 2, no. 2, pp. 141–172, 1999.

[160] G. Chalkiadakis and C. Boutilier, "Coordination in multiagent reinforcement learning: a bayesian approach," in *AAMAS*, 2003.

[161] K. Verbeeck, A. Nowé, M. Peeters, and K. Tuyls, "Multi-agent reinforcement learning in stochastic single and multi-stage games," in *ALAMAS*, 2005.

[162] M. Chakraborty, K. Y. P. Chua, S. Das, and B. Juba, "Coordinated versus decentralized exploration in multi-agent multi-armed bandits," in *IJCAI*, 2017.

[163] M. Dimakopoulou and B. V. Roy, "Coordinated exploration in concurrent reinforcement learning," in *ICML*, 2018.

[164] M. Dimakopoulou, I. Osband, and B. V. Roy, "Scalable coordinated exploration in concurrent reinforcement learning," in *NeurIPS*, 2018.

[165] G. Chen, "A new framework for multi-agent reinforcement learning - centralized training and exploration with decentralized execution via policy distillation," in *AAMAS*, 2020.

[166] A. Mahajan, T. Rashid, M. Samvelyan, and S. Whiteson, "MAVEN: multi-agent variational exploration," in *NeurIPS*, 2019.

[167] J. Schulman, S. Levine, P. Abbeel, M. I. Jordan, and P. Moritz, "Trust region policy optimization," in *ICML*, 2015.

[168] J. Oh, Y. Guo, S. Singh, and H. Lee, "Self-imitation learning," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 80, 2018, pp. 3875–3884.

[169] T. Salimans and R. Chen, "Learning montezuma's revenge from a single demonstration," *arXiv preprint arXiv:1812.03381*, 2018.

[170] A. A. Taiga, W. Fedus, M. C. Machado, A. Courville, and M. G. Bellemare, "On bonus based exploration methods in the arcade learning environment," in *ICLR*, 2020.

[171] G. Farquhar, L. Gustafson, Z. Lin, S. Whiteson, N. Usunier, and G. Synnaeve, "Growing action spaces," in *ICML*, 2020.

[172] H. Fu, H. Tang, J. Hao, Z. Lei, Y. Chen, and C. Fan, "Deep multi-agent reinforcement learning with discrete-continuous hybrid action spaces," in *IJCAI*, 2019.

[173] M. Laskin, A. Srinivas, and P. Abbeel, "CURL: contrastive unsupervised representations for reinforcement learning," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 119, 2020, pp. 5639–5650.

[174] D. Yarats, R. Fergus, A. Lazaric, and L. Pinto, "Reinforcement learning with prototypical representations," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 139, 2021, pp. 11 920–11 931.

[175] C. Bai, L. Wang, L. Han, A. Garg, J. Hao, P. Liu, and Z. Wang, "Dynamic bottleneck for robust self-supervised exploration," *Advances in Neural Information Processing Systems*, vol. 34, 2021.

[176] D. Ghosh, A. Gupta, and S. Levine, "Learning actionable representations with goal conditioned policies," in *ICLR*, 2019.

[177] P. Mirowski, M. Grimes, M. Malinowski, K. M. Hermann, K. Anderson, D. Teplyashin, K. Simonyan, A. Zisserman, R. Hadsell *et al.*, "Learning to navigate in cities without a map," in *NeurIPS*, 2018.

[178] N. Savinov, A. Dosovitskiy, and V. Koltun, "Semi-parametric topological memory for navigation," in *ICLR*, 2018.

[179] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.

[180] K. H. Kim, M. Sano, J. De Freitas, N. Haber, and D. Yamins, "Active world model learning in agent-rich environments with progress curiosity," in *ICML*, 2020.

[181] M. C. Machado, M. G. Bellemare, E. Talvitie, J. Veness, M. Hausknecht, and M. Bowling, "Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents," *Journal of Artificial Intelligence Research*, vol. 61, pp. 523–562, 2018.

[182] C. Jin, Z. Allen-Zhu, S. Bubeck, and M. I. Jordan, "Is q-learning provably efficient?" in *NeurIPS*, 2018.

[183] Q. Cai, Z. Yang, C. Jin, and Z. Wang, "Provably efficient exploration in policy optimization," in *ICML*, 2020.

[184] C. Jin, Z. Yang, Z. Wang, and M. I. Jordan, "Provably efficient reinforcement learning with linear function approximation," in *CoLT*, 2020.

[185] P. Auer, T. Jaksch, and R. Ortner, "Near-optimal regret bounds for reinforcement learning," in *NeurIPS*, 2009.

[186] T. Lattimore, "Regret analysis of the anytime optimally confident ucb algorithm," *arXiv preprint arXiv:1603.08661*, 2016.

[187] Z. C. Lipton, X. Li, J. Gao, L. Li, F. Ahmed, and L. Deng, "Bbq-networks: Efficient exploration in deep reinforcement learning for task-oriented dialogue systems," in *AAAI*, 2018.

[188] N. Tagasovska and D. Lopez-Paz, "Single-model uncertainties for deep learning," *arXiv preprint arXiv:1811.00908*, 2018.

[189] S. Thudumu, P. Branch, J. Jin, and J. J. Singh, "A comprehensive survey of anomaly detection techniques for high dimensional big data," *Journal of Big Data*, vol. 7, no. 1, pp. 1–30, 2020.

[190] A. Y. Ng, D. Harada, and S. Russell, "Policy invariance under reward transformations: Theory and application to reward shaping," in *ICML*, 1999.

[191] L. Beyer, D. Vincent, O. Teboul, S. Gelly, M. Geist, and O. Pietquin, "MULEX: disentangling exploitation from exploration in deep RL," *arXiv preprint arXiv:1907.00868*, 2019.

[192] Y. Zhang, W. Yu, and G. Turk, "Learning novel policies for tasks," in *ICML*, 2019.

[193] K. Ciosek, V. Fortuin, R. Tomioka, K. Hofmann, and R. E. Turner, "Conservative uncertainty estimation by fitting prior networks," in *ICLR*, 2020.

[194] D. A. Nix and A. S. Weigend, "Estimating the mean and variance of the target probability distribution," in *ICNN*, vol. 1. IEEE, 1994, pp. 55–60.

[195] G. Kalweit and J. Boedecker, "Uncertainty-driven imagination for continuous deep reinforcement learning," in *CoRL*, 2017, pp. 195–206.

[196] R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak, "Planning to explore via self-supervised world models," in *ICML*, 2020.

[197] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi, "Dream to control: Learning behaviors by latent imagination," in *ICLR*, 2020.

[198] A. Pacchiano, P. Ball, J. Parker-Holder, K. Choromanski, and S. Roberts, "On optimism in model-based reinforcement learning," *arXiv preprint arXiv:2006.11911*, 2020.

[199] S. Agrawal and R. Jia, "Posterior sampling for reinforcement learning: worst-case regret bounds," in *Advances in Neural Information Processing Systems*, 2017, pp. 1184–1194.

[200] S. Curi, F. Berkenkamp, and A. Krause, "Efficient model-based reinforcement learning through optimistic policy search and planning," *Advances in Neural Information Processing Systems*, vol. 33, 2020.

[201] P. Ball, J. Parker-Holder, A. Pacchiano, K. Choromanski, and S. Roberts, "Ready policy one: World building through active learning," *arXiv preprint arXiv:2002.02693*, 2020.

[202] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba, "Hindsight experience replay," in *NeurIPS*, 2017.

[203] P. Rauber, F. Mutz, and J. Schmidhuber, "Hindsight policy gradients," in *ICLR*, 2019.

[204] A. V. Nair, V. Pong, M. Dalal, S. Bahl, S. Lin, and S. Levine, "Visual reinforcement learning with imagined goals," in *NeurIPS*, 2018.

[205] M. Fang, C. Zhou, B. Shi, B. Gong, J. Xu, and T. Zhang, "Dher: Hindsight experience replay for dynamic goals," in *ICLR*, 2019.

[206] H. Liu, A. Trott, R. Socher, and C. Xiong, "Competitive experience replay," in *ICLR*, 2019.

[207] S. Sukhbaatar, I. Kostrikov, A. Szlam, and R. Fergus, "Intrinsic motivation and automatic curricula via asymmetric self-play," in *ICLR*, 2018.

[208] R. Zhao, X. Sun, and V. Tresp, "Maximum entropy-regularized multi-goal reinforcement learning," in *ICML*, 2019.

[209] M. Fang, T. Zhou, Y. Du, L. Han, and Z. Zhang, "Curriculum-guided hindsight experience replay," in *NeurIPS*, 2019.

[210] Z. Ren, K. Dong, Y. Zhou, Q. Liu, and J. Peng, "Exploration via hindsight goal generation," in *NeurIPS*, 2019.

[211] Z. D. Guo and E. Brunskill, "Directed exploration for reinforcement learning," *arXiv preprint arXiv:1906.07805*, 2019.

[212] J. Achiam, H. Edwards, D. Amodei, and P. Abbeel, "Variational option discovery algorithms," *arXiv preprint arXiv:1807.10299*, 2018.

[213] V. H. Pong, M. Dalal, S. Lin, A. Nair, S. Bahl, and S. Levine, "Skew-fit: State-covering self-supervised reinforcement learning," in *ICML*, 2020.

[214] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine, "Reinforcement learning with deep energy-based policies," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 70, 2017, pp. 1352–1361.

[215] T. Gangwani, Q. Liu, and J. Peng, "Learning self-imitating diverse policies," in *ICLR*, 2019.

[216] Y. Liu, P. Ramachandran, Q. Liu, and J. Peng, "Stein variational policy gradient," in *UAI*, 2017.

[217] Y. Guo, J. Choi, M. Moczulski, S. Bengio, M. Norouzi, and H. Lee, "Efficient exploration with self-imitation learning via trajectory-conditioned policy," *CoRR*, vol. abs/1907.10247, 2019.

[218] T. Dai, H. Liu, and A. A. Bharath, "Episodic self-imitation learning with hindsight," *CoRR*, vol. abs/2011.13467, 2020.

[219] Y. Tang, "Self-imitation learning via generalized lower bound q-learning," in *NeurIPS*, 2020.

[220] Z. Chen and M. Lin, "Self-imitation learning in sparse reward settings," *CoRR*, vol. abs/2010.06962, 2020. [Online]. Available: https://arxiv.org/abs/2010.06962

[221] J. Ferret, O. Pietquin, and M. Geist, "Self-imitation advantage learning," in *AAMAS*, 2021, pp. 501–509.

[222] K. Ning, H. Xu, K. Zhu, and S. Huang, "Co-imitation learning without expert demonstration," *CoRR*, vol. abs/2103.14823, 2021.