

Intelligent Reflecting Surface Configurations for Smart Radio Using Deep Reinforcement Learning

Wei Wang, *Member, IEEE*, and Wei Zhang, *Fellow, IEEE*

Abstract—Intelligent reflecting surface (IRS) is envisioned to change the paradigm of wireless communications from “adapting to wireless channels” to “changing wireless channels”. However, current IRS configuration schemes, consisting of sub-channel estimation and passive beamforming in sequence, conform to the conventional model-based design philosophies and are difficult to be realized practically in the complex radio environment. To create the smart radio environment, we propose a model-free design of IRS control that is independent of the sub-channel channel state information (CSI) and requires the minimum interaction between IRS and the wireless communication system. We firstly model the control of IRS as a Markov decision process (MDP) and apply deep reinforcement learning (DRL) to perform real-time coarse phase control of IRS. Then, we apply extremum seeking control (ESC) as the fine phase control of IRS. Finally, by updating the frame structure, we integrate DRL and ESC in the model-free control of IRS to improve its adaptivity to different channel dynamics. Numerical results show the superiority of our proposed joint DRL and ESC scheme and verify its effectiveness in model-free IRS control without sub-channel CSI.

I. INTRODUCTION

Metasurfaces, which consist of artificially periodic or quasi-periodic structures with sub-wavelength scales, are a new design of functional materials [1], [2]. Some extraordinary electromagnetic properties observed on metasurfaces, e.g., negative permittivity and permeability, reveal its potential in tailoring electromagnetic waves in a wide frequency range, from microwave to visible light [3]–[6]. Intelligent reflecting surface (IRS), a.k.a. reconfigurable intelligent surface (RIS), is a type of programmable metasurfaces that is capable of electronically tuning electromagnetic wave by incorporating active components into each unit cell of metasurfaces [7]–[12]. The advent of IRS is envisioned to revolutionize many industries, a major one of which is wireless communications.

Wireless communications are subject to the time-varying radio propagation environment. The effects of free space path loss, signal absorption, reflections, refractions, and diffractions caused by physical objects during the propagation of electromagnetic waves jointly render wireless channels highly dynamic [13]–[15]. IRS’s capability of manipulating electromagnetic waves in a real time manner brings infinite possibilities

to wireless communications and makes it possible for human beings to transform the design paradigm of wireless communications from “adapting to wireless channels” to “changing wireless channels” [7]. To this end, great research efforts have been spent to acquire the channel state information (CSI) of the sub-channels, i.e., the channels between wireless transceivers and IRS, which is widely regarded as the prerequisite of IRS reflection pattern (passive beamforming) design [8]. Owing to the passive nature, IRS is unable to sense the incident signal, thereby rendering the estimation process far more complicated than traditional wireless communication systems. In [9], a channel estimation scheme with reduced training overhead is proposed by exploiting the inter-user channel correlation. In [10], [11], compressed sensing based channel estimation methods are proposed to estimate the channel responses between base station, IRS and a single-antenna user at mmWave frequency band. As the proposed schemes are focusing on a single-antenna user, their extension to multiple users with array antenna might increase multi-fold the training overhead. In [16], a joint beam training and positioning scheme is proposed to estimate the parameters of the line-of-sight (LoS) paths for IRS assisted mmWave communications. The proposed random beamforming in the training stage is performed in a broadcasting manner and thus the training overhead is independent of the user number.

Despite the aforementioned endeavors to advance the CSI acquisition techniques in IRS assisted wireless communications, the practical applications of IRS are still confronting various challenges. Firstly, channel estimations for IRS assisted wireless communications demand a radical update of the existed protocols to incorporate the coordination of the transmitter, the receiver, and the IRS. It indicates that the existing wireless systems, e.g., Wi-Fi, 4G-LTE and 5G-NR, are unable to readily embrace IRS. Secondly, even if the perfect CSI is available, the real-time optimization of IRS reflection coefficients using convex and non-convex optimization techniques is computationally prohibitive [12]. Thirdly, current solutions to IRS control, which consists of CSI acquisition and IRS reflection designs, are based on the accurate modelling of IRS. However, as a type of low-cost reflective metasurfaces, IRS changes its reflection coefficient via tuning the impedance, the exact value of which is dependant on the carrier frequency of the incident signal [17], [18]. The carrier frequency can be shifted by Doppler effects and might also vary over different users, thus the mathematical modelling of IRS in the complex radio propagation environment is inherently difficult.

To tackle the aforementioned challenges, we follow the design paradigm of model-free control by treating the wireless

This work was supported in part by National Key R&D Program of China under Grant 2020YFA0711400, Shenzhen Science & Innovation Fund under Grant JCYJ20180507182451820, and the Australian Research Council’s Project funding scheme under LP160101244.

W. Wang is with Peng Cheng Laboratory, Shenzhen, China (e-mail: wangw01@pcl.ac.cn).

W. Zhang is with School of Electrical Engineering and Telecommunications, The University of New South Wales, Sydney, NSW 2052, Australia (e-mail: w.zhang@unsw.edu.au).

communication system as a (semi) black box with uncertain parameters and to optimize reflection coefficients of IRS through deep reinforcement learning (DRL) and extremum seeking control (ESC). Compared with the prevailing designs of IRS assisted wireless communication systems [11], [19]–[26], our proposed scheme is model-free. Specifically, the instantaneous (or statistical) CSI of sub-channels (i.e., Tx-Rx channel, Tx-IRS channel, and IRS-Rx channel) that constitutes the equivalent wireless channel is not required. Our design, in a true sense, treats IRS as a part of the wireless channel and requires the minimum interaction with wireless communication systems. The disentanglement of IRS configuration from wireless communication system means the improved independence of IRS and will speed up the rollout of IRS in the future. There are already some attempts towards the standalone operation of IRS [27]–[29]. In [27], [28], deep learning and deep reinforcement learning are applied to guide the IRS to interact with the incident signal given the knowledge of the sampled channel vectors. However, in order to obtain the CSI of the Tx-IRS and IRS-Rx sub-channels, the authors propose to install channel sensors on the IRS, which is, to some extent, against the initial role of IRS as a passive device. In [29], to reduce the dependence on the CSI of sub-channels, a deep learning scheme is proposed to extract the interactions between phase shifts of IRS and receiver locations. However, the training data has to be collected offline, which limits its adaptability to the more general scenarios.

In this paper, our objective is to build a model-free IRS control scheme with a higher level understanding of the radio environment, which is able to configure the IRS reflection coefficients without the CSI of the sub-channels. To this end, we adopt a typical scenario, i.e., time-division duplexing (TDD) multi-user multiple-input-multiple-output (MIMO), as an example to perform our design. To summarize, our contributions are as follows.

- We model the control of IRS as a Markov decision process (MDP) and then apply DRL, specifically, double deep Q-network (DDQN) method, to perform real-time coarse phase control of IRS. The proposed DDQN scheme outperforms the other sub-channel-CSI-independent methods, e.g., multi-armed bandit (MAB), random reflection.
- To enhance the action of DDQN, we further apply ESC as the fine phase control of IRS. Specifically, we propose a dither-based iterative method to optimize the phase shift of IRS through trial and error. We also prove that the output of the proposed dither-based iterative method is monotonically increasing.
- By updating the frame structure, we integrate DRL and ESC in the model-free control of IRS. The integrated scheme is more adaptive to various channel dynamics and has the potential to achieve better performance.

Numerical results show the superiority of our proposed DRL, ESC, and joint DRL and ESC scheme and verify their effectiveness in model-free IRS control without sub-channel CSI.

The rest of the paper is organized as follows. In Section

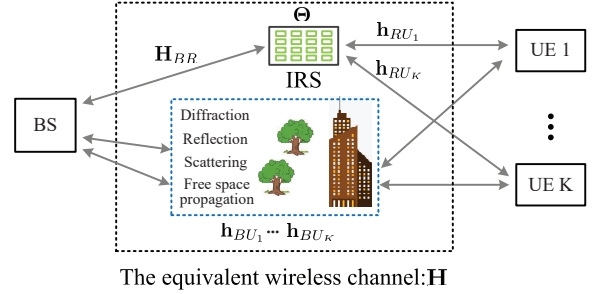


Fig. 1. The illustration of IRS assisted wireless communications

II, we introduce the system model. In Section III, we propose a DRL enabled model-free control of IRS. In Section IV, we propose a dither-based iterative method to enhance the action of DRL. In Section V, we present numerical results. Finally, in Section VII, we draw the conclusion.

Notations: Column vectors (matrices) are denoted by bold-face lower (upper) case letters, $\mathbf{x}[n]$ denotes the n -th element in the vector \mathbf{x} , \odot represents the Hadamard product, $(\cdot)^*$, $(\cdot)^T$ and $(\cdot)^H$ represent conjugate, transpose and conjugate transpose operation, respectively.

II. SYSTEM MODEL

In this section, we introduce the system model of model-free IRS control.

A. Optimal Phase Shift Vector of IRS

For IRS assisted wireless communications, the channel model between BS and a certain user k can be represented as

$$\mathbf{h}_k = \mathbf{h}_{BU_k} + \mathbf{H}_{BR} \Theta \mathbf{h}_{RU_k} \quad (1)$$

where the user k is equipped with a single antenna, $\mathbf{h}_{BU_k} \in \mathbb{C}^{N_B \times 1}$ is the channel response vector between the user k and BS, $\mathbf{H}_{BR} \in \mathbb{C}^{N_B \times N_R}$ is the channel response matrix between BS and IRS, $\mathbf{h}_{RU_k} \in \mathbb{C}^{N_R \times 1}$ is the channel response vector between the user k and IRS, $\Theta = \text{diag}\{\boldsymbol{\theta}\}$, and $\boldsymbol{\theta} \in \mathbb{C}^{N_R \times 1}$ (with $|\boldsymbol{\theta}(n)| = 1$) is the phase shift vector of the IRS. Accordingly, the multi-user channel is written as

$$\mathbf{H} = \mathbf{H}_{BU} + \mathbf{H}_{BR} \Theta \mathbf{H}_{RU} \quad (2)$$

where

$$\begin{aligned} \mathbf{H} &= [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_K] \\ \mathbf{H}_{BU} &= [\mathbf{h}_{BU_1}, \mathbf{h}_{BU_2}, \dots, \mathbf{h}_{BU_K}] \\ \mathbf{H}_{RU} &= [\mathbf{h}_{RU_1}, \mathbf{h}_{RU_2}, \dots, \mathbf{h}_{RU_K}] \end{aligned}$$

And the relationship between the aggregated equivalent channel \mathbf{H} and the sub-channels is shown in Fig. 1.

The objective of reinforcement learning based IRS configuration is to develop a widely compatible method that can be deployed in various scenarios of wireless communications

without any knowledge of the wireless system's internal working mechanism. Mathematically, the problem is formulated as

$$\begin{aligned} \max_{\boldsymbol{\theta}} \quad & P_m \\ \text{s.t.} \quad & \boldsymbol{\theta}[n] = e^{-j\varphi[n]}, \forall n \in \{1, 2, \dots, N_R\} \\ & \varphi[n] \in \mathcal{B}, \forall n \in \{1, 2, \dots, N_R\} \end{aligned} \quad (3)$$

where P_m is the performance metric of the wireless system that is to be optimized. P_m is dependent on the wireless channel \mathbf{H} , and \mathbf{H} is dependent on the reflection pattern $\boldsymbol{\theta}$. $\varphi[n]$ is the quantized phase selected from a finite set $\mathcal{B} = \{-\pi, \frac{-2^r+2}{2^r}\pi, \frac{-2^r+4}{2^r}\pi, \dots, \pi\}$ with $2^r + 1$ possible values.

It is worth mentioning that the model-free control does not need to know the exact relationship between the objective P_m and variable $\boldsymbol{\theta}$.

B. A Typical Scenario – TDD Multi-User MIMO

Without loss of generality, we use a typical scenario in wireless communication, i.e., TDD multi-user MIMO, to illustrate our design philosophy. In TDD, by exploiting the channel reciprocity, the BS can estimate the downlink channel from the pilot of the uplink channel. Thus, TDD multi-user MIMO consists of two stages (refer to Fig. 2), i.e., *uplink pilot transmission* and *downlink data transmission* [30], [31].

At uplink stage, the pilot transmits from multiple users to BS simultaneously. The received pilot signal is represented as

$$\mathbf{Y}_U = \mathbf{H}\mathbf{S} + \mathbf{N} \quad (4)$$

where $\mathbf{S} \in \mathbb{C}^{K \times K}$ is the pilot pattern, $\mathbf{N} \in \mathbb{C}^{N_B \times K}$ is the additive white Gaussian noise. Upon receiving the pilot, BS performs minimum mean square error (MMSE) estimation of the channel matrix, i.e.,

$$\hat{\mathbf{H}} = \mathbf{Y}_U \mathbf{S}^H (\mathbf{S} \mathbf{S}^H + \sigma_U^2 \mathbf{I})^{-1} \quad (5)$$

When \mathbf{S} is an unitary matrix, (5) is further expressed as

$$\hat{\mathbf{H}} = \frac{\mathbf{Y}_U \mathbf{S}^H}{1 + \sigma_U^2} \quad (6)$$

At downlink stage, data transmission with zero-forcing (ZF) precoding is performed, and the precoding matrix is represented as

$$\mathbf{M} = [\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_K]^H \quad (7a)$$

$$= \mathbf{D}_P (\hat{\mathbf{H}}^H \hat{\mathbf{H}})^{-1} \hat{\mathbf{H}}^H \quad (7b)$$

where $\mathbf{D}_P = \text{diag}(\frac{1}{\|\mathbf{m}_1\|_2}, \frac{1}{\|\mathbf{m}_2\|_2}, \dots, \frac{1}{\|\mathbf{m}_K\|_2})$ is for power normalization. The received signal of user k is given by

$$y_{D,k} = \mathbf{m}_k^H \mathbf{h}_k x_k + \sum_{l \neq k} \mathbf{m}_l^H \mathbf{h}_k x_l + n_k \quad (8)$$

where x_k is the signal intended to user k ($\mathbb{E}(x_k) = 0$ and $\mathbb{E}(|x_k|^2) = 1, \forall k \in \{1, \dots, K\}$), and $n_k \sim \mathcal{CN}(0, \sigma_k^2)$ is the additive white Gaussian noise. Thus, the signal-to-noise ratio of the k -th user is

$$\text{SINR}_k = \frac{|\mathbf{m}_k^H \mathbf{h}_k|^2}{\sum_{l \neq k} |\mathbf{m}_l^H \mathbf{h}_k|^2 + \sigma_k^2} \quad (9)$$

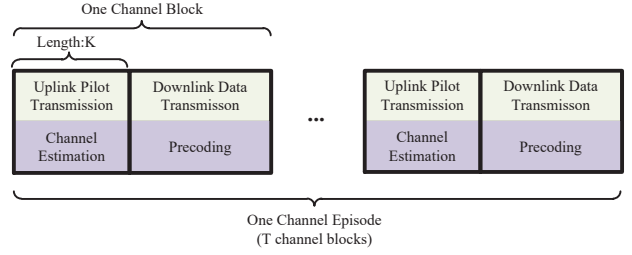


Fig. 2. The frame structure of the typical TDD multi-user MIMO (green part) and the corresponding signal processing procedures (purple part)

For a communication system, the performance metrics can be SINR, data rate, frame error rate (FER), and etc. Without loss of generality, we adopt the sum data rate as the performance metric, i.e.,

$$P_m = \sum_{k=1}^K r_k = \sum_{k=1}^K \log_2(1 + \text{SINR}_k) \quad (10)$$

C. Channel Model

We assume the Rician channel model for \mathbf{h}_{BU_k} , \mathbf{H}_{BR} and \mathbf{h}_{RU_k} . Take \mathbf{H}_{BR} as an example, it is represented as

$$\mathbf{H}_{BR} = \sqrt{\frac{K}{K+1}} \mathbf{H}_{BR,LoS} + \sqrt{\frac{1}{K+1}} \mathbf{H}_{BR,NLoS} \quad (11)$$

where $\mathbf{H}_{BR,LoS}$ denotes the deterministic LoS component, $\mathbf{H}_{BR,NLoS}$ denotes the fast fading NLoS component, and the component of which is independent and identically distributed (i.i.d.) circularly symmetric complex Gaussian random variables with zero-mean and unit variance, and K is the ratio between the power in the LoS path and the power in the NLoS paths [32].

The LoS component is position-dependent and is thus slow-time-varying; The NLoS components are caused by the multipath effects and are thus fast-time-varying [33]. Combining the characteristics of wireless channel with the setting of reinforcement learning, we introduce the following two concepts.

(1) **Channel block:** One channel block consists of the uplink pilot transmission stage and downlink data transmission stage (as shown in Fig. 2), and the channel matrix is constant during the channel block.

(2) **Channel episode:** One channel episode consists of T channel blocks (as shown in Fig. 2). The LoS component within one channel episode remains constant; The NLoS components change over time, and the NLoS components of different channel blocks are i.i.d.

III. MODEL-FREE IRS CONTROL ENABLED BY DEEP REINFORCEMENT LEARNING

In this section, we apply DRL to model-free IRS control.

A. Design Objectives

We aim to achieve stand-alone operation of the wireless communication system and the IRS, and our design includes the following characteristics.

- **Wireless Communication System:** The wireless communication system is almost unaware of the existence of the IRS, except that it needs to feed back its instantaneous performance to the IRS controller. In this regard, the uplink pilot transmission and downlink data transmission exactly follow the conventional structure in Fig. 2.
- **IRS:** The IRS is strictly regarded as part of the wireless channel and will not be jointly designed with the wireless communication system. The configuration of IRS is based on (a) the performance feedback from the wireless system and (b) its learned policy through trial-and-error interaction with the dynamic environment. And the IRS is unaware of the working mechanism of the wireless communication system.

A salient advantage of the proposed design is that the IRS can be deployed in various wireless communication applications, e.g., Wi-Fi, 4G-LTE, 5G-NR, without updating their existing protocols, which will speed up the roll-out of IRSs. Another benefit is that, by treating the existing wireless communication system as a black box, the configuration of IRS does not require the overhead-demanding channel sounding process to acquire the CSI of the subchannels, i.e., \mathbf{H}_{BU} , \mathbf{H}_{BR} , and \mathbf{H}_{RU} , that constitute the aggregated equivalent channel \mathbf{H} .

Our design is primarily based on the *reinforcement learning* technology. Specifically, the IRS and its controller are the *agent*, the wireless communication system, which comprises transmitter, wireless channel, and receiver, is the *environment*. The relationship between the different parties is given in Fig. 3. Initially, the agent takes random *actions* and the environment responds to those actions by giving rise to rewards and presenting new situations to the agent [34], [35]. Through trial-and-error interaction with the wireless communication system, the agent gradually learns the optimal policy to maximize the expected return over time. In this regard, IRS, which is capable of changing the radio environment, is analogous to the human body, and the reinforcement learning method, which guides the action of IRS, is analogous to the human brain. The integration of the IRS and the reinforcement learning method is the pathway to creating the smart radio environment.

B. Basics of Deep Reinforcement Learning

To facilitate the presentation of our design, we briefly introduce some key concepts of DRL in this subsection.

1) *Objective of Reinforcement Learning:* An MDP is specified by 4-tuple $\langle \mathcal{S}, \mathcal{A}, P, R \rangle$, where \mathcal{S} is the state space, \mathcal{A} is the action space, P is the state transition probability, and R is the immediate reward received by the agent. When an agent in the state $s \in \mathcal{S}$ takes the action $a \in \mathcal{A}$, the environment will evolve to the next state $s' \in \mathcal{S}$ with probability $P(s'|s, a) = \Pr(S_{t+1} = s' | S_t = s, A_t = a)$, and in the meantime, the agent will receive the immediate reward $R_{s \rightarrow s'}^a$. Adding the time index to $\mathcal{S}, \mathcal{A}, R$, the evolution of an MDP can be represented using the following trajectory

$$\langle S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T, S_T, \dots \rangle \quad (12)$$

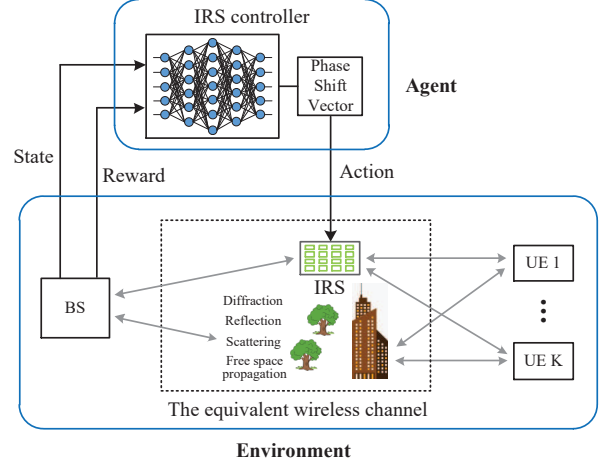


Fig. 3. The structure of model-free IRS configuration enabled by deep reinforcement learning

The agent's action is directed by the policy function

$$\pi(a|s) = \Pr(A_t = a | S_t = s) \quad (13)$$

which is the probability that the agent takes action a when the current state is s . A reinforcement learning task intends to find a policy that achieves a good return over the long run, where the return is defined as the cumulative discounted future reward, i.e.,

$$U_t = \sum_{\tau=0}^{\infty} \gamma^{\tau} R_{t+\tau+1}$$

and $\gamma \in [0, 1]$ is the discount factor for future rewards. Owing to the randomness of state transition (caused by the dynamic environment) and action selection, the return U_t is a random variable. Mathematically, the agent's goal in reinforcement learning is to find a good policy that maximizes the expected return, i.e.,

$$\max_{\pi} \mathbb{E}(U_t) \quad (14)$$

2) *Action-Value Function and Optimal Policy:* One key metric for action selection in reinforcement learning is action-value function, i.e.,

$$Q_{\pi}(s, a) = \mathbb{E}[U_t | S_t = s, A_t = a] \quad (15)$$

which is the conditional expected return for an agent to pick action a in the state s under the policy π . For any policy π and any state $s \in \mathcal{S}$, action-value function satisfies the following recursive relationship, i.e.,

$$\begin{aligned} Q_{\pi}(s, a) &= \mathbb{E}_{s'} \left[R_{s \rightarrow s'}^a + \gamma \sum_{a' \in \mathcal{A}} \pi(a'|s') Q_{\pi}(s', a') \mid s, a \right] \\ &= \sum_{s' \in \mathcal{S}} P(s'|s, a) \left(R_{s \rightarrow s'}^a + \gamma \sum_{a' \in \mathcal{A}} \pi(a'|s') Q_{\pi}(s', a') \right) \end{aligned} \quad (16)$$

where $R_{s \rightarrow s'}^a$ is the immediate reward when the environment transits from state s to state s' after taking the action a , and

Eq. (16) is the well-known Bellman equation of action-value function [34].

A policy is defined to be better than another policy if its expected return is greater than that of for all states and all actions. Thus, the optimal action-value function is

$$Q^*(s, a) := Q_{\pi^*}(s, a) = \max_{\pi} Q_{\pi}(s, a), \quad \forall s \in \mathcal{S}, a \in \mathcal{A} \quad (17)$$

With the optimal action-value function, the optimal policy is obviously

$$\pi^*(a|s) = \begin{cases} 1, & \text{if } a = \arg \max_{a \in \mathcal{A}} Q^*(s, a) \\ 0, & \text{otherwise} \end{cases} \quad (18)$$

Combining (17), (18) with (16), the Bellman optimality equation for $Q^*(s, a)$ is given by

$$Q^*(s, a) = \sum_{s' \in \mathcal{S}} P(s'|s, a) (R_{s \rightarrow s'}^a + \gamma \max_{a' \in \mathcal{A}} Q^*(s', a')) = \mathbb{E}_{s'} [R_{s \rightarrow s'}^a + \gamma \max_{a' \in \mathcal{A}} Q^*(s', a') | s_t = s, a_t = a] \quad (19)$$

With the Bellman optimality equation, the optimal policy $\pi^*(a|s)$ or the optimal action-value function $Q^*(s, a)$ can be obtained via iterative methods, i.e., policy iteration based methods and value iteration based methods [34]. Hereinafter, we will mainly focus on value iteration based methods.

3) *Temporal Difference Learning*: The aforementioned iterative methods require the complete knowledge of the environment, i.e., state transition probability $p(s'|s, a)$, reward function $R_{s \rightarrow s'}^a$, etc. However, the explicit knowledge of environment dynamics is unavailable in practice. The conditional expectation in (19) can be realized via numerically averaging over the sample sequences of states, actions, and rewards from actual interaction with the environment, e.g., temporal difference method, or Monte Carlo method.

Upon observing a new segment of the trajectory in (12), i.e., $\langle S_t = s, A_t = a, R_{t+1} = R_{s \rightarrow s'}^a, S_{t+1} = s' \rangle$, the action-value function $Q(s, a)$ updates as follows:

$$Q_{t+1}(s, a) = Q_t(s, a) + \alpha \left(R_{s \rightarrow s'}^a + \gamma \max_{a' \in \mathcal{A}} Q_t(s', a') - Q_t(s, a) \right) \quad (20)$$

where $\alpha \in (0, 1]$ is the learning rate and the following term inside the bracket is the error between the estimated Q value and the return. It means that the value function is updated in the direction of the error, and iteration will terminate when the error becomes infinitesimal.

4) *Double Deep Q-Network (DDQN)*: When the state s and the action a are both discrete, the optimal state-action function $Q^*(s, a)$ can be obtained as a lookup table, which is also known as Q-table [36], following the iterative procedures in (20). However, the size of the state (or action) space can be prohibitively large, and the state (or action) can even be continuous. In such cases, it is impractical to represent $Q(s, a)$ as a lookup table. Fortunately, the deep neural network (DNN) can be adopted to approximate the Q-table as $Q(s, a) \approx Q(s, a; \mathbf{w})$, which enables reinforcement learning to scale to more generalized decision-making problems. The

coefficients \mathbf{w} of $\tilde{Q}(s, a; \mathbf{w})$ are the weights of the DNN, and the DNN is termed as deep Q-network (DQN) [36], [37].

The trajectory segment $\langle S_t = s, A_t = a, R_{t+1} = R_{s \rightarrow s'}^a, S_{t+1} = s' \rangle$ in (12) constitutes an ‘‘experience sample’’ that will be used to train the DQN, and in accordance with (20), the loss function adopted during the training process of DQN is

$$Loss = \left(\underbrace{R_{s \rightarrow s'}^a + \gamma \max_{a' \in \mathcal{A}} \tilde{Q}(s', a'; \mathbf{w})}_{T_{DQN}} - \tilde{Q}(s, a; \mathbf{w}) \right)^2 \quad (21)$$

where T_{DQN} is the target value fed to the network.

The target T_{DQN} is dependent on the immediate reward $R_{s \rightarrow s'}^a$, as well as the output of the DQN $\tilde{Q}(s', a'; \mathbf{w})$. Such structure will inevitably result in over-estimation of the action state value (a.k.a., Q value) during the training process and thus will significantly degrade the performance of DRL. To mitigate over-estimation, we will adopt the double DQN (DDQN) structure [38], [39] in our design.

The fundamental idea of DDQN is to apply a separate target network $\tilde{Q}(s', a'; \mathbf{w}^-)$ to estimate the target value [39], and the expression of target in DDQN is

$$T_{DQN} = R_{s \rightarrow s'}^a + \gamma \tilde{Q}(s', \arg \max_{a' \in \mathcal{A}} \tilde{Q}(s', a'; \mathbf{w}); \mathbf{w}^-) \quad (22)$$

To summarize, DDQN differs from DQN in the following two aspects, i.e., (1) the optimal action is selected using the DQN $\tilde{Q}(s', a'; \mathbf{w})$ whose weights are \mathbf{w} , and (2) the Q value of the target value is taken from the target network whose weights are \mathbf{w}^- .

C. Model-Free Control of IRS Using Deep Reinforcement Learning

To apply reinforcement learning to model-free IRS configuration, we firstly model IRS assisted wireless communications as an MDP.

- *Agent*: The agent is IRS controller, which is capable of autonomously interacting with the environment via IRS to meet the design objectives.
- *Environment*: The environment refers to the things that the agent interact with, which includes BS, wireless channel, IRS, and mobile users.
- *State*: To facilitate the accurate prediction of expected next rewards and next states given an action, we define the state as $\{\mathbf{H}, \boldsymbol{\theta}\}$, which consists of two sub-states, namely the equivalent wireless channel \mathbf{H} and the reflection vector $\boldsymbol{\theta}$ of IRS.
- *Action*: The action is defined as the incremental phase shift of the current reflection pattern, i.e.,

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} \odot \Delta \boldsymbol{\theta}^{(t)} \quad (23)$$

where \odot is the Hadamard (element-wise) product, $\boldsymbol{\theta}^{(t)}$ is the reflection pattern at the t -th channel block, and $\Delta \boldsymbol{\theta}^{(t)}$ is the incremental phase shift of $\boldsymbol{\theta}^{(t)}$. We use the subset (or full set) of the discrete Fourier transform (DFT) vectors as the action set. For example, when the size of action space is 5, we set

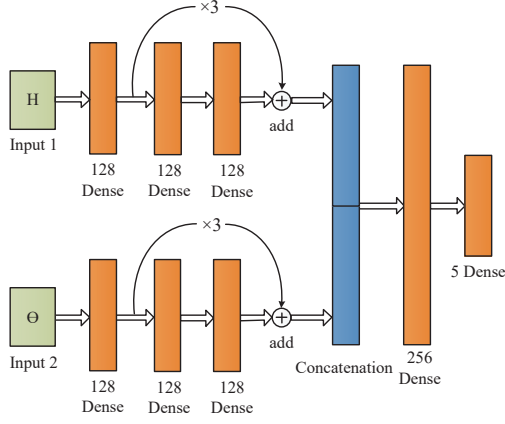


Fig. 4. Structure of the DQN

$\mathcal{A} = \left\{ \mathbf{v}\left(-\frac{6}{N_R}\right), \mathbf{v}\left(-\frac{2}{N_R}\right), \mathbf{v}(0), \mathbf{v}\left(\frac{2}{N_R}\right), \mathbf{v}\left(\frac{6}{N_R}\right) \right\}$, where $\mathbf{v}(\Psi_R)$ is the steering vector¹, i.e.,

$$\mathbf{v}(\Psi_R) = \left[1, e^{j\pi\Psi_R}, \dots, e^{j(N_R-1)\pi\Psi_R} \right]^T$$

When $\Delta\theta^{(t)} = \mathbf{v}(0)$, the sub-state θ stays unchanged, and the sub-state \mathbf{H} changes merely due to the variation of NLoS components; $\Delta\theta^{(t)} = \mathbf{v}\left(-\frac{2}{N_R}\right)$ and $\Delta\theta^{(t)} = \mathbf{v}\left(\frac{2}{N_R}\right)$ are towards the opposite directions, which enables the agent to quickly correct from a negative action; $\Delta\theta^{(t)} = \mathbf{v}\left(-\frac{6}{N_R}\right)$ and $\Delta\theta^{(t)} = \mathbf{v}\left(\frac{6}{N_R}\right)$ are used to speed up the transition of reflection pattern.

- *Reward*: The immediate reward after transition from s to s' with action a is defined as

$$R = \begin{cases} P_m, & \text{when } P_m \geq P_{th} \\ P_m - 100, & \text{when } P_m < P_{th} \end{cases} \quad (24)$$

where P_{th} is a performance threshold. When P_m is less than P_{th} , we add an penalty -100 to encourage the IRS to maximize performance, while maintaining an acceptable performance above the threshold.

Remark 1. The reasons for using incremental phase shift, rather than the absolute phase shift, as the action are two-fold. On one hand, we need to build the Markov property of the state transmission, and, on the other hand, we intend to reduce the size action space and accelerate convergence rate.

Based on the modeled MDP and the basics of DRL presented in Subsection B, we propose to maximize the expected return (cumulative discounted future reward) using Algorithm 1. Some of the key techniques applied in Algorithm 1 are explained as follows.

1) *DDQN*: Different from the naive DQN method, where the DQN $\tilde{Q}(s, a; \mathbf{w})$ (with the weights \mathbf{w}) is used to generate the target value, we use a separate target network $\tilde{Q}(s, a; \mathbf{w}^-)$ (with the weights \mathbf{w}^-) to generate the target value, and the weights of the target network are updated by $\mathbf{w}^- = \mathbf{w}$ in every N_{TNet} time intervals. The structure of the network is shown

¹Without loss of generality, we assume that the reflector array of IRS is a uniform linear array (ULA).

Algorithm 1: Double DQN based model-free IRS control for IRS-assisted wireless communications

Initialize parameters s_0, ϵ ;
Initialize the FIFO memory \mathcal{M} with the size N_m ;
Initialize the weights of the DQN \mathbf{w} and set the target network as $\mathbf{w}^- = \mathbf{w}$
for $t = 0, 1, 2, \dots$ **do**
 Input s_t to the DQN and obtain the state-action values $\tilde{Q}(s_t, a; \mathbf{w}), a \in \mathcal{A}$;
 With $\tilde{Q}(s_t, a; \mathbf{w}), a \in \mathcal{A}$, select an action a_t using ϵ -greedy policy;
 Receive the reward r_{t+1} and the estimated channel response $\hat{\mathbf{H}}_{t+1}$, and compute the next state s_{t+1} from $\hat{\mathbf{H}}_{t+1}, s_t$ and a_t .
 Store the experience tuple $\langle s_t, a_t, r_{t+1}, s_{t+1} \rangle$ to the FIFO memory \mathcal{M} ;
 If $|\mathcal{M}| \geq N_e$
 Randomly select a mini batch of N_e experience tuples $\langle s_i, a_i, r_{i+1}, s_{i+1} \rangle$ from \mathcal{M} .
 Calculate the target values $T_{DQN,i}$ for the mini batch according to (22).
 With the input $\{s_i\}$ and the output $\{\tilde{Q}(s_i, a; \mathbf{w})\}$, train the DQN, and update its weights \mathbf{w} .
 If $t \bmod N_{TNet} = 0$, update the weights of the target network, i.e., set $\mathbf{w}^- = \mathbf{w}$.
 end if
 end if
end for

in Fig. 4. Specifically, we apply the deep residual network (ResNet) [40] to process two sub-states (i.e., \mathbf{H} and θ), and then we fuse the processed information of the two sub-states using a two-layer dense network. The activation function that we use is the Swish function [41].

2) *ϵ -Greedy Policy*: Given the perfect $\tilde{Q}(s, a; \mathbf{w})$, the optimal policy is to select the action that yields the largest state-action value. However, the perfect $\tilde{Q}(s, a; \mathbf{w})$ demands for an infinite size of experiences, which is impractical and infeasible in the dynamic wireless environment. Therefore, it is necessary for the agent to keep exploring to avoid get stuck with a sub-optimal policy. To this end, we apply an ϵ -greedy policy. In ϵ -greedy policy, ϵ refers to the probability of choosing to explore, i.e., randomly select from all the possible actions, and $1 - \epsilon$ is the probability of choosing to exploit the obtained DQN in decision making. In this regard, the ϵ -greedy policy is represented as

$$\pi^\epsilon = \begin{cases} \pi^*(a/s), & w.p. 1 - \epsilon \\ P(a) = \frac{1}{|\mathcal{A}|}, & w.p. \epsilon \end{cases} \quad (25)$$

where $\pi^*(a/s)$ as the policy based on the Q-network, which is introduced in (18). In our design, ϵ is initially set to 1 and decreases exponentially at a rate of ϑ , ($0 < \vartheta < 1$) every time interval until it reaches the lower bound ϵ_{min} .

3) *Experience Replay*: Instead of training the DQN with the latest experience tuple, we store N_e recent experience tuples in the memory \mathcal{M} in ‘‘first in, first out’’ (FIFO) manner, i.e.,

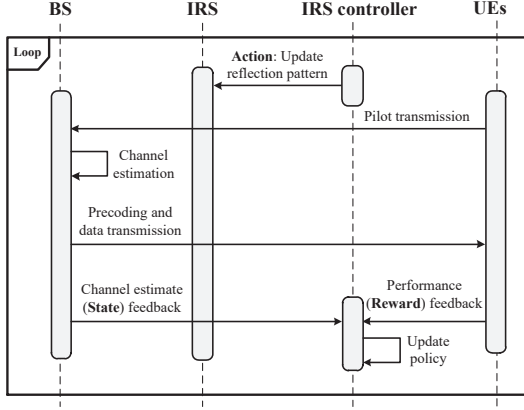


Fig. 5. Sequence diagram of the proposed model-free IRS configuration

queue data structure, and then randomly fetch a mini-batch of N_e experience samples from \mathcal{M} to train the DQN.

D. Summarizing The Work Flow of Model-Free IRS Configuration

In this subsection, we summarize the work flow of the proposed model-free IRS configuration. To this end, we plot the sequence diagram in Fig. 5.

According to Fig. 5, in a specific loop, the IRS is configured with a reflection pattern $(\theta^{(t+1)} = \theta^{(t)} \odot \Delta\theta^{(t)})$ according to the ϵ -greedy policy, and then UEs and BS perform uplink pilot transmission and downlink data transmission sequentially as if there exists no IRS. After that, BS sends the estimated channel matrix $(\hat{\mathbf{H}}^{(t+1)})$ to IRS controller, which can be fulfilled through wired communications, and UEs send back their performance metrics to the IRS controller. Finally, according to the received channel estimate $(\hat{\mathbf{H}}^{(t+1)})$ and performance feedback (P_m^{t+1}) , IRS controller derives the tuple $\langle \{\hat{\mathbf{H}}^{(t)}, \theta^{(t)}\}, \Delta\theta^{(t)}, R^{(t+1)}, \{\hat{\mathbf{H}}^{(t+1)}, \theta^{(t+1)}\} \rangle$ and stores it in the FIFO queue as the training data for the DQN $\tilde{Q}(s_t, a; w)$.

Compared with the traditional TDD multi-user MIMO, the extra efforts of incorporating IRS are merely the feedback of $\hat{\mathbf{H}}^{(t)}$ and P_m^{t+1} . The former can be easily achieved via wired communications between BS and IRS, and the latter costs negligible wireless communication resources of the mobile UEs. It is also noteworthy that IRS controller is unaware of the working mechanism of BS and UEs and does not require the CSI of the sub-channels.

IV. ENHANCING IRS CONTROL USING EXTREMUM SEEKING CONTROL

In DRL, action space is restrained for a fast convergence rate, which limits the phase freedom of IRS. To enhance the control of IRS, another model-free real-time optimization method, namely, extremum seeking control (ESC), is used to design the fine phase control of IRS.

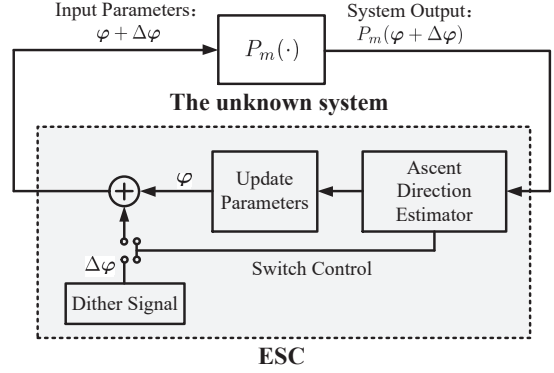


Fig. 6. Principle of the proposed ESC-inspired dither-based iterative method

A. Model-Free Control of IRS Using ESC

ESC is model-free method to realize a learning-based adaptive controller for maximizing/minimizing certain system performance metrics [42], [43]. The first application of ESC can be traced back to the work of the French engineer Leblanc in 1922 to maintain an efficient power transfer for a tram car [44]. The basic idea of ESC is to add a dither signal (e.g., sinusoidal signal [45], [46], and random noise [47]) to the system input and observing its effect on the output to obtain an approximate *implicit gradient* of a nonlinear static map of the unknown system [42], [48].

According to the design philosophy of ESC, we propose a dither-based model-free control of IRS as in Fig. 6. Our design consists of three parts, i.e., dither signal generation module, ascent direction estimation module, and parameter update module. Dither signal generation module generates random dither/perturbation signal to probe the response $P_m(\cdot)$ of the system; gradient estimation module determines the update direction of the system input according to the system performance $P_m(\varphi + \Delta\varphi)$ to guarantee the monotonic increase of the performance, and it also guides on-off switch of the random dither signal generation; parameter update module updates the system input φ according to the estimated direction.

Specifically, the iterative process in Fig. 6 runs as follows.

Step 1. Dither Signal Generation

Generate a small random dither signal through uniform random distribution², i.e.,

$$\Delta\varphi[n] = \frac{a}{2^{r-1}}, \quad a \in \mathcal{U} \left\{ -\frac{2^{r-1}}{N_R}, \frac{2^{r-1}}{N_R} \right\} \quad (26)$$

Then, add the dither signal $\Delta\varphi$ to the parameter φ , use $\varphi + \Delta\varphi$ as the input of the system, and receive the feedback of the performance metric $P_m(\varphi + \Delta\varphi)$.

Step 2. Direction Estimation and Parameter Update

Condition 1. If $P_m(\varphi + \Delta\varphi) \geq P_m(\varphi)$, adopt $\Delta\varphi$ as the direction. Update the parameter as

$$\varphi \leftarrow \varphi + \Delta\varphi, \quad (27)$$

²The parameter selection will be explained in the following context.

and update the performance metric as

$$P_m(\varphi) \leftarrow P_m(\varphi + \Delta\varphi) \quad (28)$$

Then, jump to Step 1 for the next iteration;

Condition 2. Else if $P_m(\varphi + \Delta\varphi) < P_m(\varphi)$, adopt $-\Delta\varphi$ as the direction. Update the parameter as set

$$\varphi \leftarrow \varphi - \Delta\varphi \quad (29)$$

Turn off the dither signal, use only φ as the system input, and measure the system performance $P_m(\varphi)$. Then, jump to Step 1 for the next iteration.

It is noteworthy that each iteration uses one or two time intervals, and each iteration can guarantee the monotonic increase of the performance metric P_m , which is validated in the following proposition.

Proposition 1. Each iteration in ESC based iterative process can guarantee the monotonic increase of the performance metric P_m given that the norm of the random dither signal, i.e., $\|\Delta\varphi\|$, is small enough.

Proof. To prove Proposition 1, it is essential to validate that the operation (29) in Condition 2 of Step 2 guarantees the increase of P_m , i.e., when $P_m(\varphi + \Delta\varphi) < P_m(\varphi)$, the following inequality

$$P_m(\varphi - \Delta\varphi) > P_m(\varphi)$$

holds.

To this end, we expand $P_m(\varphi + \Delta\varphi)$ using Taylor series of P_m with respect to φ , i.e.,

$$\begin{aligned} & P_m(\varphi + \Delta\varphi) \\ &= P_m(\varphi) + \frac{\partial P_m(\varphi)}{\partial \varphi^H} \Delta\varphi + \mathcal{O}(\|\Delta\varphi\|^2) \quad \text{as } \Delta\varphi \rightarrow 0 \end{aligned} \quad (30)$$

Since $\|\Delta\varphi\|$ is small, namely, $\|\Delta\varphi\| \rightarrow 0$, we adopt the first-order approximation, i.e.,

$$P_m(\varphi + \Delta\varphi) \approx P_m(\varphi) + \frac{\partial P_m(\varphi)}{\partial \varphi^H} \Delta\varphi \quad (31)$$

As it is reported by the system that $P_m(\varphi + \Delta\varphi) < P_m(\varphi)$ in Condition 2, we have

$$\frac{\partial P_m(\varphi)}{\partial \varphi^H} \Delta\varphi < 0 \quad (32)$$

Then, it is easy to verify that

$$\begin{aligned} P_m(\varphi - \Delta\varphi) &\approx P_m(\varphi) - \frac{\partial P_m(\varphi)}{\partial \varphi^H} \Delta\varphi \\ &> P_m(\varphi) \end{aligned} \quad (33)$$

□

Remark 2. As $\theta[n] = e^{-j\pi\varphi[n]}$, the operations of (27) and (29) can be written w.r.t. θ as follows

$$\theta \leftarrow \theta \odot \Delta\theta \quad (34a)$$

$$\theta \leftarrow \theta \odot \Delta\theta^* \quad (34b)$$

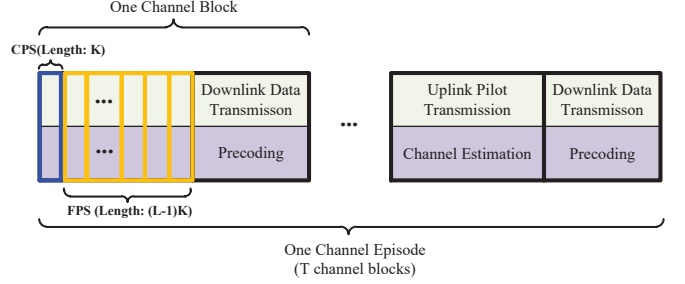


Fig. 7. The upgraded frame structure for the integration of ESC into DRL

B. Comparison with Gradient Ascent Search

To obtain further insights into the proposed dither-based method, we make a comparison with the well-known iterative algorithm – gradient ascent (descent in minimization problems) search.

For gradient ascent search algorithm, φ in each iteration is updated as follows.

$$\varphi \leftarrow \varphi + \underbrace{\gamma \frac{\partial P_m(\varphi)}{\partial \varphi}}_{\Delta\varphi} \quad (35)$$

When the step size γ is small, the iteration will almost surely guarantee the increase of $P_m(\varphi)$, because

$$\begin{aligned} P_m(\varphi + \Delta\varphi) &\approx P_m(\varphi) + \gamma \frac{\partial P_m(\varphi)}{\partial \varphi^H} \frac{\partial P_m(\varphi)}{\partial \varphi} \\ &= P_m(\varphi) + \gamma \left\| \frac{\partial P_m(\varphi)}{\partial \varphi} \right\|_2^2 \geq P_m(\varphi) \end{aligned} \quad (36)$$

Remark 3. As gradient ascent search take steps in the direction of the gradient, it is also called steepest descent. Thus, the convergence rate of gradient ascent search is faster than dither-based extremum search when the step size γ is properly selected. On the other hand, it is also noteworthy that gradient ascent search requires the exact expression of the gradient $\frac{\partial P_m(\varphi)}{\partial \varphi}$, whilst dither-based method is implemented through trial and error, which does not rely on any explicit knowledge of the wireless system's internal working mechanism.

C. Integrating ESC Into DRL

Recall that the action space for DRL is intentionally restrained for a fast convergence rate, whilst dither-based iterative method relies on the small-scale phase shift. Therefore, dither-based iterative method is complementary to the action in DRL and can be applied to enhance the action of DRL.

The **enhanced action** of DRL is defined as follows.

Step 1. Coarse phase shift (CPS) When $l = 0$, set

$$\theta_{temp}^{(t+1)} = \theta^{(t)} \odot \Delta\theta_c^{(t)} \quad (37)$$

where $\theta^{(t)}$ is the reflection pattern at the t -th channel block, $\Delta\theta_c^{(t)}$ is the coarse incremental phase shift at the t -th channel block, and $\theta_{temp}^{(t+1)}$ is the intermediate reflection pattern at the $t - 1$ -th channel block. Following the example in Section III.

C, the action set of the incremental phase shift $\Delta\theta_c^{(t)}$ is $\mathcal{A} = \left\{ \mathbf{v}\left(-\frac{6}{N_R}\right), \mathbf{v}\left(-\frac{2}{N_R}\right), \mathbf{v}(0), \mathbf{v}\left(\frac{2}{N_R}\right), \mathbf{v}\left(\frac{6}{N_R}\right) \right\}$.

Step 2. Fine phase shift (FPS)

For l from 1 to L , do

$$\boldsymbol{\theta}_{temp}^{(t+1)} \leftarrow \boldsymbol{\theta}_{temp}^{(t+1)} \odot \Delta\boldsymbol{\theta}_f \quad (38)$$

where $\Delta\boldsymbol{\theta}_f = \Delta\boldsymbol{\theta}$ (or $\Delta\boldsymbol{\theta}_f = \Delta\boldsymbol{\theta}^*$) is the ascent direction, and $\Delta\boldsymbol{\theta}$ is the random dither signal.

Remark 4. For example, when the quantization level $r = 8$, and $N_R = 32$, the step of coarse phase shift is $\frac{2\pi}{32}$, and the step of fine phase shift is $\frac{2\pi}{256}$ with the range $[-\frac{\pi}{32}, \frac{\pi}{32}]$.

To be compatible with the enhanced action, the frame structure needs to be updated as in Fig. 7. In the first K time slots, UEs transmit pilots and BS performs channel estimation with the reflection pattern in (37), and in the subsequent $(L - 1)K$ time slots, UEs repeatedly transmit pilots and BS performs channel estimation, while the reflection pattern updates as in (38). It is noteworthy that, as the performance feedback is done once per channel block, the performance metric used for the dither-based method is an approximation derived by replacing the authentic channel response \mathbf{H} in (10) with the channel estimate $\hat{\mathbf{H}}$.

Remark 5. The parameter L can be set adaptively according to the channel dynamics. For a practical wireless communication system, different values of L correspond to different modes.

V. NUMERICAL RESULTS

In this section, we present some numerical results to verify the effectiveness of our proposed model-free control of IRS³.

A. Simulation Parameters

The BS is equipped with a ULA that is placed along the direction $[1, 0, 0]$ (i.e., x-axis), and IRS is a ULA, which is placed along the direction $[0, 1, 0]$ (i.e., y-axis), UEs are equipped with a single antenna, the user number is $K = 2$, BS antenna number is $N_B = 2$, and IRS reflector number if $N_R = 32$. The element antennas/reflectors of BS and IRS are both with half wavelength spacing. The position of BS is $[0, 0, 10]$, the position of IRS is $[-2, 5, 5]$, and the UEs are uniformly distributed in the area $[0, 10] \times [0, 10]$ with the height being 1.5. The noise variance at BS side is $\sigma_B^2 = 0.1$, the noise variance at UE side is $\sigma_k^2 = 0.5, \forall k \in \{1, \dots, K\}$. Each channel episode consists of 20 channel blocks. In each channel episode, the LoS component is generated by randomly selecting the user locations within the area $[0, 10] \times [0, 10]$, and the LoS component is time-invariant within the 20 channel blocks of that channel episode.

The LoS channel between BS and IRS is

$$\mathbf{H}_{BR,LoS} = \mathbf{v}_R \mathbf{v}_B^H \quad (39)$$

³The simulation code is available at <https://github.com/WeiWang-WYS/IRSconfigurationDRL>

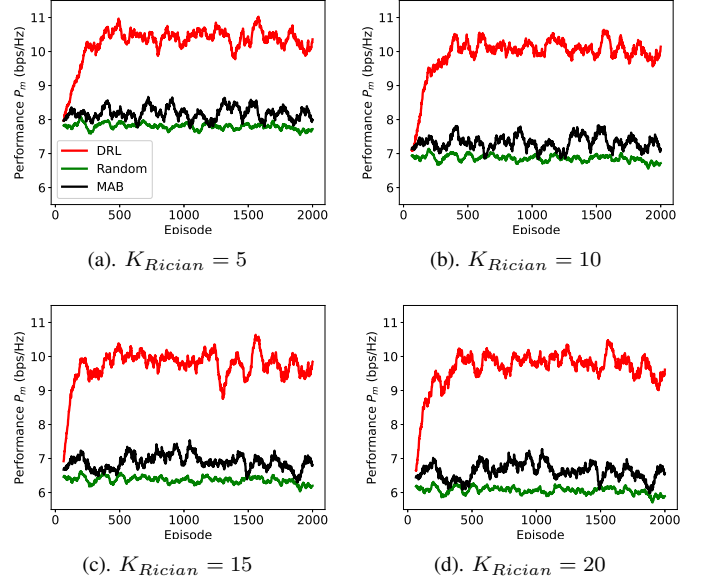


Fig. 8. Moving average of P_m for DRL under different values of Rician factor

where the steering vectors are represented as

$$\mathbf{v}_R = \mathbf{v}(\Psi_R, N_{R,y}) = \left[1, e^{j\pi\Psi_R}, \dots, e^{j(N_{R,y}-1)\pi\Psi_R} \right]^T$$

$$\mathbf{v}_B = \mathbf{v}(\Psi_B, N_{B,x}) = \left[1, e^{j\pi\Psi_B}, \dots, e^{j(N_{B,x}-1)\pi\Psi_B} \right]^T$$

and, according to [49], the directional cosines Ψ_R, Ψ_B are given by

$$\Psi_R = [0, 1, 0] \mathbf{e}_{BR} = \mathbf{e}_{BR}(2) \quad (41a)$$

$$\Psi_B = [1, 0, 0] \mathbf{e}_{BR} = \mathbf{e}_{BR}(1) \quad (41b)$$

where the direction vector \mathbf{e}_{BR} is determined by the relative position of BS and UE, i.e.,

$$\mathbf{e}_{BR} \triangleq \frac{\mathbf{p}_B - \mathbf{p}_R}{\|\mathbf{p}_B - \mathbf{p}_R\|_2} \quad (42)$$

The NLoS components are Gaussian distributed, i.e., $\mathbf{H}_{BR,NLoS}(\ell, \kappa) \in \mathcal{CN}(0, 1)$. The channel matrix \mathbf{H}_{BU} and \mathbf{H}_{RU} are generated in the same way.

B. Performance Study of DRL

In Fig. 8, we study the performance of the proposed DRL scheme (with the action set \mathcal{A}_5) under different values of Rician factor. The x-axis represents the episode, and the y-axis represents the moving average of P_m (window length is 64). As can be seen, the proposed DRL significantly outperforms the benchmark schemes, i.e., random reflection and multi-armed bandit (MAB). Different from DRL, the actions of random reflection and MAB that we used are absolute phase shift, and the action set is the DFT vectors. Although all the three schemes are independent of the sub-channel CSI, their utilizations of the other information are different. Random reflection is independent of any information and undoubtedly achieves the worst performance; MAB assumes a fixed distribution of *rewards* and explores the reward distributions of all arms. However, MAB fails to describe the state of the

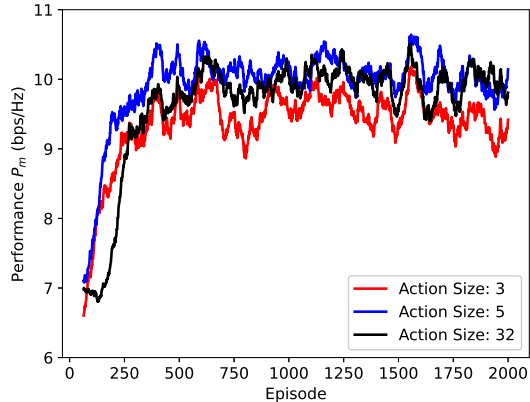
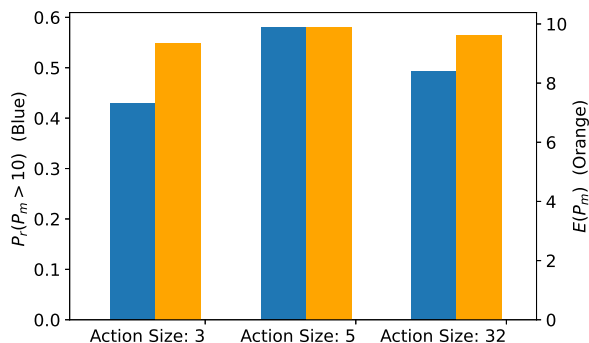
(a). Moving average of P_m for DRL(b). Average sum rate $\mathbb{E}(P_m)$ (orange color) and the probability $P_r(P_m > 10)$ (blue color)

Fig. 9. Performance of DRL with different action sets

environment and to build the connection between the action and the environment; DRL defines an appropriate *state* to represent the agent’s “position” within the environment and learns the quality of a state-action combination using the DQN from the information of *rewards* and *states*, which enables the agent to choose the best action to maximize the returns. We can also find that the performance gap of DRL and the benchmarks schemes becomes larger with the increase of Rician factor K_{Rician} , which indicates that the effectiveness of DRL is also dependent on the radio environment.

In Fig. 9, we study the impacts of action size to the performance of the proposed DRL scheme when the Rician factor is $K = 10$. In addition to the action set $\mathcal{A}_5 = \left\{ \mathbf{v}\left(-\frac{6}{N_R}\right), \mathbf{v}\left(-\frac{2}{N_R}\right), \mathbf{v}(0), \mathbf{v}\left(\frac{2}{N_R}\right), \mathbf{v}\left(\frac{6}{N_R}\right) \right\}$ defined in Section III, we adopt the action set $\mathcal{A}_3 = \left\{ \mathbf{v}\left(-\frac{2}{N_R}\right), \mathbf{v}(0), \mathbf{v}\left(\frac{2}{N_R}\right) \right\}$ and action set $\mathcal{A}_{32} = \left\{ \mathbf{v}(-1), \mathbf{v}\left(-1 - \frac{2}{N_R}\right), \dots, \mathbf{v}\left(1 - \frac{2}{N_R}\right) \right\}$ (namely, DFT matrix) as the benchmarks. From Fig. 9a, we can see that \mathcal{A}_5 is the fastest to converge, while \mathcal{A}_{32} is the slowest. Although a large action size will speed up the response rate of the agent, it will, on the other hand, demand more time for the DQN to converge. In the convergence region, we find that \mathcal{A}_5 and \mathcal{A}_{32} achieve the similar performance, while \mathcal{A}_3 ’s performance is inferior. It indicates that a well-design action

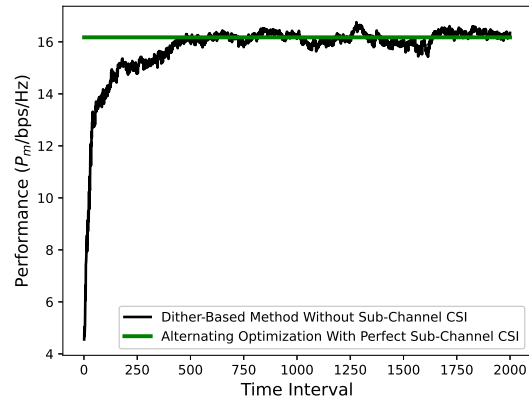


Fig. 10. Performance of ESC-inspired dither-based iterative method

set with a moderate size might be better than the small action size and the over-large action size. In Fig. 9b, the average sum rate $\mathbb{E}(P_m)$, and the probability $P_r(P_m > 10)$ are presented in the bar chart. For \mathcal{A}_5 , $\mathbb{E}(P_m) = 9.89$ bps/Hz and $P_r(P_m > 10) = 58.05\%$; for \mathcal{A}_3 , $\mathbb{E}(P_m) = 9.34$ bps/Hz and $P_r(P_m > 10) = 42.9\%$; for \mathcal{A}_{32} , $\mathbb{E}(P_m) = 9.60$ bps/Hz and $P_r(P_m > 10) = 49.35\%$. It further verifies the importance of action set design.

C. Performance Study of ESC

In Fig. 10, we study the performance of the ESC-inspired dither-based iterative method in a specific channel block. Each time interval of x -axis consists of K time slots, which is the pilot length required by the BS to estimate the multi-user channel \mathbf{H} . As can be seen that the performance metric P_m for dither-based method is almost monotonically increasing over time. Note that the performance metric used to guide the dither-based iterative method is an approximation, rather than the authentic feedback. Thus, the slight fluctuation of the performance curve is reasonable. For the purpose of comparison, we adopt the model-based method as the benchmark, in which the perfect sub-channel CSI, namely, \mathbf{H}_{BU} , \mathbf{H}_{BR} , and \mathbf{H}_{RU} , is available. The optimal reflection coefficient vector $\boldsymbol{\theta}$ is derived through solving the optimization problem (3). Recall that the difference from model-free control is that the exact relationship between the objective function P_m and the variable $\boldsymbol{\theta}$ is known in model-based methods. Due to the discrete nature of the feasible region \mathcal{B} , the optimization problem is intractable. Hence, we manage to solve it using the alternating optimization technique, which alternatively freezes $N_R - 1$ reflection coefficients and optimize only 1 reflection coefficient. According to Fig. 10, the model-free dither-based method achieves almost the same performance as the model-based alternating optimization when the time index is greater than 500. However, in practice, we have to weigh the cost of time resources against the benefits. As the dither-based method needs to sample $\hat{\mathbf{H}}$, one iteration means the cost of a unit of time resource in wireless communications. Thus, in order to balance the time allocation between pilot transmission

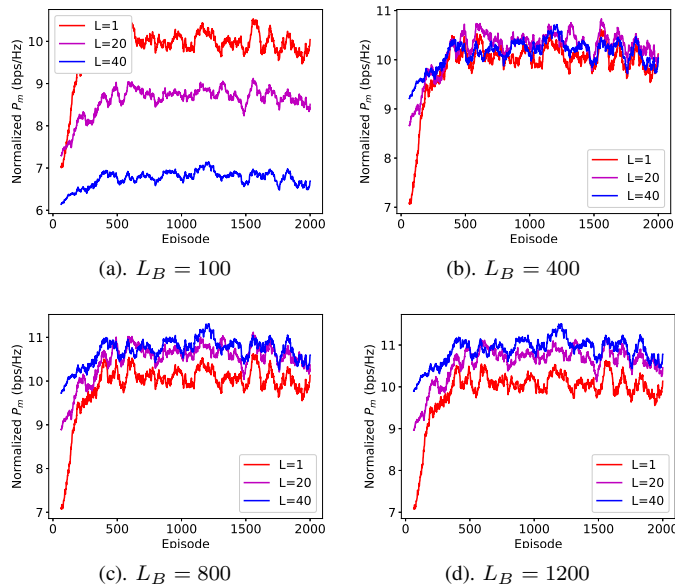


Fig. 11. Moving average of the normalized P_m for the integrated DRL and ESC in different channel dynamics

and data transmission, the time resources dedicated to dither-based method (i.e., the time for pilot transmission) should be deliberately selected according to channel dynamics (i.e., the length of a channel block).

D. Performance Study of the Integrated DRL and ESC

In Fig. 11, we study the performance of the integrated DRL and ESC method (with the action set \mathcal{A}_5) in different channel dynamics. The normalized P_m is the obtained by multiplying P_m by the coefficient $\frac{L_B-L}{L_B}$, where L_B is the channel block length and L is the training length. Take $L_B = 100$ and $L = 1$ as an example, the first $L = 1$ time interval is used for pilot transmission and the rest $L_B - L = 99$ time intervals are used for data transmission. It is also noteworthy that when $L = 1$ the scheme is DRL only, and when $L \geq 2$ the scheme is the integrated DRL and ESC. In addition, we adopt the action set \mathcal{A}_5 . From the figure, we can see that when the channel block length $L_B = 100$, DRL outperforms the integrated DRL and ESC with $L = 20, 40$. However, when the channel block length increases, the integrated DRL and ESC gradually becomes superior, which verifies its effectiveness in the slow fading channel. Therefore, the parameter L of the integrated DRL and ESC can be set adaptively to accommodate different channel dynamics.

VI. CONCLUSION

In this paper, we have proposed a model-free control of IRS that is independent of sub-channel CSI. We firstly model the control of IRS as an MDP and apply DRL to perform real-time coarse phase control of IRS. Then, we apply ESC as the fine phase control of IRS. Finally, by updating the frame structure, we integrate DRL and ESC in the model-free control of IRS to improve its adaptivity to different channel dynamics. Numerical results show the superiority of our proposed scheme in model-free IRS control without sub-channel CSI.

REFERENCES

- [1] N. Engheta and R. W. Ziolkowski, *Metamaterials: Physics and Engineering Explorations*. John Wiley & Sons, 2006.
- [2] T. J. Cui, M. Q. Qi, X. Wan, J. Zhao, and Q. Cheng, "Coding metamaterials, digital metamaterials and programmable metamaterials," *Light Sci. Appl.*, vol. 3, no. 10, pp. e218–e218, Oct. 2014.
- [3] D. Schurig, J. J. Mock, B. Justice, S. A. Cummer, J. B. Pendry, A. F. Starr, and D. R. Smith, "Metamaterial electromagnetic cloak at microwave frequencies," *Science*, vol. 314, no. 5801, pp. 977–980, Nov. 2006.
- [4] L. Liang, M. Qi, J. Yang, X. Shen, J. Zhai, W. Xu, B. Jin, W. Liu, Y. Feng, C. Zhang *et al.*, "Anomalous terahertz reflection and scattering by flexible and conformal coding metamaterials," *Adv. Opt. Mater.*, vol. 3, no. 10, pp. 1374–1380, Oct. 2015.
- [5] N. Yu, P. Genevet, M. A. Kats, F. Aieta, J.-P. Tetienne, F. Capasso, and Z. Gaburro, "Light propagation with phase discontinuities: generalized laws of reflection and refraction," *Science*, vol. 334, no. 6054, pp. 333–337, Oct. 2011.
- [6] C. Zhang, W. Chen, Q. Chen, and C. He, "Distributed intelligent reflecting surfaces-aided device-to-device communications system," *J. Comm. Inform. Networks*, vol. 6, no. 3, pp. 197–207, 2021.
- [7] M. Di Renzo *et al.*, "Smart radio environments empowered by reconfigurable AI meta-surfaces: An idea whose time has come," *EURASIP J. Wireless Commun. Netw.*, vol. 2019, no. 1, pp. 1–20, 2019.
- [8] Q. Wu, S. Zhang, B. Zheng, C. You, and R. Zhang, "Intelligent reflecting surface aided wireless communications: A tutorial," *IEEE Trans. Commun.*, vol. 69, no. 5, pp. 3313–3351, May 2021.
- [9] Z. Wang, L. Liu, and S. Cui, "Channel estimation for intelligent reflecting surface assisted multiuser communications," in *2020 IEEE Wireless Communications and Networking Conference (WCNC)*, 2020, pp. 1–6.
- [10] P. Wang, J. Fang, W. Zhang, and H. Li, "Joint active and passive beam training for IRS-assisted millimeter wave systems," *arXiv preprint arXiv:2103.05812*, 2021.
- [11] P. Wang, J. Fang, H. Duan, and H. Li, "Compressed channel estimation for intelligent reflecting surface-assisted millimeter wave systems," *IEEE Signal Process. Lett.*, vol. 27, pp. 905–909, May 2020.
- [12] Q. Wu and R. Zhang, "Intelligent reflecting surface enhanced wireless network via joint active and passive beamforming," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5394–5409, Nov. 2019.
- [13] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. Cambridge University Press, 2005.
- [14] C. Qi, P. Dong, W. Ma, H. Zhang, Z. Zhang, and G. Y. Li, "Acquisition of channel state information for mmWave massive MIMO: Traditional and machine learning-based approaches," *Sci. China Inf. Sci.*, vol. 64, no. 8, pp. 1–16, 2021.
- [15] F. Liu, J. Pan, X. Zhou, and G. Y. Li, "Atmospheric ducting effect in wireless communications: Challenges and opportunities," *J. Comm. Inform. Networks*, vol. 6, no. 2, pp. 101–109, 2021.
- [16] W. Wang and W. Zhang, "Joint beam training and positioning for intelligent reflecting surfaces assisted millimeter wave communications," *IEEE Trans. Wireless Commun.*, vol. 20, no. 10, pp. 6282–6297, Oct. 2021.
- [17] H. Yang *et al.*, "A programmable metasurface with dynamic polarization, scattering and focusing control," *Sci. Rep.*, vol. 6, no. 1, pp. 1–11, Oct. 2016.
- [18] W. Tang *et al.*, "Wireless communications with programmable metasurface: Transceiver design and experimental results," *China Commun.*, vol. 16, no. 5, pp. 46–61, May 2019.
- [19] Q. Wu and R. Zhang, "Joint active and passive beamforming optimization for intelligent reflecting surface assisted SWIPT under QoS constraints," *IEEE J. Sel. Areas in Commun.*, vol. 38, no. 8, pp. 1735–1748, Aug. 2020.
- [20] Y. Yang, S. Zhang, and R. Zhang, "IRS-enhanced OFDM: Power allocation and passive array optimization," in *2019 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2019, pp. 1–6.
- [21] X. Yu, D. Xu, D. W. K. Ng, and R. Schober, "IRS-assisted green communication systems: Provable convergence and robust optimization," *IEEE Trans. Commun.*, vol. 69, no. 9, pp. 6313–6329, Sept. 2021.
- [22] K. Zhi, C. Pan, H. Ren, and K. Wang, "Ergodic rate analysis of reconfigurable intelligent surface-aided massive MIMO systems with ZF detectors," *IEEE Commun. Lett.*, vol. 26, no. 2, pp. 264–268, Feb. 2022.
- [23] Y. Han, W. Tang, S. Jin, C.-K. Wen, and X. Ma, "Large intelligent surface-assisted wireless communication exploiting statistical CSI," *IEEE Trans. Veh. Technol.*, vol. 68, no. 8, pp. 8238–8242, Aug. 2019.

- [24] K. Zhi, C. Pan, G. Zhou, H. Ren, M. Elkhachlan, and R. Schober, "Is RIS-aided massive MIMO promising with ZF detectors and imperfect CSI?" *arXiv preprint arXiv:2111.01585*, 2021.
- [25] K. Zhi, C. Pan, H. Ren, K. Wang, M. Elkhachlan, M. Di Renzo, R. Schober, H. V. Poor, J. Wang, and L. Hanzo, "Two-timescale design for reconfigurable intelligent surface-aided massive MIMO systems with imperfect CSI," *arXiv preprint arXiv:2108.07622*, 2021.
- [26] M.-M. Zhao, Q. Wu, M.-J. Zhao, and R. Zhang, "Intelligent reflecting surface enhanced wireless networks: Two-timescale beamforming optimization," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 2–17, Jan. 2021.
- [27] A. Taha, Y. Zhang, F. B. Mismar, and A. Alkhateeb, "Deep reinforcement learning for intelligent reflecting surfaces: Towards standalone operation," in *2020 IEEE 21st International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, 2020, pp. 1–5.
- [28] A. Taha, M. Alrabeiah, and A. Alkhateeb, "Enabling large intelligent surfaces with compressive sensing and deep learning," *IEEE Access*, vol. 9, pp. 44 304–44 321, 2021.
- [29] B. Sheen, J. Yang, X. Feng, and M. M. U. Chowdhury, "A deep learning based modeling of reconfigurable intelligent surface assisted wireless communications for phase shift configuration," *IEEE Open J. Commun. Soc.*, vol. 2, pp. 262–272, 2021.
- [30] Y. Kim, G. Miao, and T. Hwang, "Energy efficient pilot and link adaptation for mobile users in TDD multi-user MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 13, no. 1, pp. 382–393, Jan. 2013.
- [31] W. Zhang, H. Ren, C. Pan, M. Chen, R. C. de Lamare, B. Du, and J. Dai, "Large-scale antenna systems with UL/DL hardware mismatch: Achievable rates analysis and calibration," *IEEE Trans. Commun.*, vol. 63, no. 4, pp. 1216–1229, Apr. 2015.
- [32] A. Goldsmith, *Wireless Communications*. Cambridge University Press, 2005.
- [33] A. K. Samingan, I. Suleiman, A. A. A. Rahman, and Z. M. Yusof, "LTF-based vs. pilot-based MIMO-OFDM channel estimation algorithms: An experimental study in 5.2 GHz wireless channel," in *2009 IEEE 9th Malaysia International Conference on Communications (MICC)*, 2009, pp. 794–800.
- [34] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT press, 2018.
- [35] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement learning: A survey," *J. Mach. Learn. Res.*, vol. 4, pp. 237–285, May 1996.
- [36] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep reinforcement learning: A brief survey," *IEEE Signal Process. Mag.*, vol. 34, no. 6, pp. 26–38, Nov. 2017.
- [37] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
- [38] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.
- [39] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double Q-learning," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, vol. 30, no. 1, 2016.
- [40] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [41] P. Ramachandran, B. Zoph, and Q. V. Le, "Searching for activation functions," *arXiv preprint arXiv:1710.05941*, 2017.
- [42] K. B. Ariyur and M. Krstić, *Real-Time Optimization by Extremum-Seeking Control*. John Wiley & Sons, 2003.
- [43] Nešić *et al.*, "A framework for extremum seeking control of systems with parameter uncertainties," *IEEE Trans. Autom. Control*, vol. 58, no. 2, pp. 435–448, Feb. 2013.
- [44] Y. Tan, W. Moase, C. Manzie, D. Nešić, and I. Mareels, "Extremum seeking from 1922 to 2010," in *Proceedings of the 29th Chinese Control Conference*, 2010, pp. 14–26.
- [45] H.-H. Wang, M. Krstić, and G. Bastin, "Optimizing bioreactors by extremum seeking," *Int. J. Adapt. Control Signal Process.*, vol. 13, no. 8, pp. 651–669, 1999.
- [46] D. Nešić, "Extremum seeking control: Convergence analysis," *Eur. J. Control*, vol. 15, no. 3-4, pp. 331–347, 2009.
- [47] D. Carnevale *et al.*, "Maximizing radiofrequency heating on FTU via extremum seeking: parameter selection and tuning," in *From Physics to Control Through an Emergent View*. World Scientific, 2010, pp. 321–326.
- [48] K. T. Atta, A. Johansson, and T. Gustafsson, "Extremum seeking control based on phasor estimation," *Syst. Control Lett.*, vol. 85, pp. 37–45, Nov. 2015.
- [49] W. Wang and W. Zhang, "Jittering effects analysis and beam training design for UAV millimeter wave communications," *IEEE Trans. Wireless Commun.*, pp. 1–1, 2021.