# Hierarchical Deep Reinforcement Learning for Age-of-Information Minimization in IRS-aided and Wireless-powered Wireless Networks

Shimin Gong, Leiyang Cui, Bo Gu, Bin Lyu, Dinh Thai Hoang, and Dusit Niyato

## Abstract

In this paper, we focus on a wireless-powered sensor network coordinated by a multi-antenna access point (AP). Each node can generate sensing information and report the latest information to the AP using the energy harvested from the AP's signal beamforming. blueWe aim to minimize the average age-of-information (AoI) by adapting the nodes' transmission scheduling and the transmission control strategies jointly. To reduce the transmission delay, an intelligent reflecting surface (IRS) is used to enhance the channel conditions by controlling the AP's beamforming bluevector and the IRS's phase shifting bluematrix. blueConsidering dynamic data arrivals at different sensing nodes, we propose a hierarchical deep reinforcement learning (DRL) framework to bluefor AoI minimization in two steps. The users' transmission scheduling is firstly determined by the outer-loop DRL approach, e.g. the DQN or PPO algorithm, and then the inner-loop optimization is used to adapt either the uplink information transmission or downlink energy transfer to all nodes. A simple and efficient approximation is also proposed to reduce the blueinner-loop rum time overhead. Numerical results verify that the hierarchical learning framework outperforms typical baselines in terms of the average AoI and proportional fairness among different nodes.

## Index Terms

Shimin Gong, Leiyang Cui, and Bo Gu are with the School of Intelligent Systems Engineering, Sun Yat-sen University, Shenzhen 518055, China (email: gongshm5@mail.sysu.edu.cn, cuily6@mail2.sysu.edu.cn, gubo@mail.sysu.edu.cn). Bin Lyu is with Key Laboratory of Ministry of Education in Broadband Wireless Communication and Sensor Network Technology, Nanjing University of Posts and Telecommunications, China (e-mail: blyu@njupt.edu.cn). Dinh Thai Hoang is with the School of Electrical and Data Engineering, University of Technology Sydney, Australia (email: hoang.dinh@uts.edu.au). Dusit Niyato is with School of Computer Science and Engineering, Nanyang Technological University, Singapore (email: dniyato@ntu.edu.sg).

AoI minimization, wireless power transfer, IRS-aided wireless network, deep reinforcement learning.

# I. INTRODUCTION

With the development of the future Internet of Things (IoT), a large portion of the emerging applications (e.g., autonomous driving, interactive gaming, and virtual reality) depends on timely transmissions and processing of the IoT devices' sensing data in the physical world, e.g., [1] and [2]. Maintaining information freshness in time-sensitive applications requires a new network performance metric, i.e., the age-of-information (AoI), which is defined as the elapsed time since the generation of the latest status update successfully received by the receiver [3]. In wireless networks, the overall time delay of each node is mainly caused by the waiting time or scheduling delay before information transmission and the in-the-air transmission delay. The waiting delay is usually determined by the multi-user scheduling strategy, while the transmission delay depends on the network capacity and channel conditions, e.g., mutual interference or channel fading effect. To minimize the transmission delay, we need to explore preferable channel opportunities and adapt the transmission control parameters accordingly, e.g., power control, channel allocation, and beamforming strategies. More recently, the intelligent reflecting surface (IRS) has been used to reduce the transmission delay by shaping the wireless channels in favor of information transmission via passive beamforming optimization [4]. The AoI-aware scheduling and transmission control in wireless networks were previously studied by the queueing theory [5]. The closed-form analysis of the AoI performance is tedious and usually difficult, typically relying on specific probability distributions of the sensor nodes' data arrivals, the data transmission/service time, and the rate of information requests at user devices. For more complex wireless networks, e.g., with users' mobility and limited energy resources, it is still very challenging to optimize multi-user scheduling policy to maximize the overall AoI performance.

In energy-constrained wireless networks, the AoI minimization depends on the optimal control of each user's energy supply and demand, especially in wireless powered communication networks. When a sensor node is scheduled more often to update its sensing information, more wireless energy transfer is required for the node to achieve self-sustainability. This may reduce the transmission opportunities and the energy transfer to other nodes, leading to the AoI-energy tradeoff study in energy-constrained wireless networks [6]. The joint optimization of multi-

user scheduling and energy management is usually formulated as a dynamic program due to the spatial-temporal couplings among wireless sensor nodes. With limited energy supply, the users' scheduling becomes more challenging to balance the tradeoff between AoI and energy consumption. The AoI minimization problem is further complicated by the unknown channel conditions and the sensor nodes' dynamic data arrivals. Without complete system information, the scheduling strategy has to be adapted according to the users' historical AoI information. Instead of the optimization approaches, the AoI minimization problems can be flexibly solved by the model-free deep reinforcement learning (DRL) approaches, e.g., [6]–[9]. The DRL approaches can reformulate the AoI minimization problem into Markov decision process (MDP). The action includes the scheduling and transmission control strategies. With a large-size IRS, the action and state spaces become high-dimensional and lead to unreliable and slow learning performance. This motivates us to design a more efficient learning approach for the AoI minimization problem.

Specifically, in this paper, we aim to minimize the average AoI in a wireless-powered network, consisting of a multi-antenna access point (AP), a wall-mounted IRS, and multiple single-antenna sensor nodes, sampling the status information from different physical processes. The IRS is used to enhance the channel conditions and reduce the transmission delay. It can assist either the downlink wireless power transfer from the AP to the sensing nodes or uplink information transmissions from nodes to the AP [4]. The joint scheduling and beamforming optimization normally leads to a high-dimensional mix-integer problem that is difficult to solve practically. Different from the conventional DRL solutions, we devise a two-step hierarchical learning framework to improve the overall learning efficiency. The basic idea is to adapt the scheduling strategy by the outer-loop DRL algorithm, e.g., the value-based deep Q-network (DQN) or the policy-based proximal policy optimization (PPO) algorithm [10], and optimize the beamforming strategy by the inner-loop optimization module. Given the outer-loop decision, the inner-loop procedures either optimize the AP's downlink energy transfer or optimize individual's uplink information transmission. Specifically, the contributions of this paper are summarized as follows:

- *The IRS-assisted scheduling and beamforming for AoI minimization:* We aim to reduce both the packet waiting and transmission delays for updating sensing information in a wireless-powered and IRS-assisted wireless network. The IRS's passive beamforming not only enhances the wireless power transfer to sensor nods, but also assists their uplink information transmissions to the AP. We formulate the AoI minimization problem by jointly

adapting the user scheduling and beamforming strategies.

- *The hierarchical DRL approach for AoI minimization:* A hierarchical DRL framework is proposed to solve the AoI minimization in two steps. The model-free DRL in the outer loop adapts the combinatorial scheduling decision according to each user's energy status and the AoI performance. Given the outer-loop scheduling, we optimize the joint beamforming strategies for either the downlink energy transfer and the uplink information transmission.

- *Policy-based PPO algorithm for outer-loop scheduling:* Our simulation results verify that the hierarchical learning framework significantly reduces the average AoI compared with typical baseline strategies. Besides, we compare both the traditional DQN and the PPO methods for the outer-loop scheduling optimization. The PPO-based hierarchical learning can improve convergence and achieve a lower AoI value compared to the DQN-based method.

Some preliminary results of this work have been presented in [11]. In this extension, we include detailed analysis about the PPO algorithm and compare it with the DQN method. We also propose a simple approximation for the inner-loop optimization to reduce the time overhead while achieving comparable AoI performance at convergence. The remainder of this paper is organized as follows. We discuss related works in Section II and present our system model in Section III. We present the hierarchical learning framework in Section IV. The inner-loop optimization and outer-loop learning procedures are detailed in Sections V and VI, respectively. Finally, we present the numerical results in Section VII and conclude the paper in Section VIII.

## II. RELATED WORKS

### A. DRL Approaches for AoI Minimization

DRL has been introduced recently as an effective solution for AoI minimization by adapting the scheduling and beamforming strategies according to time-varying traffic demands and channel conditions. The authors in [7] designed the freshness-aware scheduling solution by using the DQN method. The DQN agent continuously updates its scheduling strategy to maximize the freshness of information in the long term. The authors in [8] focused on a multi-user status update system, where a single sensor node monitors a physical process and schedules its information updates to multiple users with time-varying channel conditions. Based on the user's instantaneous ACK/NACK feedback, the DRL agent at the information source can decide on when and to which user to transmit the next information update, without any priori information on the random

channel conditions. The authors in [9] focused on a multi-access system where the base-station (BS) aims to maximize its information collected from all wireless users. Given a strict time deadline to each wireless user, the PPO algorithm was used to adapt the scheduling policy by learning the users' traffic patterns from the past experiences. The authors in [12] considered a different multi-user scheduling scheme that allows a group of sensor nodes to transmit their information simultaneously. The scheduling decision is adapted to minimize the AoI by using the double DQN (DDQN) method [13], an extension of the DQN method by using two sets of deep neural networks (DNNs) to approximate the Q-value. The AoI-energy tradeoff study in [14] revealed that the sensor nodes' energy consumption can be reduced without a significant increase in the AoI, by using DQN to adapt the content update in the caching node.

## B. RF-powered Scheduling for AoI Minimization

The wireless power transfer and energy harvesting are promising techniques to sustain the massive number of low-cost sensor nodes. However, energy harvesting is usually unreliable depending on the channel conditions. The dynamics and scarcity in energy harvesting make it more challenging for the sensor nodes' energy management and AoI minimization. The authors in [15] focused on AoI minimization in a cognitive radio network with dynamic supplies of the energy and spectrum resources. The optimal scheduling policy is derived by a dynamic programming approach, revealing a threshold structure depending on the sensor nodes' AoI states. The authors in [16] revealed that the optimal policy allows each sensor to send a status update only if the AoI value is higher than some threshold that depends on its energy level. Considering stochastic energy harvesting at sensor devices, the authors in [17] studied the AoI minimization problem with the long-term energy constraints and proposed Lyapunov-based dynamic optimization to derive an approximate solution. Generally the dynamic optimization approaches are not only computational demanding, but also relying on the availability of system information. Without knowing the dynamic energy arrivals, the authors in [18] reformulated the AoI minimization problem into MDP. The online Q-learning method was proposed to adaptively schedule the wireless devices' information update. The authors in [19] focused on the RF-powered wireless network, where the wireless devices can harvest RF power from a dedicated BS and then transmit their update packets to the BS. To minimize the long-term AoI, the DQN algorithm was used to adapt the scheduling between the downlink energy transfer

and the uplink information transmission. Both the DQN and dueling DDQN methods were used in [20] to adapt the sensing and information update policy for AoI minimization in a spectrum sharing cognitive radio system. Considering energy harvesting ad hoc networks, the authors in [21] solved the AoI minimization problem by using the advantage actor-critic (A2C) algorithm to adapt the scheduling and power allocation policy, which shows faster runtime and comparable AoI performance to the optimum. The above-mentioned works typically solve the AoI minimization by using the conventional model-free DRL methods. These methods become inflexible and unreliable due to slow convergence with the increasing state and action spaces.

*C. IRS-Assisted AoI Minimization*

The IRS's reconfigurability can be used to enhance the channel quality/capacity or reduce the transmission delay by tuning the phase shifts of the reflecting elements [4]. Only a few existing works have discussed the IRS's application for AoI minimization in wireless networks. The authors in [22] focused on a mobile edge computing (MEC) system and proposed using the IRS to minimize the workload processing delay by optimizing the IRS's passive beamforming strategy. The authors in [23] set delay constraints to the wireless users' uplink information transmissions and revealed that the IRS's passive beamforming can help reduce the wireless users' transmit power. The authors in [24] employed the IRS to enhance the AoI performance by jointly optimizing the uses' scheduling and the IRS's passive beamforming strategies. The combinatorial scheduling decision is adapted by the model-free DRL algorithm, while the passive beamforming optimization relies on the solution to the conventional semi-definite relaxation (SDR) problem. The authors in [25] employed the UAV-carrying IRS to relay information from the ground users to the BS. The AoI minimization requires the optimization of the UAV's altitude, the ground users' transmission scheduling, and the IRS's passive beamforming strategies. Comparing to [24] and [25], our work in this paper exploits the performance gains in both the uplink and downlink of the IRS-assisted system. The IRS's passive beamforming not only assists uplink information transmission but also enhances or balances the AP's downlink energy transfer to the users.

## III. SYSTEM MODEL

We consider an IRS-assisted wireless sensor network deployed in smart cities to assist information sensing and decision making, similar to that in [11]. The system consists of a multi-antenna

**TABLE I:** A list of Notations

| Notation | Description | Notation | Description |
|---|---|---|---|
| $M$ | Number of AP's antennas | $N$ | Size of the IRS |
| $K$ | Number of IoT devices | $\mathcal{T}$ | The set of time slots |
| $\psi_0(t) \in \{0,1\}$ | The AP's mode selection | $\psi_k(t) \in \{0,1\}$ | The AP's uplink scheduling |
| $\mathbf{G}(t)$ | The AP-IRS channel matrix | $\mathbf{h}_k^r(t)$ | The IRS-User channel vector |
| $\mathbf{\Theta}_d(t), \mathbf{\Theta}_u(t)$ | The IRS's beamforming strategies | $\mathbf{w}_d(t), \mathbf{w}_u(t)$ | The AP's beamforming vectors |
| $\eta$ | Energy conversion efficiency | $E_k^h(t)$ | Energy harvested by the user-$k$ |
| $E_k^c(t)$ | The user-$k$'s energy consumption | $E_k(t)$ | The user-$k$'s energy state |
| $E_{\max}$ | Maximum energy capacity | $\gamma_k(t)$ | The received SNR at the AP |
| $r_k(t)$ | The user-$k$'s uplink throughput | $\tau_k$ | The user-$k$'s uplink transmission time |
| $d_k$ | The user-$k$'s data size | $A_k(t)$ | The user-$k$'s AoI value |

AP with $M$ antennas, an IRS with $N$ reflection elements, and $K$ single-antenna IoT devices, denoted by the set $\mathcal{K} \triangleq \{1, 2, \ldots, K\}$. The system model can be straightforwardly extended to the cases with multiple AP or multiple IRSs. The sensing information is typically a small amount of data, which should be timely updated to the AP for real-time status monitoring. We assume that all sensor nodes are low-power devices and wireless powered by harvesting RF energy from the AP's beamforming signals. The wireless powered communications technology has been verified and evaluated in [26], showing that the LoRa-based sensor nodes typically have 0.5-1.5 mJ energy consumption, and require 2-5 seconds of energy harvesting time 3.2 meters away to sustain periodical information sensing and data transmissions up to 200 bytes. The IRS can be deployed on the exterior walls of buildings to enhance the channel conditions between the sensor nodes and the AP. We aim to collect all sensor nodes' data in a timely fashion by scheduling their uplink data transmissions, based on their channel conditions, traffic demands, and energy status. A list of notations is provided in Table I.

### A. Mode Selection and Scheduling

We consider a time-slotted frame structure to avoid contention between different nodes. Each data frame is equally divided into $T$ time slots allocated to different sensor nodes. Let $\mathcal{T} \triangleq \{1, 2, \ldots, T\}$ denote the set of all time slots. In each time slot, we need to firstly decide the AP's operation mode, i.e., the time slot is used for either the downlink energy transfer or the uplink data transmission. We use $\psi_0(t) \in \{0, 1\}$ to denote the AP's mode selection in each time

slot, i.e., $\psi_0(t) = 1$ indicates the downlink energy beamforming and $\psi_0(t) = 0$ represents the uplink information transmission. We further use $\psi_k(t) \in \{0, 1\}$ to denote the uplink scheduling policy, i.e., $\psi_k(t) = 1$ represents that the $k$-th sensor node is allowed to access the channel for uploading its sensing information to the AP. We require that at most one sensor node can access the channel in each time slot, which implies the following scheduling constraint:

$$\psi_0(t) + \sum_{k \in \mathcal{K}} \psi_k(t) \leq 1, \quad \forall t \in \mathcal{T}. \tag{1}$$

We denote $\mathbf{\Psi}(t) = [\psi_0(t), \psi_1(t), \ldots, \psi_K(t)]$ as the AP's scheduling policy, which depends on the sensor nodes' traffic demands, channel conditions, and energy status.

Let $\mathcal{N} \triangleq \{1, 2, \ldots, N\}$ denote the set of the IRS's reflecting elements and $\theta_n(t) \in (0, 2\pi]$ denote the phase shift of the $n$-th reflecting element in the $t$-th time slot. We define the IRS's phase shifting vector in the $t$-th time slot as $\boldsymbol{\theta}(t) = [e^{j\theta_n(t)}]_{n \in \mathcal{N}}$. Note that the IRS can set different beamforming vectors, denoted as $\boldsymbol{\theta}_d(t) \triangleq [e^{j\theta_{d,n}(t)}]_{n \in \mathcal{N}}$ and $\boldsymbol{\theta}_u(t) \triangleq [e^{j\theta_{u,1}(t)}]_{n \in \mathcal{N}}$, for the downlink and uplink phases, respectively. The channel matrix from the AP to the IRS in $t$-th time slot is given by $\mathbf{G}(t) \in \mathbb{C}^{M \times N}$. The channel vectors from the IRS and the AP to the $k$-th sensor node are denoted by $\mathbf{h}_k^r(t) \in \mathbb{C}^{N \times 1}$ and $\mathbf{h}_k^d(t) \in \mathbb{C}^{M \times 1}$, respectively. The AP can estimate the channel information by a training period at the beginning of each time slot.

### B. Downlink Energy Transfer

When $\psi_0(t) = 1$, the IRS-assisted downlink energy transfer ensures the sustainable operation of the system. Given the IRS's passive beamforming strategy $\boldsymbol{\theta}_d(t)$, the equivalent downlink channel vector $\mathbf{f}_{d,k}(t)$ from the AP to the $k$-th sensor node can be expressed as follows:

$$\mathbf{f}_{d,k}(t) = \mathbf{h}_k^d(t) + \mathbf{G}_k(t)\boldsymbol{\Theta}_d(t)\mathbf{h}_k^r(t), \tag{2}$$

where we define $\boldsymbol{\Theta}_d(t) \triangleq \text{diag}(\boldsymbol{\theta}_d(t))$ as a diagonal matrix with the diagonal element $\boldsymbol{\theta}_d(t)$. The phase shifting matrix $\boldsymbol{\Theta}_d(t)$ represents the IRS's passive beamforming strategy in the downlink energy transfer. Let $\mathbf{w}_d(t) \in \mathbb{C}^{M \times 1}$ denote the AP's transmit beamforming vector in the downlink energy transfer phase. Given the AP's transmit power $p_s$, the AP's beamforming signal is given by $\mathbf{x}(t) = \sqrt{p_s}\mathbf{w}_d(t)s_0(t)$, where $s_0(t) \in \mathbb{C}$ denotes a random complex symbol with the unit power. Then, the received signal at the $k$-th sensor node is given as $y_k(t) = \mathbf{f}_{d,k}^H(t)\mathbf{x}(t) + n_k(t)$, where $(\cdot)^H$ denotes conjugate transpose and $n_k(t)$ is the normalized Gaussian noise with zero

mean and unit power. Considering a linear energy harvesting (EH) model [27], the received signal $y_k(t)$ can be converted to energy as follows:

$$E_k^h(t) = \eta\mathbb{E}[|\mathbf{f}_{d,k}^H(t)\mathbf{x}(t)|^2] = \eta p_s|\mathbf{f}_{d,k}^H(t)\mathbf{w}_d(t)|^2, \tag{3}$$

where $\eta$ denotes the energy conversion efficiency. The energy harvested from the noise signal is assumed to be negligible. It is clear that the AP can control the energy transfer to different sensor nodes by optimizing the downlink beamforming vector $\mathbf{w}_d(t)$.

In particular, the AP can steer the beam direction toward the sensor nodes with insufficient energy supply. Besides, the IRS's passive beamforming strategy $\boldsymbol{\Theta}_d(t)$ controls the downlink channel conditions $\mathbf{f}_{d,k}(t)$ to individual receivers. Due to the broadcast nature of wireless signals, the AP's energy transfer to different sensor nodes depends on the joint beamforming strategies $(\mathbf{w}_d(t), \boldsymbol{\Theta}_d(t))$ in different time slots. A more practical non-linear EH model can be also applied to our system. In this case, the harvested power firstly increases with the received signal power and then becomes saturated as the received signal power continues to increase, e.g., [28]. This can be approximated by a piecewise linear EH model: $E_k^h(t) = \min\left\{\eta p_s|\mathbf{f}_{d,k}^H(t)\mathbf{w}_d(t)|^2, \ p_{\text{sat}}\right\}$, where $p_{\text{sat}}$ denotes the saturation power. Our algorithm in the following part adopts the linear EH model in (3) and can be easily applied to the piecewise linear model with minor modifications.

## C. Sensing Information Updates

We assume that the uplink channels are the same as the downlink channels in each time slot due to channel reciprocity, similar to that in [11] and [19]. Let $p_k(t)$ denote the transmit power of the $k$-th sensor node when it is scheduled in the $t$-th time slot, i.e., $\psi_k(t) = 1$. The signal received at the AP is given by $\mathbf{y}_k = \sqrt{p_k(t)}\mathbf{f}_{u,k}(t)s_k + \mathbf{n}_k(t)$, where $\mathbf{f}_{u,k}(t)$ denotes the uplink channel from the $k$-th sensor node to the AP and $s_k(t)$ denotes its information symbol. Similar to (2), the IRS-assisted uplink channel is given by $\mathbf{f}_{u,k}(t) = \mathbf{h}_k^d(t) + \mathbf{G}_k(t)\boldsymbol{\Theta}_u(t)\mathbf{h}_k^r(t)$, where $\boldsymbol{\Theta}_u(t) \triangleq \operatorname{diag}(\boldsymbol{\theta}_u(t))$ denotes the IRS's uplink passive beamforming strategy. Without loss of generality, the noise signal $\mathbf{n}_k(t)$ received by the AP can be normalized to the unit power. Thus, the received SNR can be characterized as $\gamma_k(t) = p_k(t)|\mathbf{f}_{u,k}^H(t)\mathbf{w}_u(t)|^2$, where $\mathbf{w}_u(t)$ represents the AP's receive beamforming vector. By using the time division protocol, the sensor nodes' uplink transmissions can avoid mutual interference. The AP can simply align its receive beamforming vector $\mathbf{w}_u(t)$ to the uplink channel $\mathbf{f}_{u,k}(t)$. As such, we can denote the received SNR as $\gamma_k(t) =$

$p_k(t)||\mathbf{f}_{u,k}(t)||^2$ and characterize the uplink throughput as $r_k(t) = \tau_k \log\big(1 + p_k(t)||\mathbf{f}_{u,k}(t)||^2\big)$, where $\tau_k \in [0,1]$ denotes the uplink transmission time. Given the data size $d_k$, we require $r_k(t) \geq d_k$ to ensure the successful transmission of the sensing information.

In each time slot, the sensor node's energy consumption is given by $E_k^c(t) = \tau_k(t)(p_k(t) + p_c)$, where $p_c$ denotes a constant circuit power to maintain the node's activity. The power consumption $\tau_k(t)p_k(t)$ in uplink data transmission is linearly proportional to the transmit power $p_k(t)$ and the transmission time $\tau_k$. The transmit power $p_k(t)$ can vary with the channel conditions to ensure the rate constraint $r_k(t) \geq d_k$. Let $E_{\max}$ denote the sensor nodes' maximum battery capacity and $E_k(t)$ be the energy state in the $t$-th time slot. Then, we have the following energy dynamics:

$$E_k(t+1) = \min\left\{\left(E_k(t) - \psi_k(t)E_k^c(t)\right)^+ + \psi_0(t)E_k^h, E_{\max}\right\}. \tag{4}$$

Here we denote $(x)^+ \triangleq \max\{x, 0\}$ for simplicity.

## IV. HIERARCHICAL LEARNING FOR AoI MINIMIZATION

In this paper, the physical sensing process is beyond our control and we only focus on the transmission scheduling and beamforming optimization over the wireless network. The sensing information can be randomly generated by the sensor nodes, depending on the energy status and the physical process under monitoring. Once new sensing data arrives, each sensor node will discard existing data in the cache and always cache the latest sensing data. From the perspective of the receiver, the sensing data from each sensor node is considered as the new information and used to replace the obsolete information at the receiver. When the node-$k$ is scheduled for uplink data transmission, the cached information will be uploaded to the AP and then the AP will replace the sensing information by the latest copy.

For each sensor node $k \in \mathcal{K}$, the caching delay depends on the AP's scheduling policy $\mathbf{\Psi}(t)$. The transmission delay can be minimized by optimizing the sensor node's transmit control strategy $(p_k(t), \tau_k(t))$ and the joint beamforming strategies $(\mathbf{w}_u(t), \mathbf{\Theta}_u(t))$ in the uplink transmissions. Let $A_k(t)$ denote the AoI value of the $k$-th sensor node, which is used to characterize the information freshness at the AP. When the node-$k$ is scheduled to update its information in the $t$-th time slot, the AP can replace the obsolete information by the new sensing information and thus update the AoI in the next time slot as $A_k(t+1) = 1$. Here we assume that the node-$k$ can successfully finish its data transmission at the end of each time slot. Otherwise, when the

node-$k$ is not scheduled, its AoI will be further increased by one, i.e., $A_k(t + 1) = A_k(t) + 1$. Therefore, the AoI of each sensor node $k \in \mathcal{K}$ will be updated by the following rules:

$$A_k(t+1) = \begin{cases} 1, & \text{if } o_k(t) = 1, \psi_k(t) = 1, r_k(t) \geq d_k, E_k(t) \geq E_k^c(t), \\ A_k(t) + 1, & \text{otherwise.} \end{cases} \quad (5)$$

Here $o_k(t) \in \{0, 1\}$ indicates the status of the caching space. When the cache is non-empty with $o_k(t) = 1$ and the node-$k$ is currently scheduled with $\psi_k(t) = 1$, the AP can update the sensing information if the uplink data transmission is successful. Given the size $d_k$ of sensing data, the scheduled node-$k$ will have a successful data transmission if it has sufficient energy, i.e., $E_k(t) \geq E_k^c(t)$, to fulfill its traffic demand, i.e., $r_k(t) \geq d_k$, where the uplink data rate $r_k(t)$ depends on the control parameters $(p_k(t), \tau_k)$ and the joint beamforming strategies $(\mathbf{w}_u(t), \boldsymbol{\Theta}_u(t))$.
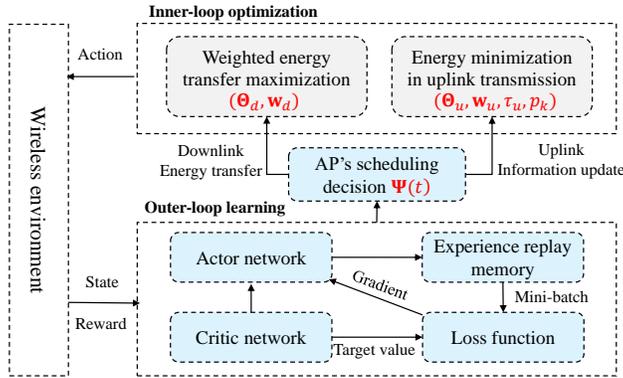
We aim to minimize the AoI by optimizing the scheduling policy $\boldsymbol{\Psi} \triangleq \{\boldsymbol{\Psi}(t)\}_{t \in \mathcal{T}}$ and the joint beamforming strategies $(\mathbf{w}, \boldsymbol{\Theta}) \triangleq (\mathbf{w}_m(t), \boldsymbol{\Theta}_m(t))_{m \in \{d,u\}, t \in \mathcal{T}}$ in both the downlink and uplink phases. Considering different priorities of the sensing information, we assign different weights to the AoI values and define the time-averaged weighted AoI as follows:

$$\bar{A}(\boldsymbol{\Psi}, \mathbf{w}, \boldsymbol{\Theta}) = \lim_{T \to \infty} \frac{1}{TK} \mathbb{E}\left[\sum_{t \in \mathcal{T}} \sum_{k \in \mathcal{K}} \lambda_k A_k(t)\right], \quad (6)$$

where the constant $\lambda_k$ indicates the delay sensitivity of different sensing information. A larger weight should be given to more critical sensing information, e.g., the safety monitoring in autonomous driving. Till this point, we can formulate the AoI minimization problem as follows:

$$\min_{\boldsymbol{\Psi}, \mathbf{w}, \boldsymbol{\Theta}} \quad \bar{A}(\boldsymbol{\Psi}, \mathbf{w}, \boldsymbol{\Theta}), \quad \text{s.t. (1)} - \text{(5)}. \quad (7)$$

Given the mode selection $\psi_0(t)$, the optimization of $(\mathbf{w}, \boldsymbol{\Theta})$ corresponds to either the downlink energy transfer or the uplink information transmission. The downlink energy transfer determines the power budgets of different sensor nodes, which should be jointly optimized with the users' scheduling policy $\{\psi_k(t)\}_{k \in \mathcal{K}}$ to improve the overall AoI performance. The problem (7) is firstly challenged by the stochasticity and high-dimensionality. The energy dynamics in (4) and the time-averaged AoI objective in (7) imply that the solutions $(\boldsymbol{\Psi}(t), \mathbf{w}(t), \boldsymbol{\Theta}(t))$ in each time slot are temporally correlated. A dynamic programming approach to solve (7) can be practically intractable due to the curse of dimensionality. The joint scheduling and beamforming optimization also lead to a high-dimensional mix-integer problem that is difficult to solve practically. The

**Fig. 1:** The hierarchical DRL framework includes the outer-loop DRL and the inner-loop optimization methods.

second difficulty of problem (7) lies in that the dynamics of the data arrival process at each sensor node can be unknown to the AP, which makes the scheduling optimization more complicated in practice. Without complete information, the AP has to adapt its scheduling policy based on the past observations of the AoI dynamics. The third difficulty lies in the combinatorial nature of the AP's scheduling policy $\mathbf{\Psi}$. Given the scheduling policy $\mathbf{\Psi}$, the challenges still exist as the joint beamforming strategies $(\mathbf{w}, \mathbf{\Theta})$ are not only coupled with each other in a non-convex form, but also introduce the competition for energy resource among different sensor nodes.

To overcome these difficulties, we devise a hierarchical learning structure for problem (7) that decomposes the optimization of $(\mathbf{\Psi}, \mathbf{w}, \mathbf{\Theta})$ into two parts. The overall algorithm framework is shown in Fig. 1, which mainly includes the outer-loop learning module for scheduling and the inner-loop optimization module for beamforming optimization. In fact, we may also apply an inner-loop learning method to adapt the beamforming strategy. However, this may still require huge action and state spaces and thus lead to slow learning performance. Instead of an inner-loop learning method, the AP can estimate the beamforming strategy $(\mathbf{w}(t), \mathbf{\Theta}(t))$ efficiently by using the optimization method, based on the AP's observation of the users' channel conditions. This motivates us to devise a hybrid solution structure that exploits the outer-loop learning and the inner-loop optimization modules. Specifically, due to the combinatorial nature of the scheduling policy $\mathbf{\Psi}(t)$, we employ the model-free DRL approach to adapt $\mathbf{\Psi}(t)$ in the outer-loop learning procedure. In each iteration, the DRL agent firstly determines $\mathbf{\Psi}(t)$ based on the past observations of the nodes' AoI dynamics. Then, the inner-loop joint optimization of $(\mathbf{w}(t), \mathbf{\Theta}(t))$ becomes much easier by using the alternating optimization (AO) and semi-definite relaxation (SDR) methods. According to the outer-loop mode selection $\psi_0(t)$, the inner-loop optimization

aims to either maximize the downlink energy transfer to all sensor nodes or fulfill the uplink information transmission of the scheduled sensor node. After the inner-loop optimization, the AP can execute the joint action $(\boldsymbol{\Psi}(t), \mathbf{w}(t), \boldsymbol{\Theta}(t))$ in the $t$-th time slot and then update the AoI state of each sensor node. The evaluation of the time-averaged AoI performance will drive the outer-loop DRL agent to adapt the scheduling decision $\boldsymbol{\Psi}(t+1)$ and the beamforming strategies $(\mathbf{w}(t+1), \boldsymbol{\Theta}(t+1))$ in the next time slot. By such a decomposition, the inner-loop optimization becomes computation-efficient, while the outer-loop learning becomes time-efficient as it only adapts the combinatorial scheduling policy with a smaller action space.

## V. INNER-LOOP OPTIMIZATION PROBLEMS

Given the scheduling decision $\psi_0(t) \in \{0, 1\}$, the AP either beamforms RF signals for downlink energy transfer or receives the sensing information from individual sensor node. In each case, the AP will optimize the joint beamforming strategies $(\mathbf{w}(t), \boldsymbol{\Theta}(t))$. In the sequel, we discuss the inner-loop optimization problems in two cases.

### A. Energy Minimization in Uplink Transmission

In the $t$-th time slot, when the $k$-th sensor node is allowed to update its sensing information with $\psi_k(t) = 1$, all other sensor nodes have to wait for scheduling in the other time slots. The sensor nodes' AoI values will be either increased by 1 or reset to 1 if the transmission is unsuccessful or successful, respectively, as shown in (5). In this case, we can minimize the energy consumption $E_k^c(t) = \tau_k(t)(p_k(t) + p_c)$ of the scheduled sensor node conditioned on the successful transmission of its sensing data, i.e., $r_k(t) \geq d_k$. This will preserve more energy for its future use. Thus, we have the following energy minimization problem:

$$\min_{\tau_k, p_k, \mathbf{w}_u, \boldsymbol{\Theta}_u} \quad \tau_k(p_k + p_c) \tag{8a}$$

$$\text{s.t.} \quad \tau_k \log(1 + p_k|\mathbf{f}_k^H \mathbf{w}_u|^2) \geq d_k, \tag{8b}$$

$$\tau_k \in (0, 1) \text{ and } \theta_{u,n} \in (0, 2\pi], \, n \in \mathcal{N}. \tag{8c}$$

The uplink transmission only considers the node-$k$'s rate constraint in (8b). Hence, the AP's receive beamforming vector $\mathbf{w}_u$ can be aligned with the uplink channel $\mathbf{f}_{u,k} = \mathbf{h}_k^d + \mathbf{G}_k \boldsymbol{\Theta}_u \mathbf{h}_k^r$.

Given $\mathbf{w}_u = \mathbf{f}_{u,k}/||\mathbf{f}_{u,k}||$, the optimal passive beamforming strategy $\mathbf{\Theta}_u$ needs to maximize the IRS-assisted channel gain $|\mathbf{f}_k^H \mathbf{w}_u|^2$ as follows:

$$\max_{\theta_{u,n} \in (0,2\pi]} ||\mathbf{h}_k^d + \mathbf{G}_k \mathbf{\Theta}_u \mathbf{h}_k^r||^2, \tag{9}$$

which can be easily solved by the SDR method similar to that in [24] and [29]. The transmission control parameters $(\tau_k, p_k)$ can be also optimized to minimize the energy consumption in (8a). Let $e_k \triangleq \tau_k p_k$ denote the node's energy consumption in RF communications. Given the optimal $\mathbf{\Theta}_u$ to (9), we can find the optimal $(\tau_k, p_k)$ by the following problem:

$$\min_{e_k, \tau_k \in (0,1)} e_k + p_c \tau_k, \quad \text{s.t.} \quad \tau_k \log\left(1 + \frac{e_k}{\tau_k}||\mathbf{f}_{u,k}^H||^2\right) \geq d_k. \tag{10}$$

Problem (10) is convex in $(\tau_k, e_k)$ and satisfies the Slater's condition, which allows us to find a closed-form solution by using the Lagrangian dual method [30]. After this, we can easily find the optimal transmit power as $p_k = e_k/\tau_k$. If the energy budget holds, i.e., $e_k + p_c \tau_k \leq E_k$, the node-$k$'s data transmission will be successful and thus we can update its AoI as $A_k(t+1) = 1$.

## B. Weighted Energy Transfer Maximization

In downlink energy transfer with $\psi_0(t) = 1$, we aim to supply sufficient energy to all sensor nodes that can sustain their uplink information transmission to minimize the time-averaged AoI performance. However, it is difficult to explicitly quantify how downlink energy transfer affects the AoI performance. Instead, our intuitive design is to transfer more energy to those sensor nodes with the relatively worse AoI conditions. If the node-$k$ has a higher AoI value, we expect to transfer more energy to the node-$k$. This allows the node-$k$ to increase its sampling frequency and report its sensing information with a shorter transmission delay, and thus reducing its AoI value in the following sensing and reporting cycles. By this intuition, in the downlink phase we aim to maximize the AoI-weighted energy transfer to all sensor nodes, formulated as follows:

$$\max_{\mathbf{w}_d, \mathbf{\Theta}_d} \sum_{k \in \mathcal{K}} v_k |\mathbf{f}_{d,k}^H \mathbf{w}_d|^2 \tag{11a}$$

$$\text{s.t.} \quad E_k^c \leq E_k + E_k^h, \quad \forall k \in \mathcal{K}, \tag{11b}$$

$$||\mathbf{w}_d|| \leq 1 \text{ and } \theta_{d,n} \in (0, 2\pi), \quad \forall n \in \mathcal{N}, \tag{11c}$$

where $E_k^c = \tau_k(p_k + p_c)$ denotes the node-$k$'s energy consumption in the uplink transmission. For each node-$k$, we define the weight parameter as $v_k = A_k + \alpha_k E_k^{-1}$, which is increasing in the

AoI value $A_k$ while inversely proportional to the energy capacity $E_k$. Thus, we prefer to transfer more energy to the sensor node with a higher AoI value and a lower energy supply, which is prioritized by a larger weight parameter $v_k$ in (11a). Such a heuristic is expected to reduce the average AoI of all sensor nodes in a long term. The constant $\alpha_k$ characterizes the tradeoff between the sensor node's energy supply and AoI performance. A larger value of $\alpha_k$ indicates that the sensor node is more sensitive to the energy insufficiency. Besides, we expect that any sensor node may need to upload its data in future time slots, but we do not know when it will be scheduled to transmit. For a conservative consideration, we impose the constraint (11b) to ensure that all sensor nodes will have sufficient energy to upload their data in the next time slot. If we remove (11b), it becomes possible that some node-$k$ may not have sufficient energy to upload its data after beamforming optimization. If this node-$k$ happens to be scheduled by the DRL agent in the next time slot, its data transmission will be unsuccessful and thus its AoI will continue increasing at the AP. Therefore, we include the constraint in (11b) as a one-step lookahead safety mechanism to ensure that every sensor has sufficient energy for data transmission when it is scheduled in the next time slot.

*1) Alternating optimization (AO) for problem* (11)*:* Given the uplink $(\mathbf{w}_u, \boldsymbol{\Theta}_u)$ and the control parameters $(\tau_k, p_k)_{k \in \mathcal{K}}$, the optimization of the downlink $(\mathbf{w}_d, \boldsymbol{\Theta}_d)$ in (11) can be decomposed by the AO method into two sub-problems, similar to that in [29]. In the first sub-problem, we optimize $\boldsymbol{\Theta}_d$ in problem (11) with the fixed $\mathbf{w}_d$. Note that only the IRS-enhanced downlink channel $\mathbf{f}_{d,k}$ relates to $\boldsymbol{\Theta}_d$ as shown in (2). We can simplify problem (11) by introducing a few auxiliary variables. Define $\mathbf{F}_k \triangleq \mathbf{G}_k \text{diag}(\mathbf{h}_k^r)$ and then we can simplify (2) as $\mathbf{f}_{d,k} = \mathbf{h}_k^d(t) + \mathbf{F}_k \boldsymbol{\theta}_d$. We further define $\bar{\boldsymbol{\theta}}_d \triangleq [\boldsymbol{\theta}_d, \zeta]^T$ where $\zeta \geq 0$ and $|\zeta| = 1$. The quadratic term in (11a) can be rewritten as $|\mathbf{f}_{d,k}^H \mathbf{w}_d|^2 = \bar{\boldsymbol{\theta}}_d^H \mathbf{R}_k \bar{\boldsymbol{\theta}}_d + (\mathbf{h}_k^d)^H \mathbf{h}_k^d$, where the matrix coefficient $\mathbf{R}_k$ is given by
$$\mathbf{R}_k = \begin{bmatrix} \mathbf{F}_k^H \mathbf{w}_d \mathbf{w}_d^H \mathbf{F}_k & \mathbf{F}_k^H \mathbf{w}_d \mathbf{w}_d^H \mathbf{h}_k^d \\ \left(\mathbf{h}_k^d\right)^H \mathbf{w}_d \mathbf{w}_d^H \mathbf{F}_k & 0 \end{bmatrix}.$$
We can further apply SDR to $\bar{\boldsymbol{\theta}}_d^H \mathbf{R}_k \bar{\boldsymbol{\theta}}_d$ by introducing the matrix variable $\boldsymbol{\Phi}_d = \bar{\boldsymbol{\theta}}_d \bar{\boldsymbol{\theta}}_d^H$. Similar transformation can be applied to the energy budget constraint in (11b). As such, the optimization of the downlink $\boldsymbol{\Theta}_d$ can be converted into the

---

**Algorithm 1** AO Method for Downlink Energy Transfer

---

**Input:** The channel information $\{\mathbf{h}_k^d(t), \mathbf{G}_k(t), \mathbf{h}_k^r(t)\}$, AoI state $A_k$, energy state $E_k$, and energy demand $E_k^c$ of each sensor node $k \in \mathcal{K}$

**Initialize:** a feasible beamforming strategy $(\mathbf{w}_d, \boldsymbol{\Theta}_d)$, $\tau \leftarrow 0$

1: $v_k \leftarrow A_k + \alpha_k E_k^{-1}$

2: $E_d^{(\tau)} \leftarrow 0$, $E_d^{(\tau+1)} \leftarrow \sum_{k \in \mathcal{K}} v_k |\mathbf{f}_{d,k}^H \mathbf{w}_d|^2$

3: **while** $||E_d^{(\tau+1)} - E_d^{(\tau)}|| \geq \epsilon$ **do**

4:      $\tau \leftarrow \tau + 1$

5:      Solve SDP (12) by the interior-point algorithm

6:      Extract the rank-one passive beamformer $\boldsymbol{\Theta}_d$

7:      Update $\mathbf{w}_d$ with the fixed $\boldsymbol{\Theta}_d$

8:      $E_d^{(\tau)} \leftarrow E_d^{(\tau+1)}$

9:      $E_d^{(\tau+1)} \leftarrow \sum_{k \in \mathcal{K}} v_k |\mathbf{f}_{d,k}^H \mathbf{w}_d|^2$

10: **end while**

---

following SDP similar to that in [29] and [31].

$$\max_{\boldsymbol{\Phi}_d \succeq 0} \quad \sum_{k \in \mathcal{K}} v_k \mathbf{Tr}\big(\mathbf{R}_k \boldsymbol{\Phi}_d\big) \tag{12a}$$

$$\text{s.t.} \quad \mathbf{Tr}\big(\mathbf{R}_k \boldsymbol{\Phi}_d\big) \geq (\eta p_s)^{-1}(E_k^c - E_k), \quad \forall\, k \in \mathcal{K}, \tag{12b}$$

$$\boldsymbol{\Phi}_d(n, n) = 1, \quad \forall\, n \in \mathcal{N}, \tag{12c}$$

where $\mathbf{Tr}(\cdot)$ denotes the matrix trace. Given the constant weight $v_k$, problem (12) becomes a conventional beamforming optimization for downlink MISO system [32], which can be solved efficiently by the interior-point algorithm. In the second sub-problem, we optimize $\mathbf{w}$ in problem (11) with the fixed $\boldsymbol{\Theta}_d$ and $(\tau_k, p_k)$. This follows a similar SDR approach as that in (12) by introducing a matrix variable $\mathbf{W}_d = \mathbf{w}_d \mathbf{w}_d^H$. Once the optimal solution $\mathbf{W}_d$ or $\boldsymbol{\Phi}_d$ is obtained, we can extract the rank-one beamformer $\mathbf{w}_d$ or $\boldsymbol{\theta}_d$ by Gaussian randomization method. We continue the iterations between $\mathbf{w}_d$ and $\boldsymbol{\Theta}_d$ until the convergence to a stable point.

*2) Simple approximation to problem* (11)*:* Given the scheduling decision $\psi_k(t)$, the inner-loop optimization estimates the beamforming strategies $(\mathbf{w}_d, \boldsymbol{\Theta}_d)$ and the transmission parameters $(\tau_k, p_k)_{k \in \mathcal{K}}$. The inner-loop optimization should be very efficient to minimize the computational

---

**Algorithm 2** Simple Approximation for $(\mathbf{w}_d, \boldsymbol{\Theta}_d)$

---

**Input:** The channel information $\{\mathbf{h}_k^d(t), \mathbf{G}_k(t), \mathbf{h}_k^r(t)\}$, AoI state $A_k$, energy state $E_k$, and energy

demand $E_k^c$ of each sensor node $k \in \mathcal{K}$

1: $v_k \leftarrow A_k + \alpha_k E_k^{-1}$

2: Solve problem (13) in the optimistic case

3: $\mathbf{f}_{d,k}^m \leftarrow \arg\min_{k \in \mathcal{K}} ||\mathbf{f}_{d,k}||^2$

4: $\mathbf{w}_d \leftarrow \mathbf{f}_{d,k}^m / ||\mathbf{f}_{d,k}^m||$

5: Solve problem (12) with the fixed $\mathbf{w}_d$

6: Extract the passive beamforming strategy $\boldsymbol{\Theta}_d$

---

overhead and processing delay in each iteration. Note that the AO algorithm for $(\mathbf{w}_d, \boldsymbol{\Theta}_d)$ can be inefficient as each iteration requires to solve the SDP problem (12) with a high computational complexity. The number of AO iterations till convergence is also unknown. Instead of the AO method, we further propose a simple approximation to problem (11) by optimizing $\boldsymbol{\Theta}_d$ with a fixed and feasible $\mathbf{w}_d$. This solution can avoid the iterations between $\mathbf{w}_d$ and $\boldsymbol{\Theta}_d$, and thus improve the inner-loop computation efficiency. In particular, we firstly consider an optimistic case in which the AP's downlink beamforming vector $\mathbf{w}_d$ can be aligned to all users' downlink channels $\mathbf{f}_{d,k}^H$, and thus we can relax problem (11) as follows:

$$\max_{\boldsymbol{\Theta}_d} \sum_{k \in \mathcal{K}} v_k ||\mathbf{f}_{d,k}||^2, \quad \text{s.t.} \quad (11b) - (11c), \tag{13}$$

which only relies on $\boldsymbol{\Theta}_d$ and can be solved by a similar SDR method for problem (9). However, problem (13) overestimates the total energy transfer to all sensor nodes. In the second step, we can reorder the sensor nodes by the channel gain $||\mathbf{f}_{d,k}||^2$ and fix the downlink beamforming vector as $\mathbf{w}_d = \mathbf{f}_{d,k}^m / ||\mathbf{f}_{d,k}^m||$, where $\mathbf{f}_{d,k}^m = \arg\min_{k \in \mathcal{K}} ||\mathbf{f}_{d,k}||^2$. This intuition ensures that we transfer more RF energy to the sensor node with the worst channel condition. In the third step, we optimize $\boldsymbol{\Theta}_d$ by solving (12) with the fixed $\mathbf{w}_d$, which provides a feasible lower bound to problem (11). As such, we only need to solve two SDPs to approximate the solution $(\mathbf{w}_d, \boldsymbol{\Theta}_d)$.

## VI. OUTER-LOOP LEARNING FOR SCHEDULING

The outer-loop DRL approach aims to update the AP's scheduling policy by continuously interacting with the network environment. We can reformulate the scheduling optimization problem

into the Markov decision process (MDP), which can be characterized by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{R})$. The state space $\mathcal{S}$ denotes the set of all system states. In each decision epoch, the AP's state $\mathbf{s}_t \in \mathcal{S}$ includes all sensor nodes' AoI values, denoted as a vector $\mathbf{A}(t) \triangleq [A_1(t), A_2(t), \ldots, A_K(t)]$, and the energy status denoted as $\mathbf{E}(t) = [E_1(t), E_2(t), \ldots, E_K(t)]$. Hence, we can define the system state in each time slot as $\mathbf{s}_t = (\mathbf{A}(t), \mathbf{E}(t))$. For each sensor node-$k$, we have $A_k(t) \geq 1$ as its AoI value can keep increasing from 1. The energy status $E_k(t)$ is upper bounded by the maximum battery capacity, i.e., $E_k(t) \in [0, E_{\max}]$. The action space $\mathcal{A}$ denotes the set of all feasible scheduling decisions $\mathbf{a}_t \triangleq \{\psi_0(t), \psi_1(t), \ldots, \psi_K(t)\} \in \{0,1\}^{K+1}$ that satisfies the inequality in (1). Given the AP's scheduling decision, we can obtain $(\mathbf{w}_d, \mathbf{\Theta}_d)$ by the inner-loop optimization and then update the AoI performance of different sensor nodes. The reward $\mathcal{R}$ assigns each state-action pair an immediate value $v_t(\mathbf{s}_t, \mathbf{a}_t)$. It also influences the DRL agent's action adaptation to maximize the long-term reward, namely, the value function $V_\pi(\mathbf{s}_0) \triangleq \sum_{t \in \mathcal{T}} \varepsilon^t v_t(\mathbf{s}_t, \mathbf{a}_t)$, where $\varepsilon \in (0,1)$ is the discount factor for cumulating the reward and $\pi(\mathbf{a}_t|\mathbf{s}_t)$ denotes the policy function mapping each state $\mathbf{s}_t$ to the action $\mathbf{a}_t$. Specifically, considering the AoI minimization in (7), we can define the reward $v_t(\mathbf{s}_t, \mathbf{a}_t)$ in the $t$-th time step as follows:

$$v_t(\mathbf{s}_t, \mathbf{a}_t) = -\frac{1}{K|\mathcal{H}_t|} \sum_{\tau \in \mathcal{H}_t} \sum_{k \in \mathcal{K}} \lambda_k A_k(\tau), \tag{14}$$

where $\mathcal{H}_t \triangleq \{t - t_o, \ldots, t - 1, t\} \subset \mathcal{T}$ denotes a set of past time slots with the length $|\mathcal{H}_t|$. Hence, the reward $v_t(\mathbf{s}_t, \mathbf{a}_t)$ can be considered as the averaged AoI of all sensor nodes in the most recent sliding window of the past time slots.

DRL approaches such as the value-based DQN and policy-based policy gradient (PG) algorithms have been demonstrated to be effective for solving MDP in large-scale wireless network by using DNNs to approximate either the value function $V_\pi(\mathbf{s}_t)$ or the policy function $\pi(\mathbf{a}_t|\mathbf{s}_t)$ [33]. In particular, the DQN method is an extension of the classic Q-learning method [34], by using DNN to approximate the Q-value function $Q_{\boldsymbol{\mu}}(\mathbf{s}_t, \mathbf{a}_t)$, where $\boldsymbol{\mu}$ denotes the DNN weight parameters for the Q-value network. Starting from the initial state $\mathbf{s}_0$, the PG algorithms directly optimize the value function $V_\pi(\mathbf{s}_0)$ by using gradient-based approaches to update the policy $\pi_{\boldsymbol{\omega}}(\mathbf{a}_t|\mathbf{s}_t)$, where $\boldsymbol{\omega}$ denotes the DNN weight parameters for the policy network. The proximal policy optimization (PPO) algorithm is recently proposed in [10] as a sample-efficient and easy-to-implement PG algorithm, striking a favorable balance between complexity, simplicity, and learning efficiency. In the sequel, we implement both the DQN and PPO algorithms and compare

their performance for the outer-loop scheduling optimization.

## A. DQN Algorithm for Scheduling

The DQN algorithm relies on two DNNs to stabilize the learning, denoted by the online Q-network and target Q-network. Given the DNN parameters $\boldsymbol{\mu}_t$ and $\boldsymbol{\mu}'_t$ for the two Q-networks, their outputs are given by $Q_{\boldsymbol{\mu}_t}(\mathbf{s}_t, \mathbf{a}_t)$ and $Q_{\boldsymbol{\mu}'_t}(\mathbf{s}_t, \mathbf{a}_t)$, respectively. At each learning episode, the DQN agent observes the system state $\mathbf{s}_t = (\mathbf{A}(t), \mathbf{E}(t))$ and chooses the best scheduling action $\mathbf{a}_t$ with the maximum value of $Q_{\boldsymbol{\mu}_t}(\mathbf{s}_t, \mathbf{a}_t)$. To enable random exploration, the DQN agent can also take a random action with a small probability. Once the action $\mathbf{a}_t$ is fixed and executed in the network, the DQN agent observes the instant reward $v_t(\mathbf{s}_t, \mathbf{a}_t)$ and records the transition to the next state $\mathbf{s}_{t+1}$. Each transition sample $(\mathbf{s}_t, \mathbf{a}_t, v_t, \mathbf{s}_{t+1})$ will be stored in the experience replay buffer. Meanwhile, the DQN agent estimates the target Q-value as $y_t = v_t(\mathbf{s}_t, \mathbf{a}_t) + \varepsilon Q_{\boldsymbol{\mu}'_t}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})$ by using the target Q-network with a different weight parameter $\boldsymbol{\mu}'_t$.
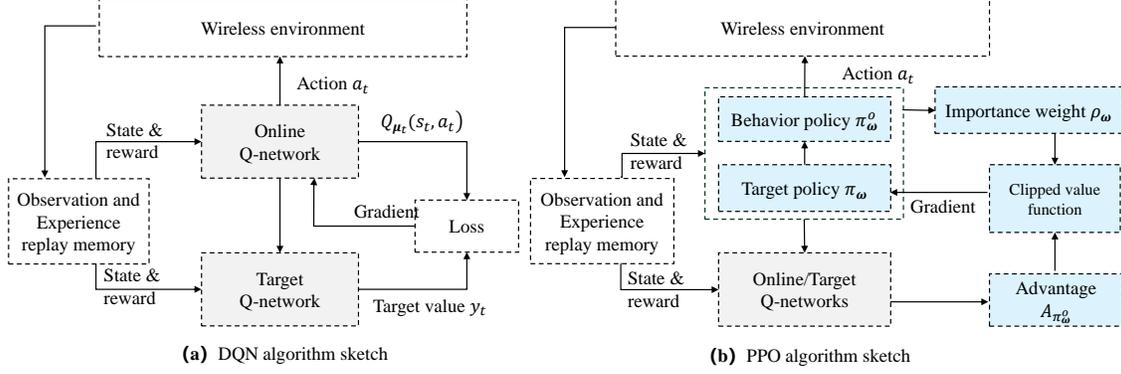
DQN's success lies in the design of the experience replay mechanism that improves the learning efficiency by reusing the historical transition samples to train the DNN in each iteration [33]. The DNN training aims to adjust the parameter $\boldsymbol{\mu}_t$ to minimize a loss function $\ell(\boldsymbol{\omega}_t)$, which is defined as the gap or more formally the temporal-difference (TD) error between the online Q-network $Q_{\boldsymbol{\mu}_t}(\mathbf{s}_t, \mathbf{a}_t)$ and the target value $y_t$, specified as follows:

$$\ell(\boldsymbol{\mu}_t) = \mathbb{E}[|y_t - Q_{\boldsymbol{\mu}_t}(\mathbf{s}_t, \mathbf{a}_t)|^2]. \tag{15}$$

The expectation in (15) is taken over a random subset (i.e., mini-batch) of transition samples from the experience replay buffer. For each mini-batch sample, we can evaluate the target value $y_t$ and generate the Q-value $Q_{\boldsymbol{\mu}_t}(\mathbf{s}_t, \mathbf{a}_t)$. The weight parameters $\boldsymbol{\mu}_t$ can be updated by the backpropagation of the gradient information. The DQN method stabilizes the learning by periodically copying the DNN parameter $\boldsymbol{\mu}_t$ of the online Q-network to the target Q-network.

## B. Policy-based Actor-Critic Algorithms

Different from the value-based DQN method, the policy-based approach improves the value function by updating the parametric policy $\pi_{\boldsymbol{\omega}}$ in gradient-based methods [33]. Given the system state $\mathbf{s}_t$, the policy $\pi_{\boldsymbol{\omega}}(\mathbf{a}_t|\mathbf{s}_t)$ specifies a probability distribution over different actions $\mathbf{a}_t \in \mathcal{A}$. The DNN training aims at updating the policy network to improve the expected value function,

**Fig. 2:** The DQN and PPO algorithms for the outer-loop scheduling optimization.

rewritten as $J(\boldsymbol{\omega}) = \sum_{\mathbf{s}\in\mathcal{S}} d(\mathbf{s})V_\pi(\mathbf{s}) = \sum_{\mathbf{s}\in\mathcal{S}} d(\mathbf{s})\sum_{\mathbf{a}\in\mathcal{A}} \pi_{\boldsymbol{\omega}}(\mathbf{a}|\mathbf{s})Q^\pi(\mathbf{s},\mathbf{a})$, where $d(\mathbf{s})$ is the stationary state distribution corresponding to the policy $\pi_{\boldsymbol{\omega}}$ and $Q^\pi(\mathbf{s},\mathbf{a})$ denotes the Q-value of the state-action pair $(\mathbf{s},\mathbf{a})$ following the policy $\pi_{\boldsymbol{\omega}}$. Now the expected value function $J(\boldsymbol{\omega})$ becomes a function of the policy parameter $\boldsymbol{\omega}$. The policy gradient theorem in [35] simplifies the evaluation of the policy gradient $\nabla_{\boldsymbol{\omega}} J(\boldsymbol{\omega})$ as follows:

$$\nabla_{\boldsymbol{\omega}} J(\boldsymbol{\omega}) = \mathbb{E}_\pi\Big[Q^\pi(\mathbf{s},\mathbf{a})\nabla_{\boldsymbol{\omega}} \ln \pi_{\boldsymbol{\omega}}(\mathbf{a}|\mathbf{s})\Big], \tag{16}$$

where the expectation is taken over all possible state-action pairs following the same policy $\pi_{\boldsymbol{\omega}}$.

For practical implementation, the policy gradient $\nabla_{\boldsymbol{\omega}} J(\boldsymbol{\omega})$ can be evaluated by sampling the historical decision-making trajectories. At each learning epoch $t$, the DRL agent interacts with the environment by the state-action pair $(\mathbf{s}_t, \mathbf{a}_t)$, collects an immediate reward $v_t$, and observes the transition to the next state $\mathbf{s}_{t+1}$. Let $\boldsymbol{\ell} \triangleq \{\mathbf{s}_1, \mathbf{a}_1, v_1, \mathbf{s}_2, \mathbf{a}_2, v_2, \mathbf{s}_3, \ldots, v_T\}$ denote the state-action trajectory as the DRL agent interacts with the environment. We can estimate the Q-value $Q^\pi(\mathbf{s}_t, \mathbf{a}_t)$ in (16) by $g_t = \sum_{\tau=t}^{T} \varepsilon^{\tau-t} v_t$. As such, we can approximate the policy gradient $\nabla_{\boldsymbol{\omega}} J(\boldsymbol{\omega})$ in each time step by the random sample $g_t \nabla_{\boldsymbol{\omega}} \ln \pi_{\boldsymbol{\omega}}(\mathbf{a}_t|\mathbf{s}_t)$ and update the policy network as $\boldsymbol{\omega} \leftarrow \boldsymbol{\omega} + \alpha_{\boldsymbol{\omega}} g_t \nabla_{\boldsymbol{\omega}} \ln \pi_{\boldsymbol{\omega}}(\mathbf{a_t}|\mathbf{s_t})$, where $\alpha_{\boldsymbol{\omega}}$ denotes the step-size for gradient update. Besides the stochastic approximation, we can also employ another DNN to approximate the Q-value $Q^\pi(\mathbf{s}_t, \mathbf{a}_t)$ in (16) and replace the Monte Carlo estimation $g_t$ by the DNN approximation $Q_{\boldsymbol{\mu}}(\mathbf{s}_t, \mathbf{a_t})$ with the weight parameter $\boldsymbol{\mu}$, similar to that in the DQN algorithm. This motivates the actor-critic framework to update both the policy network and the Q-value network [35]. The actor updates the policy network while the critic updates the Q-network by minimizing a loss function

similar to (15). We can take the derivative of the loss function and update the weight parameter as $\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} + \alpha_{\boldsymbol{\mu}}\delta_t\nabla_{\boldsymbol{\mu}}Q_{\boldsymbol{\mu}}(\mathbf{s}_t, \mathbf{a_t})$, where $\delta_t = y_t - Q_{\boldsymbol{\mu}_t}(\mathbf{s}_t, \mathbf{a}_t)$ denotes the TD error. The Q-value estimation can be also replaced by the advantage function $A_{\pi}(\mathbf{s}_t, \mathbf{a}_t) \triangleq Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) - V_{\pi}(\mathbf{s}_t)$ to reduce the variability in predictions and improve the learning efficiency.

The gradient estimation in (16) requires a complete trajectory by using the same target policy $\pi_{\boldsymbol{\omega}}$. It is actually the on-policy method that refrains us from using past experiences and limits the sample efficiency. This drawback can be avoided by a minor revision to the policy gradient:

$$\nabla_{\boldsymbol{\omega}} J(\boldsymbol{\omega}) = \mathbb{E}_{\pi_{\boldsymbol{\omega}}^o}\left[\frac{\pi_{\boldsymbol{\omega}}(\mathbf{s}, \mathbf{a})}{\pi_{\boldsymbol{\omega}}^o(\mathbf{s}, \mathbf{a})} A_{\pi_{\boldsymbol{\omega}}}(\mathbf{s}, \mathbf{a})\nabla_{\boldsymbol{\omega}} \ln \pi_{\boldsymbol{\omega}}(\mathbf{a}|\mathbf{s})\right], \tag{17}$$

where the behavior policy $\pi_{\boldsymbol{\omega}}^o$ is used to collect the training samples. This becomes the off-policy gradient that allows us to estimate it by using the past experience collected from a different and even obsolete behavior policy $\pi_{\boldsymbol{\omega}}^o$. Hence, we can improve the sample efficiency by maintaining the experience replay buffer, similar to the DQN method. To further improve training stability, the off-policy trust region policy optimization (TRPO) method imposes an additional constraint on the gradient update [36], i.e., the new policy $\pi_{\boldsymbol{\omega}}$ should not change too much from the old policy $\pi_{\boldsymbol{\omega}}^o$. Therefore, the policy optimization is to solve the following constrained problem:

$$\max_{\boldsymbol{\omega}} \quad \mathbb{E}_{\pi_{\boldsymbol{\omega}}^o}\left[\frac{\pi_{\boldsymbol{\omega}}(\mathbf{s}, \mathbf{a})}{\pi_{\boldsymbol{\omega}}^o(\mathbf{s}, \mathbf{a})} A_{\pi_{\boldsymbol{\omega}}^o}(\mathbf{s}, \mathbf{a})\right], \quad \text{s.t.} \quad \mathbb{E}_{\pi_{\boldsymbol{\omega}}^o}[D_{KL}(\pi_{\boldsymbol{\omega}}, \pi_{\boldsymbol{\omega}}^o)] \leq \delta_{KL}, \tag{18}$$

where $D_{KL}(P_1, P_2) \triangleq \int_{-\infty}^{\infty} P_1(x)\log\left(P_1(x)/P_2(x)\right)\,dx$ represents a distance measure in terms of the Kullback-Leibler (KL) divergence between two probability distributions [36]. The inequality constraint in (18) enforces that the KL divergence between two policies $\pi_{\boldsymbol{\omega}}$ and $\pi_{\boldsymbol{\omega}}^o$ are bounded by $\delta_{KL}$. The advantage function $A_{\pi_{\boldsymbol{\omega}}^o}$ in the objective of (18) is the approximation of the true advantage $A_{\pi_{\boldsymbol{\omega}}}$ corresponding to the target policy $\pi_{\boldsymbol{\omega}}$. However, the exact solution to the optimization problem (18) is not easy. Normally we require the first- and second-order approximations for both the objective and the constraint in (18).

## C. PPO Algorithm for Scheduling

The proximal policy optimization (PPO) algorithm proposed in [10] further improves the objective in (18) by limiting the probability ratio or the importance weight $\rho_{\boldsymbol{\omega}}(\mathbf{s}, \mathbf{a}) \triangleq \frac{\pi_{\boldsymbol{\omega}}(\mathbf{s},\mathbf{a})}{\pi_{\boldsymbol{\omega}}^o(\mathbf{s},\mathbf{a})}$ within a safer region $[1 - \epsilon, 1 + \epsilon]$.

$$\tilde{J}(\pi_{\boldsymbol{\omega}}) = \mathbb{E}_{\pi_{\boldsymbol{\omega}}^o}\left[\min\left\{\rho_{\boldsymbol{\omega}} A_{\pi_{\boldsymbol{\omega}}^o}, \text{clip}(\rho_{\boldsymbol{\omega}}, 1 - \epsilon, 1 + \epsilon)A_{\pi_{\boldsymbol{\omega}}^o}\right\}\right], \tag{19}$$

---

**Algorithm 3** Energy-and-AoI-Aware Scheduling and Transmission Control Algorithm

---

**Initialize:** Target policy $\pi_{\boldsymbol{\omega}}$ and behavior policy $\pi_{\boldsymbol{\omega}}^o$,

**Initialize:** $t \leftarrow 0$, $E_k(0) \leftarrow B$, $A_k(t) \leftarrow 0$

1: **for** Episode $= \{1, 2, \ldots, \text{MAX} = 3000\}$ **do**

2:     **while** $t \neq T$ **do**

3:         Observe the system state $(\mathbf{A}(t), \mathbf{E}(t))$

4:         Choose the outer-loop action $\boldsymbol{\Psi}(t)$ for scheduling

5:         **case** $\psi_0(t) = 0$: optimize uplink data transmission in (8)

6:         **case** $\psi_0(t) = 1$: optimize downlink energy transfer in (11)

7:         Execute joint action $\mathbf{a}(t) \triangleq (\boldsymbol{\Psi}(t), \mathbf{w}(t), \boldsymbol{\Phi}(t))$, evaluate $v_t(\mathbf{s}_t, \mathbf{a}_t)$

8:         Record the next state $\mathbf{s}_{t+1}$ and buffer the transition $(\mathbf{s}_t, \mathbf{a}_t, v_t, \mathbf{s}_{t+1})$

9:         $t \leftarrow t + 1$

10:     **end while**

11:     Take mini-batch form the experience replay buffer

12:     Estimate advantage $A_{\pi_{\boldsymbol{\omega}}^o}$ and update $\boldsymbol{\omega}$ to maximize (20)

13:     Update behavior policy $\pi_{\boldsymbol{\omega}}^o \leftarrow (1 - \mu)\pi_{\boldsymbol{\omega}}^o + \mu\pi_{\boldsymbol{\omega}}$

14: **end for**

---

where the function $\text{clip}(\rho_{\boldsymbol{\omega}}, 1 - \epsilon, 1 + \epsilon)$ returns $\rho_{\boldsymbol{\omega}}$ if $\rho_{\boldsymbol{\omega}} \in [1 - \epsilon, 1 + \epsilon]$ and returns $1 - \epsilon$ (or $1 + \epsilon$) if $\rho_{\boldsymbol{\omega}} < 1 - \epsilon$ (or $\rho_{\boldsymbol{\omega}} > 1 + \epsilon$). The parameter $\epsilon$ is used to control the clipping range. The approximate value function $\tilde{J}(\pi_{\boldsymbol{\omega}})$ ensures that the target policy $\pi_{\boldsymbol{\omega}}$ will not deviate too far from the behavior policy $\pi_{\boldsymbol{\omega}}^o$ for either positive or negative advantage $A_{\pi_{\boldsymbol{\omega}}^o}$. With the clipped value function $\tilde{J}(\pi_{\boldsymbol{\omega}})$, we further introduce a Lagrangian dual variable $\beta_{KL}$ and reformulate the constrained problem (18) into the unconstrained maximization as follows:

$$\max_{\boldsymbol{\omega}} \ \tilde{J}(\pi_{\boldsymbol{\omega}}) - \beta_{KL}\mathbb{E}_{\pi_{\boldsymbol{\omega}}^o}\Big[D_{KL}(\pi_{\boldsymbol{\omega}}, \pi_{\boldsymbol{\omega}}^o)\Big], \tag{20}$$

The policy parameter $\boldsymbol{\omega}$ for the new value function in (20) can be easily updated by using the stochastic gradient ascent method. The parameter $\beta_{KL}$ can be also adapted according to the difference between the measured KL divergence $\mathbb{E}_{\pi_{\boldsymbol{\omega}}^o}\Big[D_{KL}(\pi_{\boldsymbol{\omega}}, \pi_{\boldsymbol{\omega}}^o)\Big]$ and its target $\delta_{KL}$.

As shown in Fig. 2, we highlight the comparison between the DQN and the PPO algorithms for outer-loop scheduling optimization. Different from the DQN algorithm, the PPO algorithm

**TABLE II:** Parameter settings

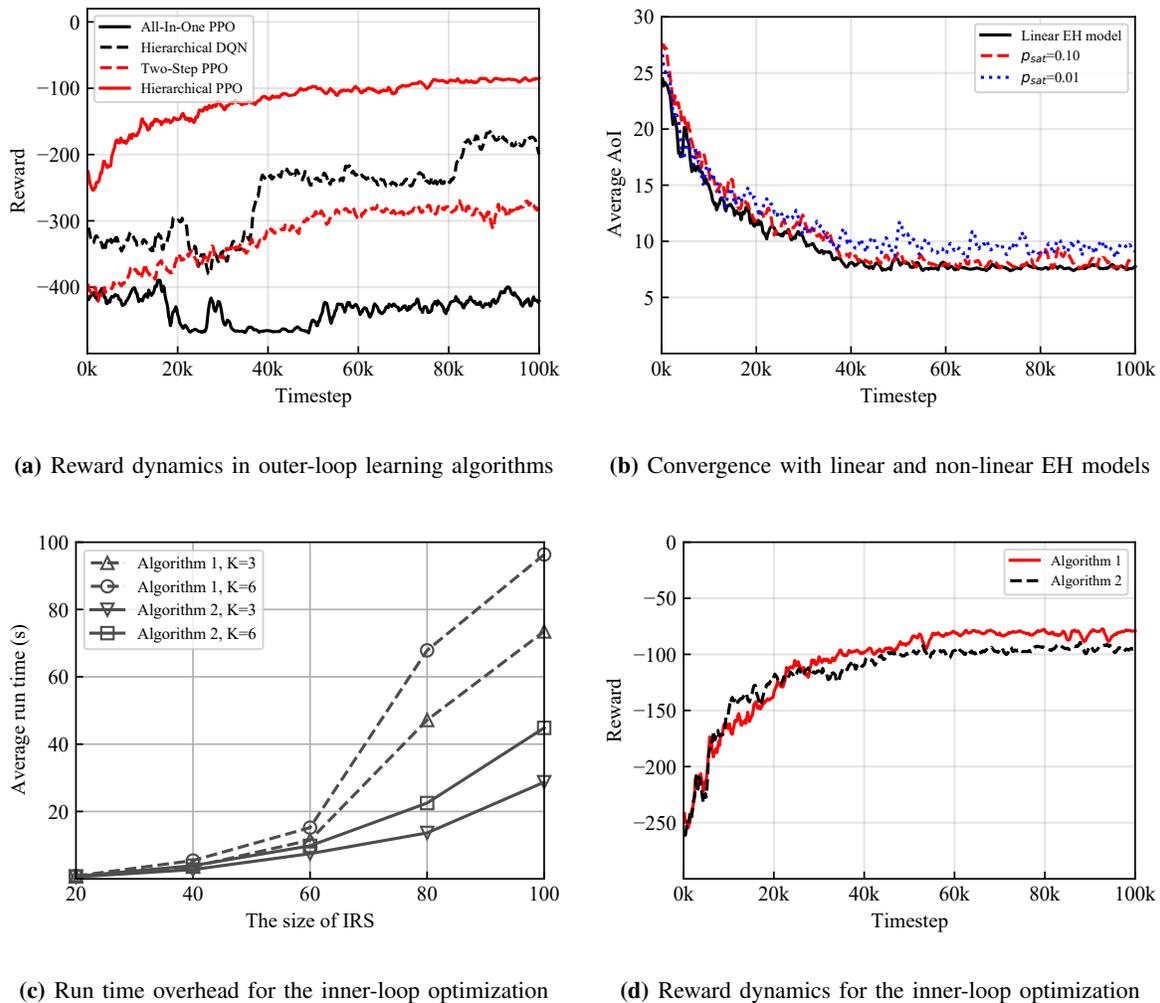| Parameters | Values | Parameters | Values |
|---|---|---|---|
| Number of DNN hidden layers | 3 | Optimizer | Adam |
| Actor's learning rate | 0.0005 | Number of AP's antennas | 4 |
| Critic's learning rate | 0.001 | AP's transmit power $p_s$ | 30dBm |
| Reward discount | 0.99 | Energy conversion efficiency $\eta$ | 0.9 |
| Number of neurons | 64 | Noise power $\sigma^2$ | $-75$dBm |
| Activation function | Tanh and Softmax | Sensor nodes' data size $D$ | 5Kbits |

employs the actor-critic structure that relies on two sets of DNNs to approximate the policy networks. The difference between the target policy network and the behavior policy network is used to generate the importance sampling weight $\rho_{\boldsymbol{\omega}}(\mathbf{s}, \mathbf{a})$. The behavior policy network is used to interact with the environment and store the transition samples in the experience replay buffer. The target policy network is used to update the DNN weight parameter $\boldsymbol{\omega}$ by sampling a mini-batch randomly from the experience replay buffer. Then we can update the target policy network by maximizing the objective in (20). The complete solution procedures are listed in Algorithm 3. Considering the preferable learning efficiency and robustness of the PPO algorithm, we integrate it in Algorithm 3 to adapt the outer-loop scheduling policy in the hierarchical learning framework. At the initialization stage, we randomly initialize the DNN weight parameters $\boldsymbol{\omega}$ for the policy network. In each learning episode, the AP collects observations $(\mathbf{A}(t), \mathbf{E}(t))$ of the system and decides the outer-loop scheduling decision $\boldsymbol{\Psi}(t)$, as shown in line 4 of Algorithm 3. Given the scheduling decision $\boldsymbol{\Psi}(\mathbf{t})$, the AP needs to optimize the joint beamforming strategies $(\mathbf{w}(t), \boldsymbol{\Theta}(t))$ for either uplink information transmission or downlink energy transfer, corresponding to lines 5-6 of Algorithm 3. Note that the solution to downlink energy transfer can be determined by either the iterative Algorithm 1 or the simplified Algorithm 2. When we determine both the outer-loop and inner-loop decision variables, we can execute the joint action $(\boldsymbol{\Psi}(t), \mathbf{w}(t), \boldsymbol{\Theta}(t))$ in the wireless system and evaluate the reward function as shown in lines 7-8 of Algorithm 3. The DNN training is based on the random mini-batch sampled from the experience replay buffer, as shown in lines 11-13 of Algorithm 3. For performance comparison, the DQN algorithm is also implemented for outer-loop scheduling. Our numerical evaluation in Section VII reveals that the PPO-based algorithm can improve the convergence performance and achieve a lower AoI value compared to the DQN-based algorithm.

## VII. SIMULATION RESULTS

In this section, we present simulation results to verify the performance gain of the proposed algorithms for the wireless-powered and IRS-assisted wireless sensor network. The $(x, y, z)$-coordinates of the AP and the IRS in meters are given by $(200, -200, 0)$ and $(0, 0, 0)$, respectively. The sensor nodes are randomly distributed in a rectangular area $[5, 35] \times [-35, 35]$ in the $(x, y)$-plane with $z = -20$. We consider that the direct channel from the AP to each sensor node-$k$ follows the Rayleigh fading distribution, i.e., $\mathbf{h}_k^d(t) = \beta_{0,k} \tilde{\mathbf{h}}_k^d(t)$, where $\tilde{\mathbf{h}}_k^d(t) \sim \mathcal{CN}(0, \boldsymbol{I})$ denotes the complex Gaussian random variable and $\beta_{0,k}$ denotes the path-loss of the direct channel. Given the distance $d_k^{\mathrm{AS}}$ from the AP to the sensor node-$k$, we have $\beta_{0,k} = 32.6 + 36.7 \log(d_k^{\mathrm{AS}})$. A similar channel model is employed in [37]. The IRS-sensor channel $\mathbf{h}_k^r(t)$ and the AP-IRS channel $\mathbf{G}_k(t)$ are modelled similarly. More detailed parameters are summarized in Table II.

### A. Improved Learning Efficiency and Convergence

In Fig. 3(a), we compare the reward performance among the hierarchical learning Algorithm 3, the hierarchical DQN and the conventional PPO, in which all decision variables $(\boldsymbol{\Psi}(t), \mathbf{w}(t), \boldsymbol{\Theta}(t))$ are adapted simultaneously. Hence, we denote the conventional PPO as the All-in-One PPO algorithm in Fig. 3. The solid line in red represents the dynamics of the AoI performance in the proposed hierarchical PPO algorithm, which spits the decision variables into two parts and optimizes the beamforming strategy $(\mathbf{w}(t), \boldsymbol{\Theta}(t))$ by the inner-loop Algorithm 1. The dotted line in red denotes hierarchical DQN algorithm, and the dash-dotted line in black denotes the conventional All-in-One PPO algorithm. It is clear that the All-in-One PPO may not converge effectively due to a huge action space in the mixed discrete and continuous domain. Compared with the PPO algorithm, the hierarchical DQN algorithm becomes unstable as shown in Fig. 3. The hierarchical PPO for joint scheduling and transmission control can reduce the action space in the outer-loop learning framework and thus achieve a significantly higher reward performance and faster convergence guided by the inner-loop optimization modules. We also implement the two-step learning algorithm (denoted as Two-Step PPO in Fig. 3) in which both inner-loop and outer-loop control variables are adapted by the PPO learning methods. The Two-Step PPO algorithm has a better convergence than that of the conventional All-in-One PPO algorithm. However, its reward performance is much inferior to that of the hierarchical learning framework, which is guided by the inner-loop optimization-driven target during the outer-loop learning.

**(a)** Reward dynamics in outer-loop learning algorithms

**(b)** Convergence with linear and non-linear EH models

**(c)** Run time overhead for the inner-loop optimization

**(d)** Reward dynamics for the inner-loop optimization

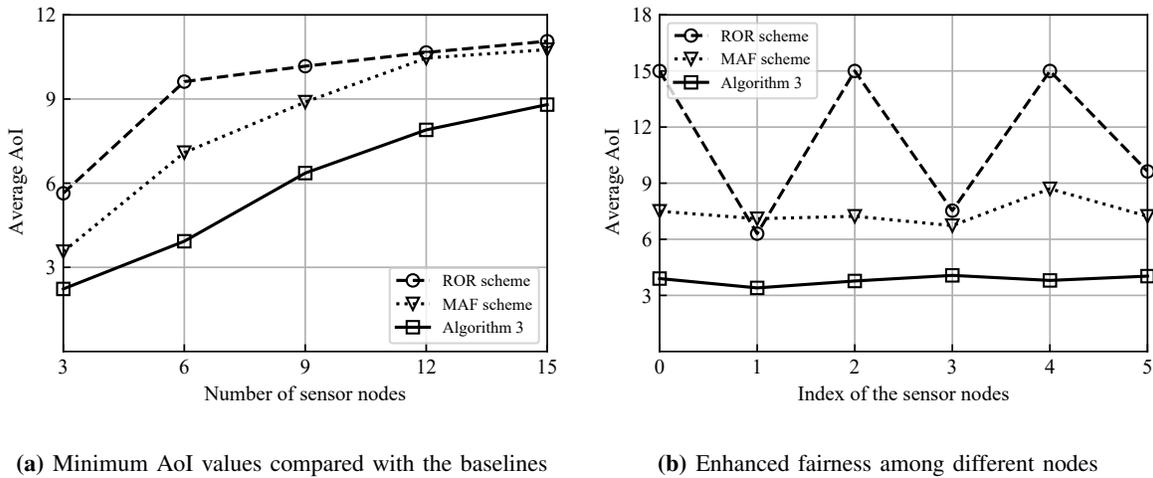**Fig. 3:** The outer-loop learning and inner-loop optimization of the hierarchical PPO algorithm.

Figure 3(b) reveals how different EH models effect the AoI performance. The linear EH model results in a slightly smaller AoI value than that of the non-linear EH model. The reason is that the linear EH model over-estimates the sensor nodes' EH capabilities. For the non-linear EH model, the harvested power will not further increase as the received signal power becomes higher than the saturation power $p_{\text{sat}}$. Fig. 3(b) shows that $p_{\text{sat}}$ also affects the AoI performance. Generally we can expect a better AoI performance with a higher saturation power $p_{\text{sat}}$. In our simulation, we also evaluate the overall reward and average AoI performances with different discount factor $\varepsilon \in \{0.99, 0.95, 0.90, 0.80\}$, which is used to accumulate the rewards in different decision epochs. The simulation results reveal that the learning with a larger discount factor becomes stable. We further show the run time and performance comparison between the iterative Algorithm 1 and

and the simplified Algorithm 2 for the inner-loop beamforming optimization. It is clear that the run time of each inner-loop optimization algorithm increases with the size of the IRS and the number of sensor nodes. Besides, we observe that Algorithm 2 significantly saves the run time by reducing the number of iterations, especially with a large-size IRS, as shown in Fig. 3(c). However, the reward performances of two algorithms are very close to each other, as shown in Fig 3(d). This implies that we can deploy Algorithm 2 preferably in practice.
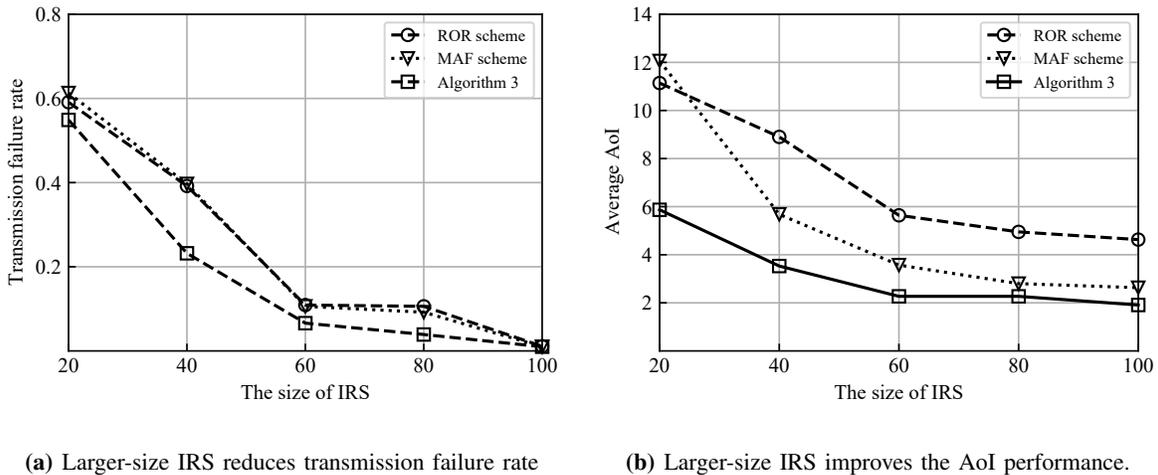
### B. Performance Gain over Existing Scheduling Policies

In this part, we develop two baselines policies to verify the performance gain of Algorithm 3. The first baseline is the round-robin (ROR) scheduling policy which periodically selects one sensor node to upload its status-update information. In each scheduling period, we jointly optimize the active and passive beamforming strategy to enhance the information transmission. The second baseline is the Max-Age-First (MAF) scheduling policy, i.e., the AP selects the sensor node with the highest AoI value to upload its sensing information [38]. Both baselines rely on the same EH policy, i.e., the AP starts downlink energy transfer only when the scheduled sensor node has insufficient energy capacity, e.g., below some threshold value. In Fig. 4(a), we show the AoI performance as we increase the number of sensor nodes. For different algorithms, we set the same coordinates for the AP and the IRS. Generally, different scheduling policies have a small AoI value with a few sensor nodes. The MAF policy performs better than the ROR policy, as it gives higher priorities to the sensor nodes with unsatisfactory AoI performance. As the number of nodes increases, Algorithm 3 always outperforms the baselines by adapting the scheduling strategy according to sensor nodes' stochastic data arrivals, as shown in Fig. 4(a).

In Fig. 4(b), we compare the fairness of scheduling by showing the average AoI of different sensor nodes. For a fair comparison, we set the same weight parameters $\lambda_k$ for all sensor nodes in (6). It can be seen from Fig. 4(b) that Algorithm 3 achieves a smaller AoI value for each sensor node. Moreover, different sensor nodes can achieve very similar AoI values, which implies the enhanced fairness in the scheduling policy by using Algorithm 3. The ROR scheme has a large deviation compared to that of the other two baselines. The reason is that the RoR scheme cannot adapt to the dynamic data arrival process. An interesting observation is that the MAF scheme also has a smaller AoI deviation among different sensor nodes. This is because that the MAF scheme always chooses the sensor node with the highest AoI value to upload its sensing

**(a)** Minimum AoI values compared with the baselines

**(b)** Enhanced fairness among different nodes

**Fig. 4:** Performance gain over baseline scheduling policies.



**(a)** Larger-size IRS reduces transmission failure rate

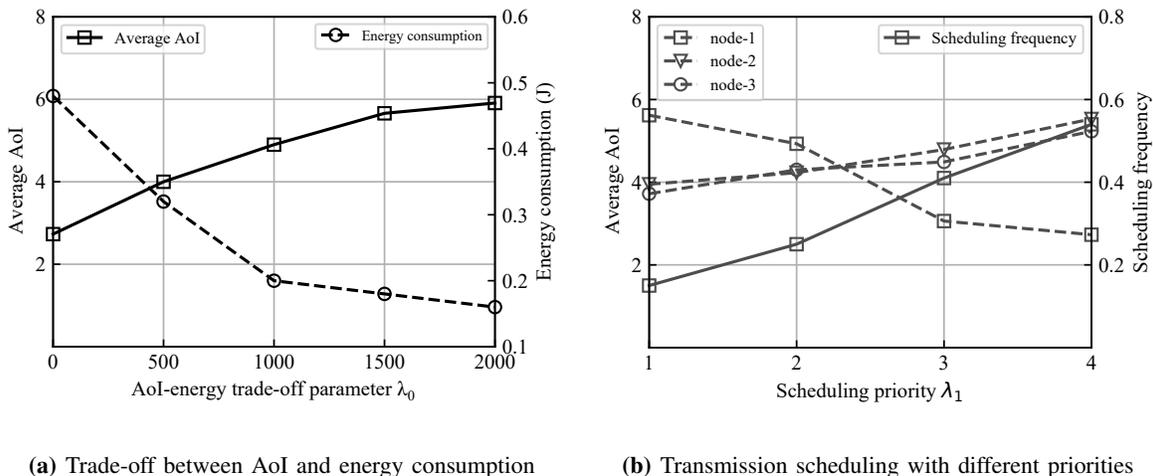**(b)** Larger-size IRS improves the AoI performance.

**Fig. 5:** Performance gain by using a larger-size IRS.

information. This can effectively prevent the AoI of some sensor nodes from being too high.

## C. Performance Gain by Using a Larger Size IRS

The use of IRS not only improves the downlink wireless energy transfer, but also enhances the uplink channels for sensing information transmissions. Both aspects implicitly improve the system's AoI performance. In this part, we intend to verify the performance gain achievable by using the IRS. In Fig. 5, we show the dynamics of the sensor nodes' transmission failure rate and the average AoI by increasing the size of IRS from 20 to 100. Specifically, we count the number of transmission failures within 30K time slots and visualize the transmission failure rate

**(a)** Trade-off between AoI and energy consumption    **(b)** Transmission scheduling with different priorities

**Fig. 6:** Energy-aware AoI minimization with different user priorities. We set $N = 80$ and $K = 3$ in the experiment.

in Fig. 5(a). It is quite intuitive that a larger-size IRS can reduce the sensor nodes' transmission failure rate, due to the IRS's reconfigurability to improve the channel quality. The increase in the IRS's size makes it more flexible to reshape the wireless channels and improve the uplink transmission success probability. This can help minimize the AoI values in a long run. However, the AoI values will not keep decreasing as the IRS's size increases. As shown in Fig. 5(b), the performance gap between our method and the baselines becomes smaller as the size increases. That is because the channel conditions become much better with a large-size IRS and thus the bottleneck of AoI performance becomes the scheduling delay, instead of the transmission delay.

### D. Trade-off Between AoI and Energy Consumption

In this part, we study the trade-off between AoI and energy consumption of the sensor nodes. Considering the AoI-energy tradeoff, we can revise the DRL agent's reward as follows:

$$v_t(\mathbf{s}_t, \mathbf{a}_t) = -\frac{1}{K|\mathcal{H}_t|} \sum_{\tau \in \mathcal{H}_t} \sum_{k \in \mathcal{K}} \Big( \lambda_k A_k(\tau) + \lambda_0 E_k^c(\tau) \Big),$$

where the AoI-energy trade-off parameter $\lambda_0$ is used to balance the sensor node's AoI and energy consumption. With a smaller $\lambda_0$, the sensor node becomes more aggressive to upload its information. This may cause energy outage due to insufficient energy supply to the sensor node. With a larger $\lambda_0$, the sensor node will focus more on its energy status and become more tolerant to the information delay. As shown in Fig. 6(a), given different $\lambda_0$, the AoI value and energy consumption have different trends of evolution. When $\lambda_0 = 0$, the average AoI can be reduced

by $53.8\%$ comparing to that with a higher $\lambda_0 = 2000$. We also show the performance gain with different priorities $\lambda_k$ for the sensor nodes. Considering three sensor nodes, we gradually increase $\lambda_1$ from 1 to 4 for the node-1 while setting fixed values for $\lambda_2 = \lambda_3 = 2$. We evaluate the scheduling frequency in 30K time slots and plot in Fig. 6(b) the change of node-$k$'s scheduling frequency, which is shown to increase linearly with respect to its priority $\lambda_1$. Fig. 6(b) also shows the change of three sensor nodes' AoI performance as we increase $\lambda_1$ for the node-1. It is clear to see that the node-1 will experience a much higher AoI value when it has a smaller priority, e.g., $\lambda_1 = 1$, than those of the other two nodes. When we gradually increase its priority, our algorithm will be more sensitive to the node-1's AoI performance and thus try to reduce its AoI by scheduling it more often, as revealed in Fig. 6(b). As such, the node-1 will take up more transmission opportunities by sacrificing the AoI performance of the other two nodes. This verifies that our algorithm can be adaptive to the change of sensor nodes' priorities.

## VIII. Conclusions

In this paper, we have focused on a wireless-powered and IRS-assisted network and aimed to minimize the overall AoI for information updates. We have formulated the AoI minimization problem as a mixed-integer program and devised a hierarchical learning framework, which includes the outer-loop model-free learning and the inner-loop optimization methods. Simulation results demonstrate that our algorithm can significantly reduce the average AoI and achieve controllable fairness among sensor nodes. More specifically, the hierarchical PPO algorithm achieves a significantly higher reward performance and faster convergence compared to the hierarchical DQN algorithm. It also outperforms typical baseline strategies in terms of the AoI performance and fairness. The performance gain can be more significant with a small size IRS.

## References

[1] M. A. Abd-Elmagid, N. Pappas, and H. S. Dhillon, "On the role of age of information in the internet of things," *IEEE Commun. Mag.*, vol. 57, no. 12, pp. 72–77, Dec. 2019.

[2] J. Lee, D. Niyato, Y. L. Guan, and D. I. Kim, "Learning to schedule joint radar-communication requests for optimal information freshness," in *proc. IEEE Intelligent Vehicles Symp. (IV)*, Jul. 2021, pp. 8–15.

[3] Y. Sun, E. Uysal-Biyikoglu, R. D. Yates, C. E. Koksal, and N. B. Shroff, "Update or Wait: How to keep your data fresh," *IEEE Trans. Inf. Theory*, vol. 63, no. 11, pp. 7492–7508, Nov. 2017.

[4] S. Gong, X. Lu, D. T. Hoang, D. Niyato, L. Shu, D. I. Kim, and Y.-C. Liang, "Toward smart wireless communications via intelligent reflecting surfaces: A contemporary survey," *IEEE Commun. Surv. Tut.*, vol. 22, no. 4, pp. 2283–2314, Jun. 2020.

[5] Z. Bao, Y. Dong, Z. Chen, P. Fan, and K. B. Letaief, "Age-optimal service and decision processes in internet of things," *IEEE Internet Things J.*, vol. 8, no. 4, pp. 2826–2841, Feb. 2021.

[6] C. Xu, Y. Xie, X. Wang, H. H. Yang, D. Niyato, and T. Q. S. Quek, "Optimal status update for caching enabled IoT networks: A dueling deep r-network approach," *IEEE Trans. Wireless Commun.*, vol. 20, no. 12, pp. 8438–8454, Dec. 2021.

[7] C. Zhou, G. Li, J. Li, Q. Zhou, and B. Guo, "FAS-DQN: Freshness-aware scheduling via reinforcement learning for latency-sensitive applications," *IEEE Trans. Comput.*, pp. 1–1, Nov. 2021.

[8] E. T. Ceran, D. Gündüz, and A. György, "A reinforcement learning approach to age of information in multi-user networks with HARQ," *IEEE J. Sel. Area. Commun.*, vol. 39, no. 5, pp. 1412–1426, May 2021.

[9] B.-M. Robaglia, A. Destounis, M. Coupechoux, and D. Tsilimantos, "Deep reinforcement learning for scheduling uplink IoT traffic with strict deadlines," in *proc. IEEE GLOBECOM*, Dec. 2021, pp. 1–6.

[10] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv:1707.06347*, 2017. [Online]. Available: http://arxiv.org/abs/1707.06347

[11] L. Cui, Y. Long, D. T. Hoang, and S. Gong, "Hierarchical learning approach for age-of-information minimization in wireless sensor networks," in *proc. IEEE Int. Symp. World Wirel. Mob. Multimed. Netw. (WoWMoM)*, Jun. 2022, pp. 1–7.

[12] J. Feng, W. Mai, and X. Chen, "Simultaneous multi-sensor scheduling based on double deep Q-learning under multi-constraint," in *proc. IEEE ICCC*, Jul. 2021, pp. 224–229.

[13] H. v. Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double $q$-learning," in *Proc. Thirtieth AAAI Conf. Artificial Intelligence*, 2016, pp. 2094–2100.

[14] X. Wu, X. Li, J. Li, P. C. Ching, and H. V. Poor, "Deep reinforcement learning for IoT networks: Age of information and energy cost tradeoff," in *proc. IEEE GLOBECOM*, Dec. 2020, pp. 1–6.

[15] S. Leng and A. Yener, "Age of information minimization for an energy harvesting cognitive radio," *IEEE Trans. Cogn. Commun. Network.*, vol. 5, no. 2, pp. 427–439, 2019.

[16] A. Arafa, J. Yang, S. Ulukus, and H. V. Poor, "Age-minimal transmission for energy harvesting sensors with finite batteries: Online policies," *IEEE Trans. Inf. Theory*, vol. 66, no. 1, pp. 534–556, 2020.

[17] X. Ling, J. Gong, R. Li, S. Yu, Q. Ma, and X. Chen, "Dynamic age minimization with real-time information preprocessing for edge-assisted IoT devices with energy harvesting," *IEEE Trans. Netw. Sci. Eng.*, vol. 8, no. 3, pp. 2288–2300, 2021.

[18] M. Hatami, M. Leinonen, and M. Codreanu, "Aoi minimization in status update control with energy harvesting sensors," *IEEE Trans. Commun.*, vol. 69, no. 12, pp. 8335–8351, 2021.

[19] M. A. Abd-Elmagid, H. S. Dhillon, and N. Pappas, "A reinforcement learning framework for optimizing age of information in RF-powered communication systems," *IEEE Trans. Commun.*, vol. 68, no. 8, pp. 4747–4760, Aug. 2020.

[20] A. H. Zarif, P. Azmi, N. Mokari, M. R. Javan, and E. Jorswieck, "AoI minimization in energy harvesting and spectrum sharing enabled 6G networks," *IEEE Trans. Green Commun. Network.*, pp. 1–1, 2022.

[21] S. Leng and A. Yener, "Learning to transmit fresh information in energy harvesting networks," *IEEE Trans Green Commun. Network.*, pp. 1–1, 2022.

[22] T. Bai, C. Pan, Y. Deng, M. Elkashlan, A. Nallanathan, and L. Hanzo, "Latency minimization for intelligent reflecting surface aided mobile edge computing," *IEEE J. Sel. Area. Commun.*, vol. 38, no. 11, pp. 2666–2682, Nov. 2020.

[23] Y. Cao, T. Lv, Z. Lin, and W. Ni, "Delay-constrained joint power control, user detection and passive beamforming in intelligent reflecting surface-assisted uplink mmwave system," *IEEE Trans. Cogn. Commun. Netw.*, vol. 7, no. 2, pp. 482–495, Jun. 2021.

[24] A. Muhammad, M. Elhattab, M. A. Arfaoui, A. Al-Hilo, and C. Assi, "Age of information optimization in a RIS-assisted wireless network," *arXiv:2103.06405*, 2021. [Online]. Available: https://arxiv.org/abs/2103.06405

[25] M. Samir, M. Elhattab, C. Assi, S. Sharafeddine, and A. Ghrayeb, "Optimizing age of information through aerial reconfigurable intelligent surfaces: A deep reinforcement learning approach," *IEEE Trans. Veh. Technol.*, vol. 70, no. 4, pp. 3978–3983, Apr. 2021.

[26] G. Cuozzo, C. Buratti, and R. Verdone, "A 2.4-GHz LoRa-based protocol for communication and energy harvesting on industry machines," *IEEE Internet Things J.*, vol. 9, no. 10, pp. 7853–7865, May 2022.

[27] J. Yao and N. Ansari, "Wireless power and energy harvesting control in IoD by deep reinforcement learning," *IEEE Trans. Green Commun. Netw.*, vol. 5, no. 2, pp. 980–989, Jun. 2021.

[28] B. Lyu, P. Ramezani, D. T. Hoang, S. Gong, Z. Yang, and A. Jamalipour, "Optimized energy and information relaying in self-sustainable IRS-empowered WPCN," *IEEE Trans. Commun.*, vol. 69, no. 1, pp. 619–633, Jan. 2021.

[29] B. Lyu, D. T. Hoang, S. Gong, D. Niyato, and D. I. Kim, "IRS-based wireless jamming attacks: When jammers can attack without power," *IEEE Wireless Commun. Lett.*, vol. 9, no. 10, pp. 1663–1667, Oct. 2020.

[30] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex optimization.* Cambridge, U.K.: Cambridge Univ. Press, 2004.

[31] Y. Zou, Y. Long, S. Gong, D. T. Hoang, W. Liu, W. Cheng, and D. Niyato, "Robust beamforming optimization for self-sustainable intelligent reflecting surface assisted wireless networks," *IEEE Trans. Cogn. Commun. Netw.*, pp. 1–1, Dec. 2021.

[32] Y. Zou, S. Gong, J. Xu, W. Cheng, D. T. Hoang, and D. Niyato, "Wireless powered intelligent reflecting surfaces for enhancing wireless communications," *IEEE Trans. Veh. Technol.*, vol. 69, no. 10, pp. 12 369–12 373, Oct. 2020.

[33] S. Gong, Y. Xie, J. Xu, D. Niyato, and Y.-C. Liang, "Deep reinforcement learning for backscatter-aided data offloading in mobile edge computing," *IEEE Netw.*, vol. 34, no. 5, pp. 106–113, Oct. 2020.

[34] N. C. Luong, D. T. Hoang, S. Gong, D. Niyato, P. Wang, Y.-C. Liang, and D. I. Kim, "Applications of deep reinforcement learning in communications and networking: A survey," *IEEE Commun. Surv. Tut.*, vol. 21, no. 4, pp. 3133–3174, May 2019.

[35] 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 *proc. Int. Conf. Learn. Represent. (ICLR)*, Jan. 2016.

[36] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *proc. Int. Conf. Mach. Learn. (ICML)*, Jul. 2015, pp. 1889–1897.

[37] T. Jiang, H. V. Cheng, and W. Yu, "Learning to reflect and to beamform for intelligent reflecting surface with implicit channel estimation," *IEEE J. Sel. Area. Commun.*, vol. 39, no. 7, pp. 1931–1945, Jul. 2021.

[38] I. Kadota, A. Sinha, E. Uysal-Biyikoglu, R. Singh, and E. Modiano, "Scheduling policies for minimizing age of information in broadcast wireless networks," *IEEE/ACM Trans. Netw.*, vol. 26, no. 6, pp. 2637–2650, Dec. 2018.