

Deep Reinforcement Learning for Wireless Scheduling in Distributed Networked Control

Gaoyang Pang, Kang Huang, Daniel E. Quevedo, *Fellow, IEEE*,
Branka Vucetic, *Fellow, IEEE*, Yonghui Li, *Fellow, IEEE*, Wanchun Liu, *Member, IEEE*

Abstract— We consider a joint uplink and downlink scheduling problem of a fully distributed wireless networked control system (WNCS) with a limited number of frequency channels. Using elements of stochastic systems theory, we derive a sufficient stability condition of the WNCS, which is stated in terms of both the control and communication system parameters. Once the condition is satisfied, there exists a stationary and deterministic scheduling policy that can stabilize all plants of the WNCS. By analyzing and representing the per-step cost function of the WNCS in terms of a finite-length countable vector state, we formulate the optimal transmission scheduling problem into a Markov decision process and develop a deep reinforcement learning (DRL) based framework for solving it. To tackle the challenges of a large action space in DRL, we propose novel action space reduction and action embedding methods for the DRL framework that can be applied to various algorithms, including Deep Q-Network (DQN), Deep Deterministic Policy Gradient (DDPG), and Twin Delayed Deep Deterministic Policy Gradient (TD3). Numerical results show that the proposed algorithm significantly outperforms benchmark policies.

Index Terms—Wireless networked control, transmission scheduling, deep Q-learning, Markov decision process, stability condition.

I. INTRODUCTION

The Fourth Industrial Revolution, Industry 4.0, is the automation of conventional manufacturing and industrial processes through flexible mass production. Automatic control in Industry 4.0 requires large-scale and interconnected deployment of massive spatially distributed industrial devices, such as sensors, actuators, machines, robots, and controllers. Eliminating the communication wires is a game-changer in reshaping traditional factories. In this regard, wireless networked control has become one of the most important technologies of Industry 4.0. It offers high-scalable and low-cost deployment capabilities and enables many industrial applications, such as smart city, smart manufacturing, smart grids, e-commerce warehouses and industrial automation systems [1], [2].

Unlike traditional cable-based networked control systems, in a large-scale wireless networked control system (WNCS), communication resources are limited. Not surprisingly, during the past decade, state estimation and transmission scheduling in WNCS have drawn a lot of attention in the research community. For example, some recent works have investigated the estimation problem under dynamic sampling periods [3] and random access protocol scheduling [4]. Focusing on the sensor-

controller communications only (the controller-actuator co-located scenario), the optimal transmission scheduling problem of over a single frequency channel for achieving the best remote estimation quality was extensively investigated in [5]–[8]. For the multi-frequency channel scenario, the optimal scheduling policy and structural results were obtained in [9]. The optimal transmission power scheduling problem of an energy-constrained remote estimation system was studied in [10]. The joint scheduling and power allocation problems of multi-plant-multi-frequency WNCSs were investigated in [11], [12] for achieving the minimum overall transmission power consumption. From a cyber-physical security perspective, denial-of-service (DoS) attackers' scheduling problems were investigated in [13], [14], for maximally deteriorating the WNCSs' performance.

In the WNCS examined in the above works, the scheduling problem is commonly reformulated into a Markov decision process (MDP) problem. The MDP problem can be solved using traditional algorithms (e.g., value iteration [2], [9], [15], policy iteration [16], and linear programming [17], [18]) and reinforcement learning algorithms (e.g., Q-learning [19] and deep Q-Networks [20], [21]). Traditional algorithms and Q-learning are effective when the system scale is small. For example, numerical results of the optimal scheduling of a two-sensor-one-frequency using value iteration system were presented in [9]. These methods involve operations over every possible state and action. As the number of actions increases, the number of state-action pairs grows, leading to a dramatic increase in the computations required. For example, when the action space is large, Q-learning necessitates constructing a Q-value table, which demands substantial storage proportional to the size of the action space, leading to the curse of dimensionality [19], [22], [23]. Deep Q-Networks (DQN) address this issue using deep neural networks (DNNs) for value function approximations, thus eliminating the need for extensive Q-value tables. Some recent works [20], [21] have applied DQN to solve multi-system-multi-frequency scheduling problems in different WNCS scenarios.

However, handling large action spaces also presents significant challenges in DQN. The size of the DNNs required for DQN expands dramatically with respect to the action space. This can hinder effective exploration of the action space during training, resulting in increased storage and computational requirements. Ma et al. [24] proposed decoupling the original problem into subproblems and developing multiple deep reinforcement learning (DRL) algorithms to solve the corresponding subproblems with reduced decision space. However, the presence of multiple learning agents creates a

G. Pang and K. Huang contributed equally to the work. G. Pang, B. Vucetic, Y. Li, and W. Liu are with School of Electrical and Computer Engineering, The University of Sydney, Australia. Emails: {gaoyang.pang; branka.vucetic; yonghui.li; wanchun.liu}@sydney.edu.au. K. Huang is with Huawei Shanghai Research Center, Shanghai, China. Email: huangkang9@huawei.com. D. E. Quevedo is with the School of Electrical Engineering and Robotics, Queensland University of Technology (QUT), Brisbane, Australia. Email: dquevedo@ieec.org.

non-stationary environment, which can destabilize learning and make convergence more difficult. To deal with large decision spaces, recent works resort to stochastic policy-based DRL algorithms [25]–[27]. A stochastic policy refers to a policy that specifies a probability distribution over actions, given the current state of the environment. Ni et al. [25] leveraged a Graph Neural Network (GNN) based DRL to learn stochastic policies, i.e., Proximal Policy Optimization (PPO) for warehouse scheduling. Yang et al. [26] developed an actor-critic DRL for learning stochastic policies with continuous action space for scheduling, power allocation, and modulation scheme adaptation. However, it has been proved that the optimal policy is deterministic for unconstrained MDP problems [28]. Here, a deterministic policy means that the action in a given state is determined by a fixed function without any randomness. Consequently, deterministic policy-based DRL methods are generally more preferable.

We note that most of the existing works only focus on WNCSs which are only partially distributed [5]–[7], [9], [11], [12], [20], [21]. For example, Chen et al. [29] focus on the sensor scheduling in the uplink transmission of a WNCS, using either DQN or Deep Deterministic Policy Gradient (DDPG). Redder et al. [30] developed a DQN-based algorithm for optimal actuator scheduling in the downlink transmission of a WNCS over a Markov fading channel. Existing works only consider optimizing either the uplink or downlink transmission for WNCSs. In a fully distributed setting, both downlink (controller-actuator) and uplink (sensor-controller) transmissions are crucial for stabilizing each plant. The joint uplink and downlink scheduling problem of fully distributed multi-plant WNCS has never been considered in the open literature.

In this paper, we investigate the transmission scheduling problem of distributed WNCS. The main contributions are summarized as follows:

- We propose a distributed N -plant- M -frequency WNCS model, where the controller schedules the uplink and downlink transmissions of all N plants and the spatial diversity of different communication links are taken into account. The controller generates sequential predictive control commands for each of the plants based on pre-designed deadbeat control laws.² Different from uplink or downlink only scheduling, a joint scheduling algorithm has a larger action space and also needs to automatically balance the tradeoff between the uplink and the downlink due to the communication resource limit. To the best of our knowledge, joint uplink and downlink transmission scheduling of distributed WNCSs has not been investigated in the open literature.
- We derive a sufficient stability condition of the WNCS in terms of both the control and communication system parameters. The result provides a theoretical guarantee that there exists at least one stationary and deterministic scheduling policy that can stabilize all plants of the

WNCS. We show that the obtained condition is also necessary in the absence of spatial diversity of different communication links.

- We construct a finite-length countable vector state of the WNCS in terms of the time duration between consecutive received packets at the controller and the actuators. Then, we prove that the per-step cost function of the WNCS is determined by the time-duration-related vector state. Building on this, we formulate the optimal transmission scheduling problem into a Markov decision process (MDP) problem with an countable state space for achieving the minimum expected total discounted cost. We propose effective action space reduction and action embedding methods that can be applied to DRL algorithms, including foundational ones like DQN and more advanced ones such as DDPG and Twin Delayed Deep Deterministic Policy Gradient (TD3), for solving the problem. Numerical results illustrate that the proposed algorithm can reduce the expected cost significantly compared to available benchmark policies.

Notations: $\sum_{m=i}^j a_m \triangleq 0$ if $i > j$. $\limsup_{K \rightarrow \infty}$ is the limit superior operator. $C_n^k \triangleq \frac{n!}{k!(n-k)!}$ and $P_n^k \triangleq \frac{n!}{(n-k)!}$ are the numbers of combinations and permutations of n things taken k at a time, respectively. $\text{Tr}(\mathbf{A})$ and $\text{rank}[\mathbf{A}]$ denote the trace and the rank of matrix \mathbf{A} , respectively.

II. DISTRIBUTED WNCS WITH SHARED WIRELESS RESOURCE

We consider a distributed WNCS system with N plants and a central controller as illustrated in Fig. 1. The output of plant i is measured by smart sensor i , which sends pre-filtered measurements (local state estimates) to the controller. The controller applies a remote (state) estimator and a control algorithm for plant i . It then generates and sends a control signal to actuator i , thereby closing the loop. A key limitation of the WNCS is that smart sensor i cannot communicate with actuator i directly due to the spatial deployment. Instead, the uplink (sensor-controller) and downlink (controller-actuator) communications for the N plants share a common wireless network with only M frequency channels, where $M < 2N$. Thus, not every node is allowed to transmit at the same time and communications need to be scheduled. As shown in Fig. 1, we shall focus on a setup where scheduling is done at the controller side, which schedules both the downlink and uplink transmissions.³

Each smart sensor has three modules: remote estimator, control algorithm and local Kalman filter. The first two modules are copies of the ones at the controller to reconstruct the control input of the plant. Such information is sent to the Kalman filter for accurate local state estimation. The details will be presented in the sequel.

²Note that deadbeat controller is commonly considered as a time-optimal controller that takes the minimum time for setting the current plant state to the origin.

³Such a setup is practical, noting that most of the existing wireless communication systems including 5G support both uplink and downlink at base stations [31].

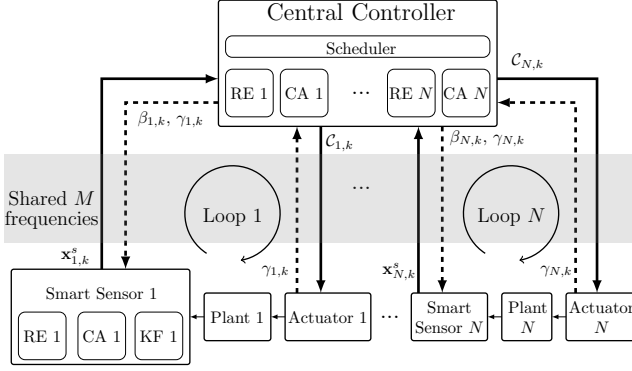


Fig. 1. A distributed networked control system with N plants sharing M frequency channels. Kalman filter, remote estimator and control algorithm are denoted as KF, RE and CA and discussed in Sections II-B, II-C and II-D, respectively.

A. Plant Dynamics

The N plants are modeled as linear time-invariant (LTI) discrete-time systems as [11], [21]

$$\begin{aligned} \mathbf{x}_{i,k+1} &= \mathbf{A}_i \mathbf{x}_{i,k} + \mathbf{B}_i \mathbf{u}_{i,k} + \mathbf{w}_{i,k}, \\ \mathbf{y}_{i,k} &= \mathbf{C}_i \mathbf{x}_{i,k} + \mathbf{v}_{i,k}, \quad i = 1, \dots, N \end{aligned} \quad (1)$$

where $\mathbf{x}_{i,k} \in \mathbb{R}^{n_i}$ and $\mathbf{u}_{i,k} \in \mathbb{R}^{m_i}$ are the state vector of plant i and the control input applied by actuator i at time k , respectively. $\mathbf{w}_{i,k} \in \mathbb{R}^{n_i}$ is the i -th plant disturbance and is an independent and identically distributed (i.i.d.) zero-mean Gaussian white noise process with covariance matrix $\mathbf{Q}_i^w \in \mathbb{R}^{n_i \times n_i}$. $\mathbf{A}_i \in \mathbb{R}^{n_i \times n_i}$ and $\mathbf{B}_i \in \mathbb{R}^{n_i \times m_i}$ are the system-transition matrix and control-input matrix for plant i , respectively. $\mathbf{y}_{i,k} \in \mathbb{R}^{p_i}$ is sensor i 's measurement of plant i at time k and $\mathbf{v}_{i,k} \in \mathbb{R}^{p_i}$ is the measurement noise, modeled as an i.i.d. zero-mean Gaussian white noise process with covariance matrix $\mathbf{Q}_i^v \in \mathbb{R}^{p_i \times p_i}$. $\mathbf{C}_i \in \mathbb{R}^{p_i \times n_i}$ is the measurement matrix of plant i . We assume that plant i is v_i -step controllable, $\forall i \in \{1, \dots, N\}$ [32], i.e., there exists a control gain $\tilde{\mathbf{K}}_i$ satisfying

$$(\mathbf{A}_i + \mathbf{B}_i \tilde{\mathbf{K}}_i)^{v_i} = \mathbf{0}. \quad (2)$$

Note that when \mathbf{A}_i is non-singular, the system (1) is v_i -step controllable if and only if

$$\text{rank} [\mathbf{A}_i^{-1} \mathbf{B}_i, \mathbf{A}_i^{-2} \mathbf{B}_i, \dots, \mathbf{A}_i^{-v_i} \mathbf{B}_i] = n_i.$$

The details of the control algorithm will be given in Section II-D.

B. Smart Sensors

Due to the noise measurement, each smart sensor runs a Kalman filter to estimate the current plant state as below [33]:⁴

⁴In this subsection, we focus on smart sensor i and the index i of each quantity is omitted for clarity

$$\begin{aligned} \mathbf{x}_{k|k-1}^s &= \mathbf{A} \mathbf{x}_{k-1}^s + \mathbf{B} \mathbf{u}_{k-1} \\ \mathbf{P}_{k|k-1}^s &= \mathbf{A} \mathbf{P}_{k-1}^s \mathbf{A}^\top + \mathbf{Q}^w \\ \mathbf{K}_k &= \mathbf{P}_{k|k-1}^s \mathbf{C}^\top (\mathbf{C} \mathbf{P}_{k|k-1}^s \mathbf{C}^\top + \mathbf{Q}^v)^{-1} \\ \mathbf{x}_k^s &= \mathbf{x}_{k|k-1}^s + \mathbf{K}_k (\mathbf{y}_k - \mathbf{C} \mathbf{x}_{k|k-1}^s) \\ \mathbf{P}_k^s &= (\mathbf{I} - \mathbf{K}_k \mathbf{C}) \mathbf{P}_{k|k-1}^s \end{aligned} \quad (3)$$

where $\mathbf{x}_{k|k-1}^s$ and \mathbf{x}_k^s are the prior and posterior state estimation at time k , respectively, and $\mathbf{P}_{k|k-1}^s$ and \mathbf{P}_k^s are the (estimation error) covariances of $\mathbf{e}_{k|k-1}^s \triangleq \mathbf{x}_k - \mathbf{x}_{k|k-1}^s$ and $\mathbf{e}_k^s \triangleq \mathbf{x}_k - \mathbf{x}_k^s$, respectively. \mathbf{K}_k is the Kalman gain at time k . We assume that each (\mathbf{A}, \mathbf{C}) is observable and $(\mathbf{A}, \mathbf{Q}^w)$ is controllable [9]. Thus, the Kalman gain \mathbf{K}_k and error covariance matrix \mathbf{P}_k^s converge to constant matrices $\hat{\mathbf{K}}$ and $\hat{\mathbf{P}}^s$, respectively, i.e., the smart sensor is in the stationary mode. As foreshadowed, the smart sensor employs a remote estimator and knowledge of the control algorithm that is applied at the controller side to obtain the control input in (3). Thus, the Kalman filter (3) is the optimal estimator of the linear system (1) in terms of the estimation mean-square error [34]. The remote estimator and control algorithm will be presented in Sections II-C and II-D, respectively.

As shown in Fig. 1, at every scheduling instant k , the smart sensor sends the local estimate \mathbf{x}_k^s (rather than the raw measurement \mathbf{y}_k) to the controller.

Before proceeding, we note that (3) leads to the following recursion for the estimation error at the sensors:

$$\begin{aligned} \mathbf{e}_{k+1|k}^s &= \mathbf{A} \mathbf{e}_k^s + \mathbf{w}_k \\ \mathbf{e}_k^s &= (\mathbf{I} - \hat{\mathbf{K}} \mathbf{C}) \mathbf{e}_{k|k-1}^s - \hat{\mathbf{K}} \mathbf{v}_k \end{aligned}$$

and hence

$$\mathbf{e}_k^s = (\mathbf{I} - \hat{\mathbf{K}} \mathbf{C}) \mathbf{A} \mathbf{e}_{k-1}^s + (\mathbf{I} - \hat{\mathbf{K}} \mathbf{C}) \mathbf{w}_{k-1} - \hat{\mathbf{K}} \mathbf{v}_k.$$

Then, the relation between the estimation errors \mathbf{e}_k^s and \mathbf{e}_{k-K}^s , $\forall K \in \mathbb{N}$, is established as

$$\mathbf{e}_k^s = \mathbf{Z}^K \mathbf{e}_{k-K}^s + \sum_{i=1}^K \mathbf{Z}^{i-1} (\mathbf{I} - \hat{\mathbf{K}} \mathbf{C}) \mathbf{w}_{k-i} - \sum_{i=1}^K \mathbf{Z}^{i-1} \hat{\mathbf{K}} \mathbf{v}_{k-i+1},$$

where $\mathbf{Z} \triangleq (\mathbf{I} - \hat{\mathbf{K}} \mathbf{C}) \mathbf{A}$.

C. Remote Estimation

At the beginning of a time slot, the controller calculates the control sequence and transmits it within the time slot based on the current plant state estimation rather than the real-time state, as the packet carrying the current plant state can only be received at the end of the current time slot due to the one-step transmission delay. Then, the controller applies a minimum mean-square error (MMSE) remote estimator for each plant taking into account the random packet dropouts and one-step transmission delay as [33]:

$$\hat{\mathbf{x}}_{i,k+1} = \begin{cases} \mathbf{A}_i \mathbf{x}_{i,k}^s + \mathbf{B}_i \mathbf{u}_{i,k} & \text{if } \beta_{i,k} = 1 \\ \mathbf{A}_i \hat{\mathbf{x}}_{i,k} + \mathbf{B}_i \mathbf{u}_{i,k} & \text{if } \beta_{i,k} = 0 \end{cases} \quad (4)$$

where $\beta_{i,k} = 1$ or 0 indicates that the controller receives sensor i 's packet or not at time k , respectively. Then, the estimation error $\mathbf{e}_{i,k}$ is obtained as

$$\mathbf{e}_{i,k} \triangleq \mathbf{x}_{i,k} - \hat{\mathbf{x}}_{i,k} = \begin{cases} \mathbf{A}_i \mathbf{e}_{i,k-1}^s + \mathbf{w}_{i,k-1} & \text{if } \beta_{i,k-1} = 1 \\ \mathbf{A}_i \mathbf{e}_{i,k-1} + \mathbf{w}_{i,k-1} & \text{if } \beta_{i,k-1} = 0 \end{cases} \quad (5)$$

From (4), the controller's current estimation depends on the most recently received sensor estimation. Let $\tau_{i,k} \in \{1, 2, \dots\}$ denote the *age-of-information (AoI)* (see [15] and reference therein) of sensor i 's packet observed at time k , i.e., the number of elapsed time slots from the latest successfully delivered sensor i 's packet before the current time k , which reflects how old the most recently received sensor measurement is. Then, it is easy to see that the updating rule of $\tau_{i,k}$ is given by

$$\tau_{i,k+1} = \begin{cases} 1 & \text{if } \beta_{i,k} = 1 \\ \tau_{i,k} + 1 & \text{otherwise.} \end{cases} \quad (6)$$

Using (5) and the AoI, the relation between the local and remote estimation error can be characterized by:

$$\mathbf{e}_{i,k} = \mathbf{A}_i^{\tau_{i,k}} \mathbf{e}_{i,k-\tau_{i,k}}^s + \sum_{j=1}^{\tau_{i,k}} \mathbf{A}_i^{j-1} \mathbf{w}_{k-j}.$$

D. Control Algorithm

Due to the fact that downlink transmissions are unreliable, actuator i may not receive the controller's control-command-carrying packets, even when transmissions are scheduled. We adopt a predictive control approach [32], [35] to provide robustness against packet failures: the controller sends a length- v_i sequence of control commands including both the current command and the predicted future commands to the actuator once scheduled; if the current packet is lost, the actuator will apply the previously received predictive command for the current time slot as the control input.

The control sequence for plant i is generated by a linear deadbeat control gain $\tilde{\mathbf{K}}_i$ as [32]

$$\mathcal{C}_{i,k} = [\tilde{\mathbf{K}}_i \hat{\mathbf{x}}_{i,k}, \tilde{\mathbf{K}}_i \Phi_i \hat{\mathbf{x}}_{i,k}, \dots, \tilde{\mathbf{K}}_i (\Phi_i)^{v_i-1} \hat{\mathbf{x}}_{i,k}] \quad (7)$$

where $\Phi_i \triangleq \mathbf{A}_i + \mathbf{B}_i \tilde{\mathbf{K}}_i$ and $\tilde{\mathbf{K}}_i$ satisfies $(\mathbf{A}_i + \mathbf{B}_i \tilde{\mathbf{K}}_i)^{v_i} = \mathbf{0}$, and v_i is the controllability index of the pair $(\mathbf{A}_i, \mathbf{B}_i)$. Note that the first element in $\mathcal{C}_{i,k}$ is the current control command and the rest are the predicted ones.

Remark 1. *It can be verified that if the current state estimation $\hat{\mathbf{x}}_{i,k}$ is perfect and the plant i is disturbance free, then the plant state $\mathbf{x}_{i,k}$ would be set to zero after applying all v_i steps of the control sequence in (7). Such a deadbeat controller is commonly considered as a time-optimal controller that takes the minimum time for setting the current plant state to the origin [36]. We note that the deadbeat control law may not be cost-optimal optimal to minimize the quadratic cost function defined in Section IV, and the optimal control law may depend on the scheduling policy. Since the current work focuses on the transmission scheduling design of the N -plant- M -frequency WNCs, the optimal joint control-scheduling problem can be*

investigated in our future work. Specifically, we aim to identify an optimal scheduling policy that minimizes the quadratic cost for a given deadbeat controller.

Accordingly, actuator i maintains a length- v_i buffer

$$\mathcal{U}_{i,k} \triangleq [\mathbf{u}_{i,k}^0, \mathbf{u}_{i,k}^1, \dots, \mathbf{u}_{i,k}^{v_i-1}]$$

to store the received control commands. If the current control packet is received, the buffer is reset with received sequence; otherwise, it is shifted one step forward as

$$\mathcal{U}_{i,k} = \begin{cases} \mathcal{C}_{i,k}, & \text{if } \gamma_{i,k} = 1 \\ [\mathbf{u}_{i,k-1}^1, \mathbf{u}_{i,k-1}^2, \dots, \mathbf{u}_{i,k-1}^{v_i}, \mathbf{0}], & \text{if } \gamma_{i,k} = 0 \end{cases} \quad (8)$$

where $\gamma_{i,k} = 1$ or 0 indicate that actuator i receives a control packet or not at time k , respectively. The first command in the buffer is applied as the control input each time

$$\mathbf{u}_{i,k} \triangleq \mathbf{u}_{i,k}^0.$$

Let $\eta_{i,k} \in \{1, 2, \dots\}$ denote the AoI of controller's packet at actuator i observed at time k , i.e., the number of passed time slots (including the current time slot) since the latest received control packet by actuator i . Its updating rule is given as

$$\eta_{i,k} = \begin{cases} 1, & \text{if } \gamma_{i,k} = 1 \\ \eta_{i,k-1} + 1 & \text{if } \gamma_{i,k} = 0 \end{cases} \quad (9)$$

From (7) and (8), and by using the deadbeat control property (2), the applied control input can be concisely written as

$$\mathbf{u}_{i,k} = \tilde{\mathbf{K}}_i (\Phi_i)^{\eta_{i,k}-1} \hat{\mathbf{x}}_{i,k+1-\eta_{i,k}}. \quad (10)$$

E. Communication Scheduler

The N -plant WNCs has N uplinks and N downlinks sharing M frequencies. Each frequency can be occupied by at most one link, and each link can be allocated to at most one frequency at a time. Let $a_{m,k} \in \{-N, \dots, 0, \dots, N\}$ denote the allocated link to frequency m at time k , where $a_{m,k} = i'$ means the frequency is allocated to $|i'|$ -th plant, where $i' > 0$ and < 0 indicate for uplink and downlink, respectively, and $i' = 0$ denotes that the frequency channel is idle.

The packet transmissions of each link are modeled as i.i.d. packet dropout processes. Unlike most of the existing works wherein transmission scheduling of WNCs assumes that different node transmissions at the same frequency channel have the same packet drop probability [20], we here consider a more practical scenario by taking into account the spatial diversity of different transmission nodes – each frequency has different dropout probabilities for different uplink and downlink transmissions. The packet success probabilities of the uplink and downlink of plant i at frequency m are given by $\xi_{m,i}^s$ and $\xi_{m,i}^c$, respectively, where

$$\begin{aligned} \xi_{m,i}^s &\triangleq \mathbb{P}[\beta_{i,k} = 1 | a_{m,k} = i], \\ \xi_{m,i}^c &\triangleq \mathbb{P}[\gamma_{i,k} = 1 | a_{m,k} = -i]. \end{aligned} \quad (11)$$

The packet success probabilities can be estimated by the controller based on standard channel estimation techniques [37], [38], and are utilized by the MDP and DRL-based solutions in Section IV-C and Section V.

Acknowledgment feedback. We note that the Transmission Control Protocol (TCP) is commonly adopted in commercial telecommunication systems for data transmission, such as 5G and WiFi [39]–[41]. TCP relies on acknowledgment feedback, where the receiver sends a one-bit acknowledgment signal back to the sender to confirm the successful receipt of data packets. This feedback mechanism ensures reliable data transmission by verifying the delivery of data. Therefore, we assume that both the uplink and downlink transmissions adopt the acknowledgment feedback scheme. In particular, the actuator sends one-bit feedback signal of $\gamma_{i,k}$ to the controller, and the controller sends $\beta_{i,k}$ as well as the received $\gamma_{i,k}$ to sensor i each time with negligible overhead. From $\{\beta_{i,k}\}$ and (4), the smart sensor knows the estimated plant state $\hat{\mathbf{x}}_{i,k}$ by the controller; then, using $\hat{\mathbf{x}}_{i,k}$, $\{\gamma_{i,k}\}$ and (10), the sensor can calculate the applied control input $\mathbf{u}_{i,k}$, which is utilized for local state estimation as mentioned in (3). Note that sending one-bit acknowledgement feedback leads to negligible overhead in contrast to sending the applied control input $\mathbf{u}_{i,k}$ from the controller to the sensors.

III. STABILITY CONDITION

Before turning to designing scheduling policies, it is critical to elucidate conditions that the WNCS needs to satisfy which ensure that there exists at least one stationary and deterministic scheduling policy that can stabilize all plants using the available network resources. Note that a policy π is a function mapping from a state to an action. A **deterministic policy** means that the function gives the same action when the input state is fixed. A **stationary policy** means the function π is time invariant, i.e., $\pi_k = \pi, \forall k$. In other words, the action only depends on the state, not the time [28]. We adopt a very commonly considered stochastic stability condition as below.

Assumption 1. *The expected initial quadratic norm of each plant state is bounded, i.e., $\mathbb{E}[\mathbf{x}_{i,0}^\top \mathbf{x}_{i,0}] < \infty, \forall i = 1, \dots, N$.*

Definition 1 (Mean-Square Stability). *The WNCS is mean-square stable under Assumption 1 if and only if*

$$\limsup_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\mathbf{x}_{i,k}^\top \mathbf{x}_{i,k}] < \infty, \forall i = 1, \dots, N. \quad (12)$$

Intuitively, the stability condition of the WNCS should depend on both the (open-loop) unstable plant systems (i.e., those where $\rho(\mathbf{A}_i) \geq 1$) and the M -frequency communication system parameters. For the tractability of sufficient stability condition analysis (i.e., to prove the existence of a stabilizing policy), we focus on policy class that groups the unstable plants into M disjoint sets $\mathcal{F}_1, \dots, \mathcal{F}_M$ and allocates them to the M frequencies, accordingly. Let $\bar{\mathcal{F}} \triangleq \{i : \rho(\mathbf{A}_i) \geq 1, i \in \{1, \dots, N\}\}$ denote the index set of all unstable plants. We have $\mathcal{F}_1 \cup \dots \cup \mathcal{F}_M = \bar{\mathcal{F}} \subseteq \{1, \dots, N\}$ and $\mathcal{F}_i \cap \mathcal{F}_j = \emptyset, \forall i \neq$

j . Then, we present the stability condition below which takes into account all potential allocations $(\mathcal{F}_1, \dots, \mathcal{F}_M)$.

Theorem 1 (Stabilizability). *Consider the index set $\{\mathcal{F}_m\}$ as introduced above and define $\rho_m^{\max} \triangleq \max_{i \in \mathcal{F}_m} \rho^2(\mathbf{A}_i)$, $\bar{\xi}_m^{\max} \triangleq \max_{i \in \mathcal{F}_m} \{\bar{\xi}_{m,i}^s, \bar{\xi}_{m,i}^c\}$, $\bar{\xi}_{m,i}^s \triangleq 1 - \xi_{m,i}^s$, $\bar{\xi}_{m,i}^c \triangleq 1 - \xi_{m,i}^c$. We then have:*

(a) *A sufficient condition under which the WNCS described by (1), (3), (4), (10) and (11) has a stationary and deterministic scheduling policy satisfying the stability condition (12) is given by*

$$\kappa \triangleq \min_{(\mathcal{F}_1, \dots, \mathcal{F}_M)} \max_{m=1, \dots, M, \mathcal{F}_m \neq \emptyset} \rho_m^{\max} \bar{\xi}_m^{\max} < 1. \quad (13)$$

(b) *For the special case that the packet error probabilities of different links are identical at the same frequency (i.e., where no spatial diversity exists), $\bar{\xi}_{m,1}^s = \bar{\xi}_{m,1}^c = \dots = \bar{\xi}_{m,N}^s = \bar{\xi}_{m,N}^c, \forall m = 1, \dots, M$, the condition (13) is necessary and sufficient.*

Proof. The sufficient condition (13) is derived in two steps: 1) the construction of a stationary and deterministic scheduling policy and 2) the proof of the condition under which the constructed policy leads to a bounded average cost. Since a plant with $\rho(\mathbf{A}_i) < 1$ does not need any communication resources for stabilization, in the following, we only focus on the unstable plants in $\bar{\mathcal{F}} \neq \emptyset$.

We construct a multi-frequency persistent scheduling policy: the unstable plants are grouped into M sets $\mathcal{F}_1, \dots, \mathcal{F}_M$, corresponding to the M frequencies. In each frequency, the controller schedules the uplink of the first plant, say plant i , persistently until success. It then persistently schedules the downlink until success, and waits for $v_i - 1$ steps for applying all the control commands in the actuator's buffer. It then schedules the uplink of the second plant and so on and so forth. This procedure is repeated ad-infinitum. The reason we choose such a policy is because of its tractability and the tightness of the sufficient stability condition that we will derive. The detailed proof is included in [42]. \square

Example 1. *Consider a WNCS with $N = 3$ and $M = 2$, and $\rho^2(\mathbf{A}_i) \geq 1, i = 1, 2, 3$. The packet error probabilities are $\bar{\xi}_{1,1}^s = 0.1, \bar{\xi}_{1,1}^c = 0.3, \bar{\xi}_{1,2}^s = 0.2, \bar{\xi}_{1,2}^c = 0.1, \bar{\xi}_{1,3}^s = 0.2, \bar{\xi}_{1,3}^c = 0.4, \bar{\xi}_{2,1}^s = 0.1, \bar{\xi}_{2,1}^c = 0.3, \bar{\xi}_{2,2}^s = 0.2, \bar{\xi}_{2,2}^c = 0.1, \bar{\xi}_{2,3}^s = 0.2, \bar{\xi}_{2,3}^c = 0.4$. There are eight plant-grouping schemes at the two frequencies $(\mathcal{F}_1, \mathcal{F}_2)$, i.e., $(\{1, 2, 3\}, \emptyset), (\{1, 2\}, \{3\}), (\{1\}, \{2, 3\}), (\{2, 3\}, \{1\}), (\{2\}, \{1, 3\}), (\{3\}, \{1, 2\}), (\{1, 3\}, \{2\})$ and $(\emptyset, \{1, 2, 3\})$. If $\rho^2(\mathbf{A}_1) = 3, \rho^2(\mathbf{A}_2) = 2$ and $\rho^2(\mathbf{A}_3) = 1$, the stability condition is satisfied as $\kappa = 0.9 < 1$; if $\rho^2(\mathbf{A}_1) = 1, \rho^2(\mathbf{A}_2) = 2$ and $\rho^2(\mathbf{A}_3) = 3$, the condition is unsatisfied as $\kappa = 1.2 > 1$.*

Remark 2. *Theorem 1 captures the stabilizability of the WNCS scheduling problem in terms of both the dynamic system parameters, $\mathbf{A}_i, \forall i \in \bar{\mathcal{F}}$, and the wireless channel conditions, i.e., $\{\bar{\xi}_{m,i}^s, \bar{\xi}_{m,i}^c\}, i \in \bar{\mathcal{F}}, m = 1, \dots, M$. Once (13) holds, there is at least one stationary and deterministic policy that stabilized all plants of the WNCS. If the channel quality of*

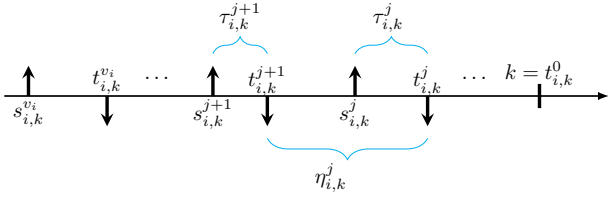


Fig. 2. Illustration of the state parameters of plant i .

different links does not differ much at the same frequency, then the sufficient stabilizability condition is tight. To the best of our knowledge, this is the first stabilizability condition established for N -plant- M -frequency WNCS with uplink and downlink scheduling in the literature.

IV. ANALYSIS AND MDP DESIGN

As a performance measure of the WNCS, we consider the expected (infinite horizon) total discounted cost (ETDC) given by

$$J = \sum_{k=0}^{\infty} \vartheta^k \sum_{i=1}^N \mathbb{E}[\mathbf{x}_{i,k}^{\top} \mathbf{S}_i^x \mathbf{x}_{i,k} + \mathbf{u}_{i,k}^{\top} \mathbf{S}_i^u \mathbf{u}_{i,k}], \quad (14)$$

where $\vartheta \in (0, 1)$ is the discount factor, and a smaller ϑ means the future cost is less important. \mathbf{S}_i^x and \mathbf{S}_i^u are positive definite weighting matrices for the system state and control input of plant i , respectively. Thus, it is important to find a scheduling policy that can minimize the design objective (14).

Note that an ETDC minimization problem is commonly obtained by reformulating it into an MDP and solving it by classical policy and value iteration methods [28]. Theoretically speaking, the MDP solution provides an optimal deterministic and stationary policy, which is a mapping between its state and the scheduling action at each time step. However, the optimal MDP solution is intractable due to the uncertainties involved and the curse of dimensionality. Thus, we seek to find a good approximate MDP solution by DRL. In the following, we aim to formulate the scheduler design problem into an MDP, and then present a DRL solution in Section V.

From (14), the per-step cost of the WNCS depends on state $\mathbf{x}_{i,k}$ and $\mathbf{u}_{i,k}$, which have continuous (uncountable) state spaces. Furthermore, $\mathbf{x}_{i,k}$ is not observable by the controller. To design a suitable MDP problem with a countable state space, we first need to determine an observable, discrete state of the MDP (in Section IV-C) and investigate how to represent the per-step cost function in (14), i.e., $\sum_{i=1}^N \mathbb{E}[\mathbf{x}_{i,k}^{\top} \mathbf{S}_i^x \mathbf{x}_{i,k} + \mathbf{u}_{i,k}^{\top} \mathbf{S}_i^u \mathbf{u}_{i,k}]$, in terms of the state.

A. MDP State Definition

We introduce event-related time parameters in the following. Let $t_{i,k}^j, j = 1, 2, \dots, v_i + 1$, denote the time index of the j -th latest successful packet reception at actuator i prior to the current time slot k , and $t_{i,k}^0 \triangleq k$. Let $s_{i,k}^j, j = 0, \dots, v_i$, denote the time index of the latest successful sensor i 's transmission prior to $t_{i,k}^j$, where $s_{i,k}^j \leq t_{i,k}^j$, as illustrated in Fig 2.

Then, we define a sequence of variables, $\tau_{i,k}^j, j = 0, \dots, v_i$, as

$$\tau_{i,k}^j \triangleq t_{i,k}^j - s_{i,k}^j, \quad (15)$$

to record the estimation quality at k and at the v_i successful control transmissions, where $\tau_{i,k}^0 \triangleq \tau_{i,k}$ was defined above (6). Similarly, we define

$$\eta_{i,k}^j \triangleq t_{i,k}^j - t_{i,k}^{j+1}, \quad j = 0, \dots, v_i, \quad (16)$$

denoting the time duration between consecutive successful controller's transmissions, where $\eta_{i,k}^0 \triangleq \eta_{i,k}$ was defined above (9). From (15) and (16), $\tau_{i,k}^j$ and $\eta_{i,k}^j$ can be treated as the AoI of the sensor's and the controller's packet of plant i , respectively, observed at $t_{i,k}^j$.

From (15), (6), (16), and (9), the updating rules of $\tau_{i,k}^j$ and $\eta_{i,k}^j$ can be obtained as

$$\tau_{i,k+1}^j = \begin{cases} 1, & \text{if } \beta_{i,k} = 1 \\ \tau_{i,k}^j + 1, & \text{if } \beta_{i,k} = 0 \end{cases} \text{ for } j = 0 \\ \tau_{i,k+1}^j = \begin{cases} \tau_{i,k}^{j-1}, & \text{if } \gamma_{i,k+1} = 1 \\ \tau_{i,k}^j, & \text{if } \gamma_{i,k+1} = 0 \end{cases} \text{ for } j = 1, \dots, v_i \quad (17)$$

$$\eta_{i,k}^j = \begin{cases} 1, & \text{if } \gamma_{i,k} = 1 \\ \eta_{i,k-1}^j + 1, & \text{if } \gamma_{i,k} = 0 \end{cases} \text{ for } j = 0 \\ \eta_{i,k}^j = \begin{cases} \eta_{i,k-1}^{j-1}, & \text{if } \gamma_{i,k} = 1 \\ \eta_{i,k-1}^j, & \text{if } \gamma_{i,k} = 0 \end{cases} \text{ for } j = 1, \dots, v_i \quad (18)$$

Now we define the AoI-related vector state of plant i as

$$\mathbf{s}_{i,k} \triangleq (\tau_{i,k}^0, \dots, \tau_{i,k}^{v_i}, \eta_{i,k}^0, \dots, \eta_{i,k}^{v_i}). \quad (19)$$

B. MDP Cost Function

We will show that the per-step cost function of plant i in (14), i.e., $\mathbb{E}[\mathbf{x}_{i,k}^{\top} \mathbf{S}_i^x \mathbf{x}_{i,k} + \mathbf{u}_{i,k}^{\top} \mathbf{S}_i^u \mathbf{u}_{i,k}]$, is determined by the vector state $\mathbf{s}_{i,k}$. We focus on plant i (the analytical method is identical for the other plants), and shall omit the plant index subscript of $\mathbf{x}_{i,k}$, $\mathbf{s}_{i,k}$, $\mathbf{u}_{i,k}$, \mathbf{S}_i^x and \mathbf{S}_i^u in the remainder of this subsection, where the corresponding parameters are \mathbf{x}_k , \mathbf{s}_k , \mathbf{u}_k , \mathbf{S}^x and \mathbf{S}^u .

Taking (10) into (1), the plant state evolution can be rewritten as

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\tilde{\mathbf{K}}(\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0 - 1} \hat{\mathbf{x}}_{k+1 - \eta_k^0} + \mathbf{w}_k.$$

By using this backward iteration for $k - t_k^v$ times and the definition of η_k^j and t_k^j , we have

$$\begin{aligned}
\mathbf{x}_k &= (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0} \mathbf{x}_{t_k^1} + (\mathbf{A}^{\eta_k^0} - (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0}) \mathbf{e}_{t_k^1} + \sum_{i=1}^{\eta_k^0} \mathbf{A}^{i-1} \mathbf{w}_{k-i} \\
&= (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0} \times \\
&\quad \left((\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^1} \mathbf{x}_{t_k^2} + (\mathbf{A}^{\eta_k^1} - (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^1}) \mathbf{e}_{t_k^2} + \sum_{i=1}^{\eta_k^1} \mathbf{A}^{i-1} \mathbf{w}_{t_k^1-i} \right) \\
&\quad + (\mathbf{A}^{\eta_k^0} - (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0}) \mathbf{e}_{t_k^1} + \sum_{i=1}^{\eta_k^0} \mathbf{A}^{i-1} \mathbf{w}_{k-i} \\
&= \dots \\
&= (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0 + \eta_k^1 + \dots + \eta_k^{v-1}} \mathbf{x}_{t_k^v} + \mathbf{w}' + \mathbf{e}' \\
&= \mathbf{w}' + \mathbf{e}',
\end{aligned}$$

where the last equation is due to the deadbeat control property (2). The quantity \mathbf{w}' is related to the plant disturbance and \mathbf{e}' is determined by the controller's estimation errors at the successful control packet transmissions, i.e., $\mathbf{e}_{t_k^1}, \dots, \mathbf{e}_{t_k^v}$, as given below:

$$\begin{aligned}
\mathbf{w}' &= \sum_{i=1}^{\eta_k^0} \mathbf{A}^{i-1} \mathbf{w}_{k-i} + (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0} \sum_{i=1}^{\eta_k^1} \mathbf{A}^{i-1} \mathbf{w}_{k-i} + \\
&\quad \dots + (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0 + \dots + \eta_k^{v-2}} \sum_{i=1}^{\eta_k^{v-1}} \mathbf{A}^{i-1} \mathbf{w}_{t_k^{v-1}-i} \\
&= \sum_{j=1}^v \left((\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\sum_{m=0}^{j-2} \eta_k^m} \sum_{i=t_k^j}^{t_k^{j-1}-1} \mathbf{A}^{t_k^{j-1}-1-i} \mathbf{w}_i \right), \\
\mathbf{e}' &= (\mathbf{A}^{\eta_k^0} - (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0}) \mathbf{e}_{t_k^1} + \\
&\quad (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0} (\mathbf{A}^{\eta_k^1} - (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^1}) \mathbf{e}_{t_k^2} + \\
&\quad \dots + (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^0 + \dots + \eta_k^{v-2}} (\mathbf{A}^{\eta_k^{v-1}} - (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^{v-1}}) \mathbf{e}_{t_k^v}.
\end{aligned}$$

By analyzing the correlation between the sequences of plant disturbance and estimation noise, we have the following result.

Proposition 1. *The per-step cost function about the plant state is a deterministic function of the vector state \mathbf{s}_k (see (19)) as*

$$J^x(\mathbf{s}_k) \triangleq \mathbb{E}[\mathbf{x}_k^\top \mathbf{S}^x \mathbf{x}_k] = \text{Tr}(\mathbf{S}^x V(\mathbf{s}_k)), \quad (20)$$

where $V(\mathbf{s}_k) \triangleq \mathbb{E}[\mathbf{x}_k \mathbf{x}_k^\top]$ is the plant state covariance given in (21). In the latter equation, for $0 \leq i \leq j \leq v$, we have

$$\Delta_k(i, j) \triangleq s_k^i - s_k^j = \sum_{n=i}^{j-1} \eta_k^n + \tau_k^j - \tau_k^i.$$

Then, building on the system dynamics (1), the local estimate (3), the remote estimate (4), and the control input (10), and by comprehensively analyzing the effect of the correlations between plant disturbance and estimation noise on the control input covariance, we obtain the per-step cost function about the control input as below.

Proposition 2. *The per-step cost function about the control input at k is a deterministic function of the vector state \mathbf{s}_k as*

$$\begin{aligned}
J^u(\mathbf{s}_k) &\triangleq \mathbb{E}[\mathbf{u}_k^\top \mathbf{S}^u \mathbf{u}_k] \\
&= \text{Tr} \left(\left(\tilde{\mathbf{K}} (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_{k-1}} \right)^\top \mathbf{S}^u \left(\tilde{\mathbf{K}} (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_{k-1}} \right) \hat{V}(\mathbf{s}_{k+1-\eta_k}) \right), \quad (22)
\end{aligned}$$

where $\mathbf{s}_{k+1-\eta_k}$ can be directly obtained by \mathbf{s}_k . $\hat{V}(\mathbf{s}_k) \triangleq \mathbb{E}[\hat{\mathbf{x}}_k \hat{\mathbf{x}}_k^\top]$ is the covariance of the remote estimate given in (23), where

$$\begin{aligned}
\tilde{\mathbf{D}}_k &\triangleq \mathbf{D}_k - \mathbf{A}^{\tau_k^0} \mathbf{Z} \Delta_k(0, v), \\
\tilde{\mathbf{E}}_{(k,n)}^i &\triangleq \tilde{\mathbf{E}}_{(k,n)}^i - \mathbf{A}^{\tau_k^0} \mathbf{Z} \Delta_k(0, n) - \tau_k^n + i (\mathbf{I} - \hat{\mathbf{K}} \mathbf{C}), \\
\tilde{\mathbf{F}}_{(k,n)}^i &\triangleq \tilde{\mathbf{F}}_{(k,n)}^i + \mathbf{A}^{\tau_k^0} \mathbf{Z} \Delta_k(0, n) + i \hat{\mathbf{K}}.
\end{aligned}$$

The proofs of Propositions 1 and 2 are given in [42] due to the space limitation.

Remark 3. *From Propositions 1 and 2, the per-step cost of the plant is determined by the finite-length vector state \mathbf{s}_k . However, expressions for the cost functions are involved due to the sequential predictive control and the command buffer adopted at the actuator, as well as the consideration of plant disturbance, and local and remote estimation errors.*

By substituting plant index i into the vector state and the cost functions (20) and (22), the above results have opened the door to address the optimal scheduling problem of the N -plant- M -frequency WNCs with ETDC in (14) as a decision making problem with a countable state space.

C. Resulting MDP

From (11), (17) and (18), given the current state $\mathbf{s}_{i,k}$ and the current transmission scheduling action related to plant i , the next state, $\mathbf{s}_{i,k+1}$, is independent of all previous states and actions, satisfying the Markov property. Thus, the transmission scheduling problem of the WNCs can be formulated as an MDP:

- The state of the MDP at time k is $\mathbf{s}_k \triangleq (\mathbf{s}_{1,k}, \dots, \mathbf{s}_{N,k})$, where $\mathbf{s}_{i,k}$ is as defined in (19), $i = 1, \dots, N$. The state space is $\mathcal{S} = \underbrace{\mathbb{N}^{2v_1+2} \times \dots \times \mathbb{N}^{2v_N+2}}_{N \text{ terms}}$.
- The action at time k , $\mathbf{a}_k \triangleq [a_{1,k}, a_{2,k}, \dots, a_{M,k}] = \pi(\mathbf{s}_k)$, is the transmission link allocation at each frequency, where $a_{m,k} \in \{-N, \dots, N\}$, $m = 1, \dots, M$, and $a_{m,k} \neq a_{m',k}$ if $a_{m,k} a_{m',k} \neq 0$. Then, the action space $\mathcal{A} \subseteq \{-N, \dots, N\}^M$ has the cardinality of $|\mathcal{A}| = \sum_{m=0}^M \mathbf{C}_M^m \mathbf{P}_{2N}^m$.
- The state transition probability $P(\mathbf{s}_{k+1} | \mathbf{s}_k, \mathbf{a}_k)$ can be obtained directly from the state updating rules in (11), (17) and (18).
- The per-step cost of the MDP is the sum cost of each plant in Propositions 1 and 2 as

$$c(\mathbf{s}_k) \triangleq \sum_{i=1}^N c_i(\mathbf{s}_{i,k}), \quad (24)$$

$$\begin{aligned}
V(\mathbf{s}_k) &= \mathbf{D}_k \hat{\mathbf{P}}^s (\mathbf{D}_k)^\top + \sum_{n=1}^{v-1} \left(\sum_{i=0}^{\Delta_k(n,n+1)-1} \check{\mathbf{E}}_{(k,n)}^i \mathbf{Q}_w (\check{\mathbf{E}}_{(k,n)}^i)^\top \right) + \sum_{i=0}^{\eta_k^0 + \tau_k^1 - 1} \mathbf{A}^i \mathbf{Q}_w (\mathbf{A}^i)^\top + \sum_{n=1}^{v-1} \left(\sum_{i=0}^{\Delta_k(n,n+1)-1} \check{\mathbf{F}}_{(k,n)}^i \mathbf{Q}_w (\check{\mathbf{F}}_{(k,n)}^i)^\top \right), \\
\mathbf{D}_k &\triangleq \sum_{j=1}^v \left((\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}}) \sum_{m=0}^{j-2} \eta_k^m (\mathbf{A}^{\eta_k^{j-1}} - (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^{j-1}}) \mathbf{A}^{\tau_k^j} \mathbf{Z}^{\Delta_k(j,v)} \right), \\
\check{\mathbf{E}}_{(k,n)}^i &\triangleq (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}}) \sum_{m=0}^{n-1} \eta_k^m \mathbf{A}^{\tau_k^{n+i}} + \sum_{j=1}^n \left((\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}}) \sum_{m=0}^{j-2} \eta_k^m (\mathbf{A}^{\eta_k^{j-1}} - (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^{j-1}}) \mathbf{A}^{\tau_k^j} \mathbf{Z}^{i+\Delta_k(j,n)} (\mathbf{I} - \hat{\mathbf{K}}\mathbf{C}) \right), \\
\check{\mathbf{F}}_{(k,n)}^i &\triangleq \sum_{j=1}^n \left((\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}}) \sum_{m=0}^{j-2} \eta_k^m (\mathbf{A}^{\eta_k^{j-1}} - (\mathbf{A} + \mathbf{B}\tilde{\mathbf{K}})^{\eta_k^{j-1}}) \mathbf{A}^{\tau_k^j} \mathbf{Z}^{\Delta_k(j,n)+i} \hat{\mathbf{K}} \right)
\end{aligned} \tag{21}$$

$$\begin{aligned}
\hat{V}(\mathbf{s}_k) &= \tilde{\mathbf{D}}_k \hat{\mathbf{P}}^s \tilde{\mathbf{D}}_k^\top + \sum_{n=1}^{v-1} \left(\sum_{i=0}^{\Delta_k(n,n+1)-1} \check{\mathbf{E}}_{(k,n)}^i \mathbf{Q}_w (\check{\mathbf{E}}_{(k,n)}^i)^\top \right) + \sum_{i=0}^{\Delta_k(0,1)-1} (\mathbf{A}^{\tau_k^0+i} - \mathbf{A}^{\tau_k^0} \mathbf{Z}^i (\mathbf{I} - \hat{\mathbf{K}}\mathbf{C})) \mathbf{Q}_w (\mathbf{A}^{\tau_k^0+i} - \mathbf{A}^{\tau_k^0} \mathbf{Z}^i (\mathbf{I} - \hat{\mathbf{K}}\mathbf{C}))^\top \\
&+ \sum_{n=1}^{v-1} \left(\sum_{i=0}^{\Delta_k(n,n+1)-1} \check{\mathbf{F}}_{(k,n)}^i \mathbf{Q}_w (\check{\mathbf{F}}_{(k,n)}^i)^\top \right) + \sum_{i=0}^{\Delta_k(0,1)-1} (\mathbf{A}^{\tau_k^0} \mathbf{Z}^i \hat{\mathbf{K}}) \mathbf{Q}_w (\mathbf{A}^{\tau_k^0} \mathbf{Z}^i \hat{\mathbf{K}})^\top
\end{aligned} \tag{23}$$

where $c_i(\mathbf{s}_{i,k}) \triangleq J_i^x(\mathbf{s}_{i,k}) + J_i^u(\mathbf{s}_{i,k})$.

- The discount factor is $\vartheta \in (0, 1)$.
- The scheduling problem of the N -plant- M -frequency system can be rewritten as

$$\min_{\pi} \mathbb{E}^{\pi} \left[\sum_{k=0}^{\infty} \vartheta^k c(\mathbf{s}_k) \right]. \tag{25}$$

Remark 4. The discounted MDP problem above with an infinite state space can be numerically solved to some extent by using classic policy or value iteration methods with a truncated state space \mathcal{S}_L

$$\mathcal{S}_L \triangleq \underbrace{\mathbb{N}_L^{2v_1+2} \times \dots \times \mathbb{N}_L^{2v_N+2}}_{N \text{ terms}}, \text{ where } \mathbb{N}_L = \{1, \dots, L\}.$$

The computation complexity of relative value iteration algorithm is given as $\mathcal{O}(|\mathcal{S}_L| |\mathcal{A}|^2 K)$ [43], where K is the number of iteration steps, and the state space and action space sizes are $|\mathcal{S}| = L^{2 \sum_{i=1}^N v_i + 2N}$ and $|\mathcal{A}| = \sum_{m=0}^M \mathbf{C}_M^m \mathbf{D}_{2N}^m$, respectively. However, the sizes of both the state and action spaces are considerably large even for relatively small N and M , leading to numerical difficulties in finding a solution. In the literature of WNCS, even for the (simpler) N -plant- M -channel remote estimation system, only the $M = 1$ case has been found to have a numerical solution [9]. To tackle the challenge for larger scale WNCS deployment, we will use DRL methods exploiting deep neural networks for function approximation in the following.

Remark 5. Although the MDP problem formulation of the WNCS assumes static wireless channels with fixed packet drop probabilities, it can be extended to a Markov fading channel scenario. A Markov fading channel can have multiple channel states with different packet drop probabilities, and the channel state transition is modeled by a Markov chain [33]. In this scenario, the state of the MDP problem should also include the channel state, and the state transition probability needs to take into account both the AoI state and the Markov channel state transition probabilities. The details of WNCS scheduling

over Markov fading channels can be investigated in our future work.

V. WNCS SCHEDULING WITH DEEP REINFORCEMENT LEARNING

Building on the MDP framework, DRL is widely applied in solving decision making problems with pre-designed state space, action space, per-step reward (cost) function and discount factor for achieving the maximum long-term reward. The main difference is that DRL does not exploit the state transition probability as required by MDP, but records and utilizes many sampled data sequences, including the current state \mathbf{s} , action \mathbf{a} , reward r and next state \mathbf{s}' , to train deep neural networks for generating the optimal policy [44]. To find deterministic policies⁵, the most widely considered DRL algorithms are DQN [45], DDPG [46], and TD3 [47]. In this study, we applied three algorithms to address the scheduling problem. DQN, which is the focus of our detailed algorithmic presentation, is the foundational algorithm in the field of deep reinforcement learning, primarily designed for environments with discrete action spaces. DDPG extends the ideas from DQN by adapting them to continuous action spaces using policy gradient methods. TD3, further building on DDPG, introduces improvements in terms of stability and performance. DQN serves as a foundation for understanding the basic principles that are also applicable to the more complex algorithms like DDPG and TD3, and is well-suited for our problem (Section IV-C) with the discrete and finite action space \mathcal{A} . In this regard, we choose to describe the DQN algorithm in detailed technical depth to avoid redundancy. Detailed pseudocode for DDPG and TD3, while omitted from the main text to conserve space, is available in [46], [47] for interested readers. Although DDPG and TD3 are more typically applied to continuous action spaces, they can be adapted for discrete

⁵The present work focuses on deterministic scheduling policies as it has been proved that the optimal policy is deterministic for unconstrained MDP problems [28].

actions through the action embedding method introduced in this section.

In general, we consider an episodic training scenario, where the system is operating in a simulated environment. Hence, one can learn and gather samples over multiple episodes to train a policy. Then, after verifying the properties of the final policy, one can “safely” deploy it [46]. We note that if the DRL policy does lead to instability (i.e., an unbounded expected cost), one needs to adjust the hyperparameters of the deep neural network (e.g., the initial network parameters and the number of layers and neurons) and start retraining. Such readjusting of hyperparameters of a deep learning agent is commonly adopted in practice [45].

In the following, we present a deep-Q-learning approach for scheduler design that requires a sampled data sequence $\{(s, \mathbf{a}, r, s')\}$. We note that the data sampling and the deep Q-learning are conducted offline. In particular, given the current state and action pair (s, \mathbf{a}) , the reward r can be obtained immediately from Propositions 1 and 2, and the next state s' affected by the scheduling action and the packet dropouts can be sampled easily based on the packet error probabilities of each channel (11) and state transition rules (17) and (18). Furthermore, the trained policy needs to be tested offline to verify the WNCS’s stability before an online deployment.

A. Deep Q-Learning Approach

We introduce the state-action value function given a policy $\pi(\cdot)$, which is also called the Q function [28]:

$$Q^\pi(s, \mathbf{a}) \triangleq \mathbb{E} \left[\sum_{k=0}^{\infty} \vartheta^k r_k \mid s_0 = s, \mathbf{a}_0 = \mathbf{a}, \mathbf{a}_k = \pi(s_k), \forall k > 0 \right],$$

where $r_k \triangleq -c(s_k)$ can be treated as the negative cost function in (24) at k . Then, by dropping out the time index k , the Q function of the maximum average discounted total reward achieving policy $\pi^*(\cdot)$ satisfies the Bellman equation as

$$Q^*(s, \mathbf{a}) = \mathbb{E} \left[r + \vartheta \max_{\mathbf{a}' \in \mathcal{A}} Q^*(s', \mathbf{a}') \mid s, \mathbf{a} \right]. \quad (26)$$

The optimal stationary and deterministic policy can be written as

$$\mathbf{a}^* = \pi^*(s) \triangleq \arg \max_{\mathbf{a} \in \mathcal{A}} Q^*(s, \mathbf{a}).$$

Solving the optimal Q function is the key to find the optimal policy, but is computationally intractable by conventional methods as discussed in Remark 4. In contrast, deep Q-learning methods approximate $Q^*(s, \mathbf{a})$ by a function $Q(s, \mathbf{a}; \theta)$ parameterized by a set of neural network parameters, θ , (including both weights and biases), and then learns θ to minimize the difference between the left- and right-hand sides of (26) [48]. Deep Q-learning can be easily implemented by the most well-known machine learning framework TensorFlow with experience replay buffer, ϵ -greedy exploration and mini-batch sampling techniques [49]. Building on this, the approach for solving problem (25) is given as Algorithm 1.

Algorithm 1 Deep Q-learning for transmission scheduling in WNCS

- 1: Initialize experience replay buffer \mathcal{B} to capacity K
 - 2: Initialize multi-layer fully connected neural network Q with vector input s , $|\mathcal{A}|$ outputs $\{Q(s, \mathbf{a}_1; \theta_0), \dots, Q(s, \mathbf{a}_{|\mathcal{A}|}; \theta_0)\}$ and random parameter set θ_0
 - 3: **for** episode = 1, \dots , E **do**
 - 4: Randomly initialize s_0
 - 5: **for** $t = 0, 1, \dots, T$ **do**
 - 6: With probability ϵ select a random action \mathbf{a}_t , otherwise select $\mathbf{a}_t = \arg \max_{\mathbf{a} \in \mathcal{A}} Q(s_t, \mathbf{a}; \theta_t)$
 - 7: Execute \mathbf{a}_t , and obtain r_t and s_{t+1}
 - 8: Store $(s_t, \mathbf{a}_t, r_t, s_{t+1})$ in \mathcal{B}
 - 9: Sample random mini-batch of l transitions $(s_t, \mathbf{a}_t, r_t, s_{t+1})$ from \mathcal{B} as $\tilde{\mathcal{B}}$
 - 10: Set $z_j = r_j + \vartheta \max_{\mathbf{a}' \in \mathcal{A}} \hat{Q}(s_{j+1}, \mathbf{a}'; \theta_t)$ for each sample in $\tilde{\mathcal{B}}$
 - 11: Perform a mini-batch gradient descend step to minimize the Bellman error, i.e., $\min_{\hat{\theta}} \sum_{\tilde{\mathcal{B}}} (z_j - Q(s_j, \mathbf{a}_j; \hat{\theta}))^2$
 - 12: Update $\theta_{t+1} = \hat{\theta}$
 - 13: **end for**
 - 14: **end for**
-

B. Deep Q-Learning with Reduced Action Space

In practice, the training of deep Q-network (DQN) converges and the trained DQN performs well for simple decision-making problems with a bounded reward function, small input state dimension and small action space, see e.g., [49], where the state input is a length-4 vector and there are only 2 possible actions. However, for the problem of interest in the current work, the reward function is unbounded (due to the potential of consecutive packet dropouts), the state s is of high dimension, and action-space $|\mathcal{A}|$ is large. Hence, the convergence of the DQN training is not guaranteed for all hyper-parameters, which include the initialization of θ , the number of neural network layers, the number of neurons per layer, and the replay buffer and mini-batch sizes.⁶ Even if the convergence is achieved, due to the complexity of the problem, it may take very long training time and only converge to a local optimal parameter set θ (which is also affected by the choice of hyper-parameters), leading to a worse performance than some conventional scheduling policies. Whilst the hyper-parameters can in principle be chosen to enhance performance, no appropriate tuning guidelines exist for the problem at hand.

To overcome the computational issues outlined above, in the following we will reduce the action space. This will make the training task simpler and enhance the convergence rate. To be more specific, for each plant system, instead of considering all the possible actions to schedule either or both of the uplink and downlink communications at each time instant, we restrict to the two modes *downlink only* and *uplink only*. Switching between these two modes of operation occurs once a scheduled transmission of that plant is successful.

⁶Note that when $N = 3$, $M = 3$ and $v_1 = \dots = v_N = 2$, the length of the state is 18 and the action space size is 229 based on Section IV-C.

The principle behind this is the intuition that the controller requires the information from the sensor first and then performs control, and so on and so forth. The above concept can be incorporated into the current framework by adding the downlink/uplink indicator $s_{i,k}^{\text{link}} \in \{-1, 1\}$, $\forall i = 1, \dots, N$, into state \mathbf{s}_k leading to the aggregated state

$$\check{\mathbf{s}}_k \triangleq [\mathbf{s}_k, s_{1,k}^{\text{link}}, \dots, s_{N,k}^{\text{link}}],$$

where $s_{i,k}^{\text{link}} = 1$ or -1 indicates the uplink or downlink of plant i can be scheduled at k . The state updating rule is

$$s_{i,k+1}^{\text{link}} = \begin{cases} -1 & \text{if } s_{i,k}^{\text{link}} = 1 \text{ and } \beta_{i,k} = 1 \\ 1 & \text{if } s_{i,k}^{\text{link}} = -1 \text{ and } \gamma_{i,k} = 1 \\ s_{i,k}^{\text{link}} & \text{otherwise.} \end{cases}$$

The new scheduling action at the M frequencies is denoted as $\check{\mathbf{a}}_k \triangleq [\check{a}_{1,k}, \dots, \check{a}_{M,k}]$, where $\check{a}_{m,k} \in \{0, \dots, N\}$, $\forall m \in \{1, \dots, M\}$, and $\check{a}_{m,k} \neq \check{a}_{m',k}$ if $\check{a}_{m,k} \check{a}_{m',k} \neq 0$. If $\check{a}_{m,k} = i$ and $s_{i,k}^{\text{link}} = 1$, then the uplink of plant i is scheduled on frequency m at time k . Compared with the original action \mathbf{a}_k in Section IV-C, the size of the new action space is reduced significantly to $|\check{\mathcal{A}}| = \sum_{m=0}^M C_M^m P_N^m$. For example, when $N = 3$, $M = 3$ and $v_1 = v_2 = v_3 = 2$, the action space size is reduced from 229 to 34. By replacing \mathbf{s} and \mathbf{a} with $\check{\mathbf{s}}$ and $\check{\mathbf{a}}$, respectively, in Algorithm 1, we can apply the deep Q-learning method to find a policy with a reduced action space at a faster rate. It should be noted that the reduced action space still rapidly increases with the increment of N and M . This makes it difficult for the DQN-based algorithm to solve the scheduling problem of larger-scale WNCS, such as $N = 10$ and $M = 10$ with a discrete action space of 234662231.

C. Action Embedding for DDPG and TD3

To further handle the issue of large discrete action space, the discrete actions can be embedded into a continuous action space. This approach uses a vector of continuous values to represent the priority of allocating channels to plants. Specifically, each plant is assigned a continuous value. Plants are then allocated channels based on the ranking of these continuous values. For instance, the plant with the highest value is allocated to the first channel, the second highest to the second channel, and so on. This method significantly reduces the complexity of the discrete action space, making it feasible to handle larger N and M . For example, the number of output neurons is 6 when $N = 6$ and $M = 4$, while the original discrete action space is 1045. The continuous action space allows for more efficient learning and optimization by leveraging the capabilities of DRL algorithms designed for continuous spaces. We will illustrate numerical results of using DRL for solving the scheduling problem in the following section.

VI. NUMERICAL RESULTS

The plant system matrices \mathbf{A}_i are randomly generated, in which case $\rho(\mathbf{A}_i) \in (1.0, 1.1)$. The control input matrices are all set to $\mathbf{B}_i = [1 \ 1]^T$. The measurement matrix \mathbf{C}_i is equal to the identity matrix. The covariance matrices are $\mathbf{Q}_i^w =$

$\mathbf{Q}_i^v = 0.1\mathbf{I}$. Each plant is 2-step controllable, i.e., $v_i = 2$. The packet success probabilities of each link, i.e., $\xi_{m,i}^s$ and $\xi_{m,i}^c$, $i, m = 1, 2, 3$, are generated randomly and drawn uniformly from $(0.5, 1.0)$. The weighting terms \mathbf{S}_i^x and \mathbf{S}_i^u are chosen as identity matrices and the discount factor is $\vartheta = 0.95$. The resulting WNCS satisfies the stability condition established in Theorem 1.

We adopt the following hyper-parameters for the DRL algorithms. The input state dimension of DQN is set to $2N(v_i + 1)$. We use three hidden layers with 300, 200, 100 neurons, respectively. The output layer of DQN has $|\check{\mathcal{A}}| = \sum_{m=0}^M C_M^m P_N^m$ outputs (actions) for the problem with reduced action space (see Section V-B). The activation function in each hidden layer is the ReLU [49]. The activation function in the output layer of DDPG and TD3 is the Sigmoid [49]. The experience reply memory has a size of $K = 100000$, and the size of mini-batch is 64. So in each time step, 64 data is sampled from the reply memory for training. The exploration parameter ϵ of DQN decreases from 1 to 0.01 at the rate of 0.999 after each time step. The exploration method for DDPG and TD3 is based on Gaussian action noise with a standard deviation of 0.2 and a mean value of 0. We adopt the ADAM optimizer to update the DNNs. Each training takes $E = 500$ episodes, each including $T = 500$ time steps. We test the trained policy for 100 episodes and compare the empirical average costs $\frac{1}{T} \sum_{k=1}^T \sum_{i=1}^N c(\mathbf{s}_{i,k})$.

We use Algorithm 1 (DQN) and its adapted version, i.e., DDPG and TD3, to solve the scheduling problem with reduced action space and compare it with the original DQN, DDPG and TD3, and three heuristic benchmark policies. 1) **Random policy**: randomly choose M out of the $2N$ links and randomly allocate them to the M frequencies at each time. 2) **Round-robin policy**: the $2N$ links are divided into M groups, and the links in each of the groups $m \in \{1, \dots, M\}$ are allocated to the corresponding frequency m ; the transmission scheduling at each frequency follows the round-robin fashion [20] at each time. Note that in the simulation, we test all link-grouping combinations and only present the one with the lowest average cost. 3) **Greedy policy**: At each time, the M frequencies are randomly allocated to M of the $2N$ links with the largest AoI value in $\{\tau_{1,k}, \eta_{1,k}, \dots, \tau_{N,k}, \eta_{N,k}\}$. Recall that $\tau_{i,k}$ and $\eta_{i,k}$ denote the time duration since the last received sensing and control packets of plant i , respectively.

We first consider a N -plant- M -frequency WNCS with $N = 5$ and $M = 5$, i.e., 10 links sharing 5 frequencies with a reduced discrete action space of 1546. Then, we will extend the network to a larger scale with $N = 10$ and $M = 10$, i.e., 20 links sharing 10 frequencies with a reduced discrete action space of 234662231.

Fig. 3 showcases the comparison of three DRL algorithms and the three heuristic benchmark policies when $N = 5$ and $M = 5$. All the DRL algorithms with the reduced action space, including DQN, DDPG, and TD3, achieve a lower long-term average cost than the best heuristic benchmark policy, i.e., Greedy Policy. We note that DRL is a dynamic decision-making algorithm designed to maximize the cumulative reward

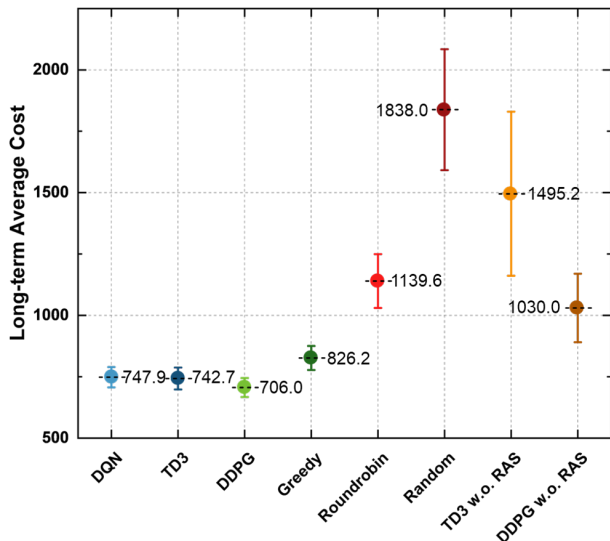


Fig. 3. The long-term average performance of DRL and benchmark algorithms over a system with $N = 5$ and $M = 5$. Values are means \pm standard error of mean. RAS denotes reduced action space.

over time, inherently aligning with the goal of optimizing long-term costs. It evaluates actions based on their expected future rewards, considering both immediate and delayed consequences. Heuristic benchmark policies are based on static decision rules and lack the ability to adapt to changing environments or evolving system dynamics. They often prioritize immediate rewards or costs, making decisions that optimize short-term performance without considering long-term implications. They might select actions that yield immediate gains but could lead to suboptimal long-term outcomes, such as increased future costs. Thus, DRL algorithms show positive performance.

In addition, Fig. 3 demonstrates that the well-trained DQN-based policy reduces the long-term average cost of greedy policy by 9.5% from 826.2 to 747.9. The DDPG and TD3 algorithms based on the action embedding method achieve similar or slightly better performance compared to the DQN-based algorithm. This is attributed to the large discrete action space, which makes it difficult for the DQN-based algorithm to explore adequately during training and to get the optimal policy after training. The DDPG algorithm reduces the long-term average cost of the greedy policy by 14.5% from 826.2 to 706.0. TD3 is an extension of DDPG, which is designed to reduce overestimation bias and variance in the learning process by introducing delayed policy updates, target policy smoothing, and twin critic networks. These features can lead to overly conservative behaviour in WNCs with high stochasticity, which potentially slow down the learning process or lead to underfitting where the system fails to adequately capture the beneficial actions amidst noise. Thus, DDPG, as a more aggressive approach than TD3.

Fig. 3 also shows that the DQN algorithm without reduced action space (RAS) has an action space size of 63591, which is unable to solve the original MDP problem with such a large action space. The DDPG and TD3 algorithms without RAS

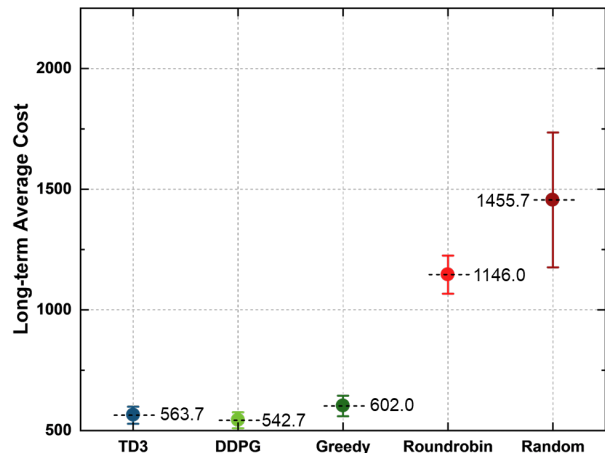


Fig. 4. The long-term average performance of DRL and benchmark algorithms over a system with $N = 10$ and $M = 10$. Values are means \pm standard error of mean.

cannot outperform the ones with RAS. These results illustrate the effectiveness of the proposed action space reduction method.

Fig. 4 showcases the comparison of three DRL algorithms and the three heuristic benchmark policies when $N = 10$ and $M = 10$. In this case, the DQN-based algorithm is unable to solve the MDP, as the reduced discrete action space is still too large. Results show that the DDPG and TD3 algorithms based on the action embedding method effectively address the challenge of the rapidly growing action space, making it more suitable for larger-scale WNCs problems. The DDPG algorithm reduces the long-term average cost of the greedy policy by 9.9% from 542.7 to 602.0.

VII. CONCLUSIONS

We have investigated the transmission scheduling problem of the fully distributed WNCs. A sufficient stability condition of the WNCs in terms of both the control and communication system parameters has been derived. We have proposed advanced DRL algorithms to solve the scheduling problem, which performs much better than benchmark policies. For future work, we will develop a distributed DRL approach for enhancing the scalability of the scheduling algorithm. Furthermore, we will consider co-design problems of the estimator, the controller, and the scheduler for distributed WNCs with time-varying channel conditions.

REFERENCES

- [1] P. Park, S. Coleri Ergen, C. Fischione, C. Lu, and K. H. Johansson, "Wireless network design for control systems: A survey," *IEEE Commun. Surveys Tuts.*, vol. 20, no. 2, 2018.
- [2] K. Huang, W. Liu, Y. Li, B. Vucetic, and A. Savkin, "Optimal downlink-uplink scheduling of wireless networked control for Industrial IoT," *IEEE Internet Things J.*, Mar. 2020.
- [3] J. Deng, S. Liu, M. Zhang, and C. Ye, "Sequential fusion estimation for multisensor multirate systems with measurement outliers," *Asian J. Control*, vol. 25, no. 6, pp. 4909–4919, 2023.
- [4] S. Liu, X. Zhao, E. Tian, and G. Wei, "Distributed recursive filtering under random access protocols: A multirate strategy," *Int. J. Robust Nonlinear Control*, vol. 32, no. 12, pp. 7132–7148, 2022.

- [5] L. Zhao, W. Zhang, J. Hu, A. Abate, and C. J. Tomlin, "On the optimal solutions of the infinite-horizon linear sensor scheduling problem," *IEEE Trans. Autom. Control*, vol. 59, no. 10, pp. 2825–2830, 2014.
- [6] D. Han, J. Wu, H. Zhang, and L. Shi, "Optimal sensor scheduling for multiple linear dynamical systems," *Automatica*, vol. 75, pp. 260–270, 2017.
- [7] A. S. Leong, S. Dey, and D. E. Quevedo, "Sensor scheduling in variance based event triggered estimation with packet drops," *IEEE Trans. Autom. Control*, vol. 62, no. 4, pp. 1880–1895, Apr. 2017.
- [8] J. Wei and D. Ye, "Double threshold structure of sensor scheduling policy over a finite-state markov channel," *IEEE Trans. Cybern.*, pp. 1–10, 2022.
- [9] S. Wu, X. Ren, S. Dey, and L. Shi, "Optimal scheduling of multiple sensors over shared channels with packet transmission constraint," *Automatica*, 2018.
- [10] L. Peng, X. Cao, and C. Sun, "Optimal transmit power allocation for an energy-harvesting sensor in wireless cyber-physical systems," *IEEE Trans. Cybern.*, vol. 51, no. 2, pp. 779–788, 2021.
- [11] K. Gatsis, M. Pajic, A. Ribeiro, and G. J. Pappas, "Opportunistic control over shared wireless channels," *IEEE Trans. Autom. Control*, vol. 60, pp. 3140–3155, Dec. 2015.
- [12] M. Eisen, M. M. Rashid, K. Gatsis, D. Cavalcanti, N. Himayat, and A. Ribeiro, "Control aware radio resource allocation in low latency wireless control systems," *IEEE Internet Things J.*, vol. 6, no. 5, pp. 7878–7890, 2019.
- [13] K. Wang, W. Liu, and T. J. Lim, "Deep reinforcement learning for joint sensor scheduling and power allocation under dos attack," in *Proc. IEEE ICC*, pp. 1968–1973, 2022.
- [14] R. Gan, Y. Xiao, J. Shao, and J. Qin, "An analysis on optimal attack schedule based on channel hopping scheme in cyber-physical systems," *IEEE Trans. Cybern.*, vol. 51, no. 2, pp. 994–1003, 2021.
- [15] K. Huang, W. Liu, M. Shirvanimoghaddam, Y. Li, and B. Vucetic, "Real-time remote estimation with hybrid ARQ in wireless networked control," *IEEE Trans. Wireless Commun.*, vol. 19, no. 5, pp. 3490–3504, 2020.
- [16] L. Liu, A. Chattopadhyay, and U. Mitra, "On solving MDPs with large state space: Exploitation of policy structures and spectral properties," *IEEE Trans. Comm.*, vol. 67, no. 6, pp. 4151–4165, 2019.
- [17] X. Chen, W. Chen, J. Lee, and N. B. Shroff, "Delay-optimal buffer-aware scheduling with adaptive transmission," *IEEE Trans. Comm.*, vol. 65, no. 7, pp. 2917–2930, 2017.
- [18] D. Han, W. Chen, and Y. Fang, "Joint channel and queue aware scheduling for latency sensitive mobile edge computing with power constraints," *IEEE Trans. Wireless Comm.*, vol. 19, no. 6, pp. 3938–3951, 2020.
- [19] Y. Li, A. S. Mehr, and T. Chen, "Multi-sensor transmission power control for remote estimation through a SINR-based communication channel," *Automatica*, vol. 101, pp. 78–86, 2019.
- [20] A. S. Leong, A. Ramaswamy, D. E. Quevedo, H. Karl, and L. Shi, "Deep reinforcement learning for wireless sensor scheduling in cyber-physical systems," *Automatica*, 2020.
- [21] B. Demirel, A. Ramaswamy, D. E. Quevedo, and H. Karl, "DeepCAS: A deep reinforcement learning algorithm for control-aware scheduling," *IEEE Control Syst. Lett.*, vol. 2, no. 4, pp. 737–742, 2018.
- [22] J. Zhu, Y. Song, D. Jiang, and H. Song, "A new deep-Q-learning-based transmission scheduling mechanism for the cognitive Internet of Things," *IEEE Internet Things J.*, vol. 5, no. 4, pp. 2375–2385, 2017.
- [23] F. Naeem, S. Seifollahi, Z. Zhou, and M. Tariq, "A generative adversarial network enabled deep distributional reinforcement learning for transmission scheduling in internet of vehicles," *IEEE Trans. Intell. Transp. Syst.*, vol. 22, no. 7, pp. 4550–4559, 2020.
- [24] Y. Ma, X. Hao, J. Hao, J. Lu, X. Liu, T. Xialiang, M. Yuan, Z. Li, J. Tang, and Z. Meng, "A hierarchical reinforcement learning based optimization framework for large-scale dynamic pickup and delivery problems," *Adv. Neural Inf. Proces. Syst.*, vol. 34, pp. 23609–23620, 2021.
- [25] F. Ni, J. Hao, J. Lu, X. Tong, M. Yuan, J. Duan, Y. Ma, and K. He, "A multi-graph attributed reinforcement learning based optimization algorithm for large-scale hybrid flow shop scheduling problem," in *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, pp. 3441–3451, 2021.
- [26] H. Yang and X. Xie, "An actor-critic deep reinforcement learning approach for transmission scheduling in cognitive internet of things systems," *IEEE Syst. J.*, vol. 14, no. 1, pp. 51–60, 2019.
- [27] F. D. Hildebrandt, B. W. Thomas, and M. W. Ulmer, "Opportunities for reinforcement learning in stochastic dynamic vehicle routing," *Comp. Oper. Res.*, vol. 150, p. 106071, 2023.
- [28] D. P. Bertsekas *et al.*, *Dynamic programming and optimal control: Vol. 1 (3rd ed.)*. Athena scientific Belmont, 2005.
- [29] J. Chen, W. Liu, D. E. Quevedo, Y. Li, and B. Vucetic, "Structure-enhanced deep reinforcement learning for optimal transmission scheduling," in *Proc. ICC*, pp. 1212–1218, IEEE, 2023.
- [30] A. Redder, A. Ramaswamy, and D. E. Quevedo, "Deep reinforcement learning for scheduling in large-scale networked control systems," *IFAC-PapersOnLine*, vol. 52, no. 20, pp. 333–338, 2019.
- [31] 3GPP, "5G NR Overall description Stage 2," Technical Specification (TS) 38.300, 3rd Generation Partnership Project (3GPP), 09 2018, Version 15.2.0.
- [32] B. Demirel, V. Gupta, D. E. Quevedo, and M. Johansson, "On the trade-off between communication and control cost in event-triggered dead-beat control," *IEEE Trans. Autom. Control*, vol. 62, no. 6, pp. 2973–2980, 2017.
- [33] W. Liu, D. E. Quevedo, Y. Li, K. H. Johansson, and B. Vucetic, "Remote state estimation with smart sensors over Markov fading channels," *IEEE Trans. Autom. Control*, vol. 67, no. 6, pp. 2743–2757, 2022.
- [34] T. Kailath and P. Hall, *Linear Systems*. Information and System Sciences Series, Prentice-Hall, 1980.
- [35] W. Liu, D. E. Quevedo, Y. Li, and B. Vucetic, "Anytime control under practical communication model," *IEEE Trans. Autom. Control*, vol. 67, no. 10, pp. 5400–5407, 2022.
- [36] J. O'Reilly, "The discrete linear time invariant time-optimal control problem—an overview," *Automatica*, no. 2, 1981.
- [37] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. Cambridge University Press, 2005.
- [38] S. Coleri, M. Ergen, A. Puri, and A. Bahai, "Channel estimation techniques based on pilot arrangement in ofdm systems," *IEEE Trans. Broadcast.*, pp. 223–229, 2002.
- [39] L. Schenato, B. Sinopoli, M. Franceschetti, K. Poolla, and S. S. Sastry, "Foundations of control and estimation over lossy networks," *Proc. IEEE*, vol. 95, pp. 163–187, Jan 2007.
- [40] S. R. Pokhrel, M. Panda, H. L. Vu, and M. Mandjes, "TCP performance over Wi-Fi: Joint impact of buffer and channel losses," *IEEE Trans. Mob. Comput.*, vol. 15, no. 5, pp. 1279–1291, 2016.
- [41] M. Polese, R. Jana, and M. Zorzi, "TCP and MP-TCP in 5G mmwave networks," *IEEE Internet Comput.*, vol. 21, no. 5, pp. 12–19, 2017.
- [42] W. Liu, K. Huang, D. E. Quevedo, B. Vucetic, and Y. Li, "Deep reinforcement learning for wireless scheduling in distributed networked control," available: <https://download.arxiv.org/pdf/2109.12562v1>, 2022.
- [43] L. I. Sennott, *Stochastic dynamic programming and the control of queueing systems*. John Wiley & Sons, 2009.
- [44] D. P. Bertsekas, *Reinforcement Learning and Optimal Control*. Belmont, MA: Athena Scientific, 2019.
- [45] V. Mnih, K. Kavukcuoglu, D. Silver, *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [46] 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," *Proc. ICLR*, 2015.
- [47] S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *Proc. ICML*, pp. 1587–1596, PMLR, 2018.
- [48] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [49] TensorFlow, "Train a deep Q network with TF-agents." https://www.tensorflow.org/agents/tutorials/1_dqn_tutorial, 2021.