Coexistence between Task- and Data-Oriented Communications: A Whittle's Index Guided Multi-Agent Reinforcement Learning Approach

Ran Li, Chuan Huang, Xiaoqi Qin, Shengpei Jiang, Nan Ma, and Shuguang Cui

Abstract-We investigate the coexistence of task-oriented and data-oriented communications in a IoT system that shares a group of channels, and study the scheduling problem to jointly optimize the weighted age of incorrect information (AoII) and throughput, which are the performance metrics of the two types of communications, respectively. This problem is formulated as a Markov decision problem, which is difficult to solve due to the large discrete action space and the time-varying action constraints induced by the stochastic availability of channels. By exploiting the intrinsic properties of this problem and reformulating the reward function based on channel statistics, we first simplify the solution space, state space, and optimality criteria, and convert it to an equivalent Markov game, for which the large discrete action space issue is greatly relieved. Then, we propose a Whittle's index guided multi-agent proximal policy optimization (WI-MAPPO) algorithm to solve the considered game, where the embedded Whittle's index module further shrinks the action space, and the proposed offline training algorithm extends the training kernel of conventional MAPPO to address the issue of time-varying constraints. Finally, numerical results validate that the proposed algorithm significantly outperforms state-of-the-art age of information (AoI) based algorithms under scenarios with insufficient channel resources.

Index Terms—Task-oriented communications, data-oriented communications, age of incorrect information (AoII), Whittle's index, Whittle's index guided multi-agent proximal policy optimization (WI-MAPPO), age of information (AoI)

I. INTRODUCTION

In the past decades, the spotlight of Internet of things (IoT) is shifting towards enabling autonomous networked control applications that require timely status updates [1], e.g., environmental monitoring [2], emergency detection [3], and healthcare systems [4]. Under such scenarios, IoT devices

This work was submitted in part to 2022 IEEE Global Communications Conference.

R. Li is with the School of Science and Engineering (SSE) and the Future Network of Intelligence Institute (FNii), The Chinese University of Hong Kong, Shenzhen 518172, China, (e-mail: ranli2@link.cuhk.edu.cn).

C. Huang and S. Cui are with the School of Science and Engineering (SSE) and Future Network of Intelligence Institute (FNii), The Chinese University of Hong Kong, Shenzhen 518172, China, and with Peng Cheng Laboratory, Shenzhen 518066, China, (e-mails: huangchuan@cuhk.edu.cn; shuguangcui@cuhk.edu.cn).

X. Qin is with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing 100876, China, (e-mail: xiaoqiqin@bupt.edu.cn).

S. Jiang is with the SF Technology, Shenzhen 518052, China, (e-mail: philip.jiang@sfmail.sf-express.com).

N. Ma is with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing 100876, China, and with the Department of Broadband Communication, Peng Cheng Laboratory, Shenzhen 518066, China, (e-mail: manan@bupt.edu.cn).

are deployed to obtain real-time awareness of a monitored physical process by continuous sampling and data uploading. Considering the time-varying nature of environment, the freshness of status updates is of critical importance to the performance of subsequent tasks. Age of information (AoI) [5] is proposed as a performance metric for such task-oriented communications, which is defined as the time elapsed since the latest received status update was generated at the monitoring device. Note that conventional resource provisioning schemes follow the data-oriented design philosophy of maximizing network-level throughput. Therefore, in order to accommodate time-sensitive data transmission in task-oriented communications, one straightforward approach is to reserve spectrum for devices with critical status updates to guarantee the timely delivery of data. However, such separate design degrades the spectrum efficiency due to the intermittent nature of monitoring devices. Therefore, it is important to establish a flexible coexistence strategy between task-oriented and dataoriented communications. This problem, if left unsolved, will jeopardize the utility of networked control systems.

The focus of recent research on resource scheduling strategies for time-sensitive applications is optimizing AoI among multiple devices. The authors in [6] studied the scenario with multiple IoT devices monitoring multiple processes and then transmitting their statuses to one central base station (BS) over one channel, and a near-optimal algorithm to minimize the average AoI was proposed. The authors in [7] extended the previous work to the multiple channels case, and proposed a low-complexity algorithm based on the Lagrangian relaxation method [8]. The authors in [9] studied the AoI minimization problem in a time-framed system, where multiple statuses can be transmitted during one time frame, and compared the performances among the randomized policy, the max-weight policy [10], and the Whittle's index algorithm. In order to strike a balance between the achievable throughput and age under sparse sampling at source nodes, various strategies have been proposed [11]–[13]. The authors in [11] studied the scenario with multiple IoT devices, and aimed to minimize their average AoIs while simultaneously satisfy the constraints on their throughputs. In [12], IoT was utilized to simultaneously monitor one process with a single monitoring device and collect data from multiple traditional devices, and the average AoI at the monitoring device and the throughputs of the traditional devices were evaluated under the ALOHA protocol. The authors in [13] considered the most general scenario with multiple monitoring devices acquiring statuses from multiple

processes and multiple traditional devices collecting data, and a Lyapunov drift based scheduling policy was proposed to jointly optimize the average AoI and throughput.

Recently, age of incorrect information (AoII) [14] is proposed as a more advanced performance metric that captures not only the aging of status updates, but more importantly the change of context of the monitored process at source node. To elaborate, the authors in [14], [15] discussed the scenario where a base station (BS) in IoT predicts the state of a process. Apparently, the BS prefers to collect new status information from the monitoring device when the state of the considered process indeed changes, not when the currently reserved status at the BS is not fresh. This indicates that AoII is more compatible with the prediction tasks than AoI.

There has been growing research interest in designing AoIIbased scheduling policy for monitoring devices [14], [16]-[19]. Specifically, the authors in [16], [17] discussed the scenario where the target process to be predicted simply follows the binary distribution, and the optimal policy for the case with one target process and a low-complexity suboptimal policy for the case with multiple target processes were developed. One step further, the authors in [14], [18], [19] discussed a more complex prediction scenarios, where the target processes have more than two states. Specifically, the authors in [18], [19] discussed the scenario where the state transitions of the processes only occur between the adjacent states, and a Whittle's index based suboptimal policies were proposed for the case with multiple target processes. The author in [14] considered the scenario where the state transitions between any two states are allowable, and derived the optimal policy for the case with one target process. However, the optimization framework in [14] cannot be directly extended to the case with multiple target processes and no existing work has touched this case.

In this paper, we employ AoII as the performance metric for task-oriented communications and investigate the coexistence strategy of task-oriented and data-oriented communications by jointly optimizing the average AoII and the throughput. Specifically, we study a general scenario: multiple processaware monitoring devices monitor and transmit the status updates of multiple random processes, respectively, which have more than two states and the state transition between any two states is allowable; multiple traditional devices purse high throughput and each round of the data transmissions for the traditional devices may last for multiple time slots; limited channel resources are available in IoT for the data transmissions of the monitoring and traditional devices. We summarize our contributions as follows:

• We formulate the joint optimization problem as a Markov decision problem, which is challenging due to the large solution space, large state space, and average optimality criteria. Moreover, it has a large discrete action space and time-varying action constraints induced by the stochastic availability of channels, where existing algorithms to solve Markov decision problems cannot efficiently address the problems of this type. To overcome these challenges, we first analyze the intrinsic properties of this Markov decision problem, and prove that there exist

stationary policies to achieve its optimum. Next, based on this stationary feature, we reformulate the reward function and transform the original problem to an equivalent form with a much smaller state space and a simplified state transitions. Then, we validate the existence of the Blackwell policies for the equivalent problem, replace the average optimality criteria by the discounted version, and prove that the optimal policies under this discounted criteria also optimize the problem with average optimality criteria. Finally, we convert the above discounted Markov decision problem as an equivalent Markov game by treating each channel as individual agent, for which the large discrete action space issue is relieved. Remarkably, all these problem simplifications and conversions are theoretically validated to be equivalent.

• We propose a multi-agent reinforcement learning algorithm, namely the Whittle's index guided multi-agent proximal policy optimization (WI-MAPPO), to efficiently solve the proposed Markov game. Specifically, WI-MAPPO deploys a Whittle's index guided action fusion module to further shrink the action space of the Markov game. To design this module, we first prove that Whittle's index for the monitoring devices exists by validating their indexabilities, and then utilize an efficient exhausted searching algorithm to approximate the Whittle's index. Finally, we construct this module by generating a sufficient large Whittle's index table. Moreover, we modify both the actor network and the probability ratio derivation of the training algorithm for multi-agent proximal policy optimization (MAPPO) to train the proposed WI-MAPPO. By doing so, the training procedure will not violate the time-varying constraints and ensure even faster and more accurate estimations on the advantage functions of MAPPO. Remarkably, although the time-varying constraints issue shrinks the solution space of the considered Markov game and is doomed to induce loss of optimality compared with the non-constrained Markov game, it is validated that the proposed WI-MAPPO can greatly narrow this optimality gap and achieve almost the same performances for the constrained and non-constrained Markov games.

The remainder of this paper is organized as follows. Section II introduces the system model and formulates the joint scheduling problem. Section III presents the proposed algorithm. Section IV evaluates the performance of the proposed algorithm. Finally, Section V concludes this paper.

II. SYSTEM MODEL AND PROBLEM FORMULATION

Consider a slotted IoT system as shown in Fig. 1, which consists of a BS and a set of IoT devices. The set of devices consists of two types, including I monitoring devices to enable real-time situational awareness at BS and J traditional devices for data collection. The monitoring devices are deployed at I different monitoring spots to obtain real-time perceptions about I different random process $\{X_i(t)\}_{i=1}^{I}$ by periodical data sampling and uploading, and the scheduling performance for monitoring devices is quantified by AoII, which captures



Figure 1: Joint schedule of the task-oriented and data-oriented communications.

the semantic attributes of information in terms of the relevance of transmitted data to the subsequent task. As for traditional devices, the scheduling performance can be quantified by the achievable network-level throughput. We assume an OFDMAbased system with M sub-channels. The data transmission for traditional devices may last for multiple consecutive time slots depending on the size of sampled packets, while the transmission for monitoring devices is assumed to be completed within one time slot due to small packet size for status update information. Notably, we also call the i^{th} monitoring device as the i^{th} device and the j^{th} traditional device as the $(I + j)^{\text{th}}$ device.

A. System model

The joint scheduling between the task-oriented and dataoriented communications is formulated as a Markov decision problem and we introduce its state, action, transitions, and reward as follows.

1) State: The state contains AoII vector, channel gains, and channel availability.

AoII vector: AoII counts the age of the incorrect prediction and increases as long as the predicted state is incorrect. Particularly, it is defined as follows [14].

Definition 2.1(AoII): The AoII at the i^{th} monitoring device in the t^{th} time slot is denoted as $x_i(t) \in \mathbb{Z}_{\geq 0}$, where $\mathbb{Z}_{\geq 0}$ is the set of all non-negative integers, and recursively defined by

$$x_i(t+1) \triangleq \begin{cases} 0 & \bar{X}_i(t+1) = X_i(t+1) \\ x_i(t) + 1 & \bar{X}_i(t+1) \neq X_i(t+1). \end{cases}$$
(1)

where $\bar{X}_i(t)$ is the prediction on $X_i(t)$ made by the BS.

Then, the *AoII vector* at *I* monitoring devices is defined as $\boldsymbol{x}(t) \triangleq [x_1(t), x_2(t), \cdots, x_I(t)]^T$.

Channel gains: Denote the channel coefficient and the channel gain of the link between the BS and the j^{th} traditional device over the m^{th} channel at the t^{th} time slot as $h_{j,m}(t)$ and $g_{j,m}(t)$, respectively, i.e., $g_{j,m}(t) = |h_{j,m}(t)|^2$. Then, define the *channel gains* at the t^{th} time slot as matrix $\boldsymbol{G}(t) \in \mathcal{G}^{J \times M}$, where the $(j,m)^{\text{th}}$ entry of $\boldsymbol{G}(t)$ is $g_{j,m}(t)$, i.e., $[\boldsymbol{G}(t)]_{(j,m)} \triangleq g_{j,m}(t)$, and \mathcal{G} is the value space of $g_{j,m}(t)$ and considered as a finite set.

Channel availability: Each data transmission for the jth traditional device is considered to consume $T_i \in \mathbb{Z}^+$ consecutive time slots with \mathbb{Z}^+ being the set of all positive integers, during which the occupied channel is not available for new data transmission. Specifically, denote the availability condition of the m^{th} channel for the data transmission of the j^{th} traditional device at the tth time slot as $b_{i,m}(t) \in \{0, 1, \dots, T_i - 1\}$: if the m^{th} channel is not transmitting data for the j^{th} traditional device at the t^{th} time slot, $b_{i,m}(t)$ is set to 0; otherwise, $b_{i,m}(t)$ is equal to the number of the remaining time slots for the release of the m^{th} channel. Then, define the *channel availability* as a J-by-M-dimension matrix B(t), where the $(j,m)^{\text{th}}$ entry of $\boldsymbol{B}(t)$ equals $b_{j,m}(t)$, i.e., $[\boldsymbol{B}(t)]_{(j,m)} \triangleq b_{j,m}(t)$. Note that the status update transmission of monitoring devices is assumed to be completed within one time slot, and thus the allocated channel will always be released for new data transmission at the next time slot.

Denote the state of the considered Markov decision problem at the t^{th} time slot as s(t). Obviously, $s(t) = (\boldsymbol{x}(t), \boldsymbol{G}(t), \boldsymbol{B}(t))$ holds and the state space $\mathcal{S} = \mathbb{Z}_{\geq 0}^{I \times 1} \times \mathcal{G}^{J \times M} \times \mathcal{B}$ is countable, where \mathcal{B} is the value space of $\boldsymbol{B}(t)$. 2) Action: Denote the scheduling decision for the m^{th} channel at the t^{th} time slot as $a_m(t) \in \{0, 1, \dots, I + J\}$. Specifically, $a_m(t) = 0$ means that the m^{th} channel is not scheduled to start a new data transmission; and $a_m(t) > 0$ means to transmit data for the $a_m(t)^{\text{th}}$ device over the m^{th} channel. Then, denote the scheduling decision for all the channels at the t^{th} time slot as $\boldsymbol{a}(t) \triangleq [a_1(t), a_2(t), \cdots, a_M(t)]^T$, and apparently $\boldsymbol{a}(t)$ is the action of the considered Markov decision problem.

Remarkably, the actions in $\{0, 1, \dots, I + J\}^{M \times 1}$ are not always allowable. If the m^{th} channel has been reserved for data transmission in previous slots, i.e., $\sum_{j=1}^{J} b_{j,m}(t) > 0$ holds, it cannot be scheduled for any new data transmission. That is, the action is constrained by

$$\sum_{j=1}^{J} b_{j,m}(t) a_m(t) = 0, \ \forall \ m \in \{1, 2, \cdots, M\}.$$
 (2)

Then, we denote the allowable action space at state s(t), which contains all the actions satisfying the constraints in (2), as $\mathcal{A}_{s(t)}$.

3) Transitions: The transitions are to update AoII vector, channel gains, and channel availability.

AoII vector: The transitions of AoII vector needs the knowledge of the transitions of $\{X_i(t)\}_{i=1}^I$. Particularly, this paper considers that the state space of $\{X_i(t)\}$, denoted as \mathcal{X}_i , contains $|\mathcal{X}_i|$ real numbers, the self-transition probability of $\{X_i(t)\}$, denoted as $\Pr\{X_i(t+1) = x | X_i(t) = x\}$ with x being any state in \mathcal{X}_i , equals p_i , and the probability of transition to any other state, denoted as $\Pr\{X_i(t+1) = x | X_i(t) = y\}$ with $x \neq y$ and y being any other state in \mathcal{X}_i , equals $q_i \triangleq \frac{1-p_i}{|\mathcal{X}_i|-1}$.

Now, we study the transitions of $x_i(t)$ and first consider the case that we choose to transmit data for the i^{th} monitoring device over some channel at the t^{th} time slot, i.e., $\sum_{m=1}^{M} \mathbb{1}_i (a_m(t)) > 0$, where the indicator function $\mathbb{1}_i(x)$ equals 1 if x equals i and otherwise, it equals 0. For this case, state $X_i(t)$ is transmitted to the BS during the t^{th} time slot and the BS makes the prediction for the $(t+1)^{\text{th}}$ time slot by $\bar{X}_i(t+1) = X_i(t)$. Apparently, $\bar{X}_i(t+1) = X_i(t+1)$ holds with probability p_i since $X_i(t+1) = X_i(t)$ holds with probability p_i . Therefore, based on (1), when $\sum_{m=1}^{M} \mathbb{1}_i (a_m(t)) > 0$, we have

$$x_i(t+1) = \begin{cases} 0 & \text{with probability } p_i \\ x_i(t) + 1 & \text{with probability } 1 - p_i. \end{cases}$$
(3)

Then, we consider the other case that we choose not to transmit data for the i^{th} monitoring device over any channel, i.e., $\sum_{m=1}^{M} \mathbb{1}_i (a_m(t)) = 0$. For this case, the BS has to inherit its previous prediction, i.e., $\bar{X}_i(t+1) = \bar{X}_i(t)$, and $x_i(t)$ updates itself based on the following rules:

• When $x_i(t) = 0$, $\overline{X}_i(t+1) = X_i(t+1)$ holds with probability p_i , since $\overline{X}_i(t+1) = \overline{X}_i(t) = X_i(t)$ holds for sure and $X_i(t) = X_i(t+1)$ holds with probability p_i . Then, based on (1), when $\sum_{m=1}^M \mathbb{1}_i(a_m(t)) = 0$ and $x_i(t) = 0$, we have

$$x_i(t+1) = \begin{cases} 0 & \text{with probability } p_i \\ x_i(t) + 1 & \text{with probability } 1 - p_i, \end{cases}$$
(4)

• When $x_i(t) > 0$, $\overline{X}_i(t+1) = X_i(t+1)$ holds with probability q_i , since $\overline{X}_i(t+1) = \overline{X}_i(t)$ holds for sure and $\overline{X}_i(t) = X_i(t+1)$ holds with probability q_i . Then, based on (1), when $\sum_{m=1}^M \mathbb{1}_i (a_m(t)) = 0$ and $x_i(t) > 0$, we have

$$x_i(t+1) = \begin{cases} 0 & \text{with probability } q_i \\ x_i(t) + 1 & \text{with probability } 1 - q_i. \end{cases}$$
(5)

Channel gains: For the link between the BS and any traditional device, the channel coefficient of this link is modeled as a stationary ergodic process, and so is the channel gain. Particularly,

$$\Pr\{g_{j,m}(t+1) = g'|g_{j,m}(t) = g\} = \Pr\{g_{j,m}(1) = g'|g_{j,m}(0) = g\} \triangleq \Pr_{j,m}\{g'|g\}, \quad (6)$$

holds for all $t \in \mathbb{Z}^+$, $j \in \{1, 2, \dots, J\}$, $m \in \{1, 2, \dots, M\}$, and $g, g' \in \mathcal{G}$, where $\Pr_{j,m}\{g'|g\}$ is a constant and represents the probability for $g_{j,m}(t)$ transiting from g to g'. We also consider $\mathcal{G} \triangleq \{g_1, g_2, \dots, g_{|\mathcal{G}|}\}$ as a finite real number set with $0 < g_1 \leq g_2 \leq \dots \leq g_{|\mathcal{G}|}$.

Channel availability: We update the *channel availability* in the following four cases: If the m^{th} channel is currently transmitting data for the j^{th} traditional device, i.e., $b_{j,m}(t) > 0$, the remaining time for the release of the m^{th} channel decreases by one at the next time slot, i.e., $b_{j,m}(t+1) = b_{j,m}(t) - 1$; if the m^{th} channel is reserved by another traditional device, i.e., $b_{j,m}(t) = 0$ and $\sum_{j' \neq j} b_{j',m}(t) > 0$, $b_{j,m}(t+1)$ remains 0; if the m^{th} channel is currently available and is about to transmit data for the j^{th} traditional device at the t^{th} time slot, i.e., $\sum_{j=1}^{J} b_{j,m}(t) = 0$ and $a_m(t) = I + j$, the m^{th} channel will be released after $T_j - 1$ time slots counting from the $(t+1)^{\text{th}}$ time slot, i.e., $b_{j,m}(t+1) = T_j - 1$; finally, if the m^{th} channel is available and not going to transmit data for the j^{th} traditional device, i.e., $\sum_{j=1}^{J} b_{j,m}(t) = 0$ and $a_m(t) \neq I + j$, $b_{j,m}(t+1)$ remains 0. To summary, we have

$$b_{j,m}(t+1) = \begin{cases} b_{j,m}(t) - 1 & b_{j,m}(t) > 0, \\ 0 & b_{j,m}(t) = 0, \sum_{j' \neq j} b_{j',m}(t) > 0 \\ T_j - 1 & \sum_{j=1}^J b_{j,m}(t) = 0, a_m(t) = I+j \\ 0 & \sum_{j=1}^J b_{j,m}(t) = 0, a_m(t) \neq I+j. \end{cases}$$
(7)

4) Reward: The reward of the whole system consists of the AoIIs at I monitoring devices and the throughputs of J traditional devices. Specifically, the throughput of transmitting data for the j^{th} traditional device over the m^{th} channel at the t^{th} time slot is computed as

$$u_{j,m}(t) \triangleq W_m \log\left(1 + \frac{g_{j,m}(t)P}{N}\right),$$
 (8)

where W_m is the bandwidth of the m^{th} channel, P is the transmission power at the transmitters of the traditional devices, and N is the noise power at the receiver of the BS. Notably, the above throughput exists only if $\mathbb{1}_{I+j}(a_m(t))=1$ or $b_{j,m}(t)>0$ holds, when the m^{th} channel starts to transmit data or has been reserved for transmitting data for the j^{th} traditional device at the t^{th} time slot. The reward of the whole system is defined as the summation of all the AoIIs and throughputs at the t time slot, i.e.,

$$r(t) \triangleq -\sum_{i=1}^{I} w_{i} x_{i}(t) + \sum_{j=1}^{J} w_{I+j}$$

$$\sum_{m=1}^{M} \left(\mathbb{1}_{I+j} \left(a_{m}(t)\right) + \mathcal{G}\left(b_{j,m}(t)\right)\right) u_{j,m}(t),$$
(9)

where w_i , $i \in \{1, 2, \dots, I + J\}$, measures the importance of the *i*th device, and the indicator function $\mathcal{I}(x)$ equals 1 if x is positive and otherwise, it equals 0.

B. Problem formulation

1) Markov decision problem formulation

This paper aims to maximize the long-term average reward for (9). However, the objective $\max_{\{a(t)\}} \lim_{T\to\infty} \mathbb{E}_{\Pr\{x'|x,a\},\Pr\{G'|G\}} \left[\frac{1}{T} \sum_{t=1}^{T} r(t)\right]$ may not exist under some assignments of $\{a(t)\}$ (cf. [20, Example 8.1.1]), where the expectation is taken with respect to the *AoII vector* and the *channel gains*. Therefore, we utilize *liminf* average optimality criteria and formulate the joint scheduling problem between the task-oriented and data-oriented communications as

(P1)
$$\max_{\{a(t)\}} \quad \liminf_{T \to \infty} \mathbb{E}_{\Pr\{x' \mid x, a\}, \Pr\{G' \mid G\}} \left[\frac{1}{T} \sum_{t=1}^{T} r(t) \right]$$

s.t. (2), (3), (4), (5), (6), (7).

To solve problem (P1), we need to find the optimal policy $\pi \triangleq (\pi_1, \pi_2, \dots, \pi_t, \dots)$, where π_t is the optimal decision rule at the t^{th} time slot and maps the history of the states and actions $h(t) \triangleq (s(1), a(1), \dots, s(t-1), a(t-1), s(t))$ to the optimal distribution of the current action a(t), i.e., $\pi_t : \mathcal{H}(t) \times \mathcal{A}_{s(t)} \to [0, 1]$ with $\mathcal{H}(t)$ being the set of all histories h(t).

It can be checked that problem (**P1**) has upper-bounded rewards and countable states, and thus the stationary policies achieving its optimum exist if certain conditions are satisfied [21]. Then, we obtain the following proposition.

Proposition 2.1: There exist stationary policies to achieve the optimal value of problem (**P1**).

Based on Proposition 2.1, the *liminf* average optimality criteria for problem (**P1**) can be replaced by *lim*, since the Cesaro limit always exists for stationary policies [20]. Moreover, the state occurrence probabilities are constants under stationary policies [20], where the state occurrence probability of state s is defined as $\lim_{T\to\infty} \sum_{t=1}^{T} \mathbb{1}_s(s(t))/T$ [20]. Based on this property, we can simplify problem (**P1**) to an equivalent Markov decision problem (**P2**), which is given as

• State: $\hat{s}(t) \triangleq (\boldsymbol{x}(t), \boldsymbol{G}(t), \boldsymbol{b}(t))$, where $\boldsymbol{b}(t) \triangleq [b_1(t), b_2(t), \cdots, b_M(t)]^T \in \mathbb{Z}_{\geq 0}^{M \times 1}$ is defined as the *channel* release time. Here, $b_m(t)$ is equal to the number of the remaining time slots for the release of the m^{th} channel and defined as $b_m(t) \triangleq \sum_{j=1}^J b_{j,m}(t), m \in \{1, 2, \cdots, M\}$. Apparently, $\boldsymbol{b}(t)$ can be derived from $\boldsymbol{B}(t)$ and thus $\hat{\boldsymbol{s}}(t)$ can be derived from $\boldsymbol{s}(t)$. The new state space is

 $\hat{\mathcal{S}} \triangleq \mathbb{Z}_{\geq 0}^{I \times 1} \times \mathcal{G}^{J \times M} \times \mathcal{B}$, where \mathcal{B} is the value space of $\boldsymbol{b}(t)$;

- Action: a(t), which is constrained by (2);
- *Transitions:* (3), (4), (5), (6), and

$$b_m(t+1) = \begin{cases} b_m(t) - 1 & b_m(t) > 0\\ T_{a_m(t) - I} - 1 & b_m(t) = 0, a_m(t) > I\\ 0 & b_m(t) = 0, a_m(t) \le I; \end{cases}$$
(10)

where (10) is derived by (7) and $b_m(t) \triangleq \sum_{j=1}^J b_{j,m}(t)$. • *Reward:* $\hat{r}(t)$ is defined as

$$\hat{r}(t) \triangleq -\sum_{i=1}^{I} w_i x_i(t) + \sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbb{1}_{I+j} \left(a_m(t) \right) \bar{u}_{j,m}(t),$$
(11)

where, $\bar{u}_{j,m}(t)$ is defined as

$$\bar{u}_{j,m}(t) \triangleq \sum_{\tau=t}^{t+T_j-1} \mathbb{E}_{\Pr\{\boldsymbol{G}'|\boldsymbol{G}\}} \left[u_{j,m}(\tau) |_{g_{j,m}(t)} \right]; \quad (12)$$

• *Problem formulation:* with the *lim* average optimality criteria, problem (P2) is reformulated as

(P2)
$$\max_{\pi \in \Pi^{S}} \lim_{T \to \infty} \mathbb{E}_{\pi, \Pr\{\boldsymbol{x}' | \boldsymbol{x}, \boldsymbol{a}\}, \Pr\{\boldsymbol{G}' | \boldsymbol{G}\}} \left[\frac{1}{T} \sum_{t=1}^{T} \hat{r}(t) \right]$$
(13)
s.t. (2), (3), (4), (5), (6), (10),

where Π^{S} is the set of all stationary policies and the expectation is taken with respect to the policy π , the *AoII vector*, and the *channel gains*.

Apparently, problem (P2) has a much simpler structure than problem (P1), where the dimension of the state space is reduced from I + 2JM to I + JM + M, the transitions in (10) evolves much simpler than the ones in (7), and the reward in $\hat{r}(t)$ no longer involves the high-dimension matrix B(t) compared with the reward in r(t). Moreover, compare problems (P1) and (P2), we obtain the following proposition.

Proposition 2.2: Problem (P2) is equivalent to problem (P1). Particularly, for any stationary policy optimizing problem (P1), there exists another stationary policy that optimizes problem (P2), and the inverse holds, too. Moreover, the optimal values of the two problems are the same.

Proof: Please see Appendix B.

However, problem (P2) is still difficult to be solved since it is a Markov decision problem with the *lim* average optimality criteria, and the existing algorithms addressing such problems, e.g., relative value iteration algorithm [22], still suffer from the curse of dimensionality. To tackle this issue, we refer [23, Proposition 4.1.7] and verify the existence of the Blackwell policies in problem (P2). If the Blackwell policies do exist, the optimal policies of problem (P2) can be found by solving a discounted version of it, which is much easier to be addressed. Particularly, the discounted version of problem (P2) is given 6

as

(P3)
$$\max_{\pi \in \Pi^{S}} \lim_{T \to \infty} \mathbb{E}_{\pi, \Pr\{\boldsymbol{x}' | \boldsymbol{x}, \boldsymbol{a}\}, \Pr\{\boldsymbol{G}' | \boldsymbol{G}\}} \left[\sum_{t=1}^{T} \alpha^{t-1} \hat{r}(t) \right]$$
(14)
s.t. (2), (3), (4), (5), (6), (10),

where $\alpha \in (0, 1)$ is the discount factor. Then, we derive the relationship between problems (P2) and (P3) in the following proposition.

Proposition 2.3: When the discount factor $\alpha \in (0,1)$ is sufficiently close to 1, there exist stationary policies to simultaneously achieve the optimal values of problems (P2) and (P3).

Proof: Please see Appendix C.

Remark 2.1: With Proposition 2.3, we can solve problem (P2) by first selecting a proper discount factor α and then deriving a stationary optimal policy of problem (P3) with the chosen discount factor. However, how to solve problem (P3) is still challenging due to its three features: 1) large state space; 2) large and discrete action space; and 3) multiple time-varying action constraints in (2). However, conventional approaches including dynamic programming [22] and Lyapunov drift optimization [10], can barely deal with Markov decision problems with the first two features. Modern DRL algorithms [24] can neither efficiently address the Markov decision problems with large and discrete action space¹, and most of them deploy the trial and error mechanism [24], where the trial part would always violate the time-varying constraints in (2) and thus the training procedure would be terminated.

2) Markov game formulation

To overcome the three challenges in problem (P3), we treat each channel as one virtual agent, which can first observe the system state and then independently determine the device for data transmission over itself. Therefore, problem (P3) can be reformulated as an equivalent Markov game (P4), which has the same state, transition, and optimality criteria as those of problem (P3) and also contains

- Observation $\mathbf{s}_m(t) \triangleq (\mathbf{x}(t), \mathbf{g}_m(t), b_m(t))$ at the m^{th} agent: $\mathbf{g}_m(t)$ contains the channel gains for transmitting data for J traditional devices over the m^{th} channel, respectively, and is defined by $\mathbf{g}_m(t) \triangleq [g_{1,m}(t), g_{2,m}(t), \cdots, g_{J,m}(t)]^T$;
- Action $a_m(t)$ at the m^{th} agent: it is now individually constrained by $b_m(t)a_m(t) = 0$;
- Reward $r_m(t)$ for the m^{th} agent: it is set to be $\hat{r}(t)$ in (11).

Remark 2.2: On the design of the agent observation, although the whole state $\hat{s}(t)$ is observable for each agent, only $s_m(t)$ is reserved as the m^{th} agent's observation. The reasons are: 1) the dimensionality of $\hat{s}(t)$ is too large and difficult to be dealt with; 2) except the elements in $s_m(t)$, the rest ones in $\hat{s}(t)$ are weekly correlated to the m^{th} agent; 3) the union of $s_m(t)$ among all agents covers all the elements in state $\hat{s}(t)$ and thus the algorithms developed based on $s_m(t)$ can achieve the same optimum with the algorithms based on $\hat{s}(t)$ (see Fig. 4 in Section IV).

Markov game (P4) is obviously equivalent to problem (P3) since all its agents cooperatively optimize the common reward $\hat{r}(t)$ in problem (P3). Meanwhile, it faces similar challenges to problem (P3). Fortunately, the existing multiagent reinforcement learning (MARL) algorithms can well address the challenge from the large state space issue by approximating the optimal policy with neural network (NN), and relieve the challenge from the large discrete action space issue by deploying decentralized learning mechanism, where each agent needs only to determine its own action in a much smaller action space $\{0, 1, \dots, I + J\}$.

However, there are still two obstacles for MARL algorithms to solve (**P4**): 1) the action space for each agent, i.e., $\{0, 1, \dots, I + J\}$, is still large and would slow down the training procedure for MARL algorithms to a great extent (see Fig. 4 in Section IV); 2) the action space for each agent is constrained by one time-varying constraint, which again collides with the trial and error mechanism deployed in MARL algorithms;

III. WHITTLE'S INDEX GUIDED MULTI-AGENT PROXIMAL POLICY OPTIMIZATION

To address Markov game (P4), we propose WI-MAPPO, which mainly comprises one Whittle's index guided action fusion (WIAC) module and multiple proximal policy optimization (PPO) modules. Particularly, the WIAC module calculates the Whittle's index for I monitoring devices based on their AoIIs, and accordingly determines the priorities for their data transmissions. Then, all agents only need to transmit data for the group of monitoring devices with the highest priority, and thus the action space for each individual agent is greatly shrunk. Meanwhile, we modify both the actor network and the probability ratio derivation of the training algorithm for multi-agent proximal policy optimization (MAPPO) to train the proposed WI-MAPPO, which perfectly addresses the timevarying constraint issue.

In the following, we first introduce the structure of the proposed algorithm. Then, we present the offline training algorithm in details. Finally, we briefly introduce the online applying algorithms for solving Markov game (**P4**).

A. Structure of proposed algorithm

As illustrated in Fig. 2, the main body of the proposed algorithm consists of one observation derivation (OD) module, M PPOs, and one WIAC module.

1) OD module: This module derives the agent observations $s_1(t), s_2(t), \dots, s_M(t)$ from the state $\hat{s}(t)$ based on the definition $s_m(t) \triangleq (x(t), g_m(t), b_m(t))$.

¹The most advanced algorithm to address Markov decision problem with large discrete action space is the Wolpertinger policy [25], which adds an action-embedding module right after the deep deterministic policy gradient (DDPG ^[26]) algorithm and directly discretizes the continuous-valued action generated by DDPG. However, the Wolpertinger policy has a poor interpretability and could generate really large training variance even when the action space is small [25].



Figure 2: Structure of the proposed WI-MAPPO algorithm.

2) *PPOs:* Each agent utilizes a PPO module to determine the device for data transmission over its channel. Specifically, The PPO utilized by the m^{th} agent is named as PPO_m and it has a simple structure:

 Actor network: It contains a fully connected NN parameterized by θ_1 . Particularly, this NN takes $s_m(t)$ as the input and thus has I + J + 1 nodes at the input layer; the output layer has J + 2 nodes and represents the probabilities of transmitting data for all Jtraditional devices, transmitting data for one monitoring device, and not starting new transmission, respectively, which are denoted by $\pi_{m,\theta_1}(1|s_m(t)), \pi_{m,\theta_1}(2|s_m(t)),$ \cdots , and $\pi_{m,\theta_1}(J+2|\boldsymbol{s}_m(t))$. The output of the actor network, which is named as the PPO action, is denoted as $\bar{a}_m(t) \in \{1, 2, \dots, J+2\}$. Particularly, when the m^{th} channel is currently available, i.e., $b_m(t) = 0$, $\bar{a}_m(t)$ is equal to a discrete random variable X, whose probability mass function is $\Pr\{X = j\} = \pi_{m,\theta_1}(j|s_m(t)), \forall j \in$ $\{1, 2, \dots, J+2\}$. And when the m^{th} channel is occupied, i.e., $b_m(t) > 0$, $\bar{a}_m(t)$ is equal to J+2. To summary, we have

$$\bar{a}_m(t) = \begin{cases} J+2 & b_m(t) > 0\\ X & b_m(t) = 0. \end{cases}$$
(15)

Remarkably, trivial action selection in the actor network of PPO simply follows the second line in (15) [27]. While in our design, we introduce the first line of (15) to manually change the action $\bar{a}_m(t)$ to not starting new transmission, i.e., $\bar{a}_m(t) = J + 2$, when the m^{th} channel is currently occupied. Such modification will benefit the offline training procedure discussed in III-B and we will explain the reasons in Remark 3.2.

• Critic network: It contains a fully connected NN parameterized by θ_2 . Specifically, this NN takes $\hat{s}(t)$ as the input and thus has I + JM + M nodes at the input layer; the output layer of the NN has only one node and gives the estimation on the maximum total discounted reward starting from state $\hat{s}(t)$, which is also called as the value function of state $\hat{s}(t)$ and denoted by $V_{m,\theta_2}(\hat{s}(t))$;

• Experience buffer: It stores the experiences generated in the offline training procedure, where the experiences are the five-component tuples $(\hat{s}(t), s_m(t), \bar{a}_m(t), \pi_{m,\theta_1}(\bar{a}_m(t)|s_m(t)), r_m(t))$.

Remark 3.1: On the design of the actor network, it is reasonable to set the output dimension of the actor network's NN as I + J + 1 and let the output values represent the probabilities of transmitting data for all I + J devices, and not starting new transmission, respectively, at state $s_m(t)$. By doing so, the mth agent can directly determine its action by sampling from the action distribution generated by the trained actor network's NN. However, high dimensionality of the output, which is equal to I + J + 1, requires large NN, and the convergence of the training procedure among multiple agents might be very difficult. In the proposed algorithm, we design the WIAC module to figure out the group of monitoring device with the highest potential to minimize their long-term average AoIIs, by which each agent only needs to determine whether it is willing to transmit data for one monitoring device. If it is, the agent will select any one in the group figured out by the WIAC module for data transmission. Therefore, the actor network in the proposed algorithm only requires an output with J + 2 dimensions.

3) WIAC module: This module determines the action a(t) based on state $\hat{s}(t)$ and M PPO actions $\bar{a}(t) \triangleq [\bar{a}_1(t), \bar{a}_2(t), \dots, \bar{a}_M(t)]^T$. In this subsection, we first introduce the design intuition of this module, and then introduce the explicit method to compute a(t).

Design intuition: WIAC module first solves the following problem (**P5**): Given the current AoII values at I monitoring devices, which group of the monitoring devices should be selected for status update transmissions over limited number of channels to minimize the long-term average weighted AoIIs? Problem (**P5**) is a restless multi-armed bandit (RMAB) prob-

lem and a typical Whittle's index algorithm is deployed to solve it [8]. Particularly, we first model I monitoring devices as I individual agents and the goal of each agent is to minimize its own long-term average AoII. Then, by studying the AoII evolution of each agent, the maximum offer that each agent is willing to pay for hiring one channel at the current time slot can be derived, which is named as the Whittle's index for this agent. Finally, the group of agents with the highest Whittle's indices will be selected for data transmissions. Remarkably, the Whittle's index algorithm is validated to be a near-optimal algorithm to solve RMABs [8].

However, to apply this algorithm, the existence of the Whittle's index should be guaranteed. Therefore, in the following, we first validate the existence of the Whittle's index for problem (**P5**), and then derive a Whittle's index table with an exhausted search algorithm. Finally, to derive the final action a(t), we first check the values of $\bar{a}(t)$ and obtain the number of channels willing to transmit data for one monitoring device, i.e., $\sum_{m=1}^{M} \mathbb{1}_{J+1}(\bar{a}_m(t))$. Then, we look up the Whittle's index table and transmit data for the group of the monitoring devices with the highest Whittle's indices over these channels, where the size of this group is also equal to $\sum_{m=1}^{M} \mathbb{1}_{J+1}(\bar{a}_m(t))$.

Existence of Whittle's index: The following proposition validates the existence of the Whittle's index for problem (**P5**).

Proposition 3.1: There exists Whittle's index for problem (P5).

Sketch of proof: To validate the existence of Whittle's index, we first decouple problem (**P5**) to I sub-problems, where the i^{th} sub-problem is to minimize the average AoII at the i^{th} monitoring device. Next, we analyze the properties of these sub-problems and validate that the optimal policies for these sub-problems are of threshold type. Then, based on the "threshold" feature on the optimal policy, we prove the indexability for the decoupled sub-problems, which validates the existence of the Whittle's index for problem (**P5**). Please check Appendix D for more details.

Derivation of Whittle's index table: The Whittle's index for the i^{th} monitoring device with its AoII being x, notated as $I_i(x)$, is defined as the additional cost C that makes both transmitting data and not transmitting data for the i^{th} monitoring device equally desirable, i.e.,

$$f_i(x, C) = f_i(x+1, C),$$
 (16)

where the additional cost C and the average cost function $f_i(x, C)$ are introduced in Appendix D. Since to derive the closed-form formulation of $I_i(x)$ by solving (16) is very difficult, we use exhausted search algorithm to obtain $I_i(x)$ with the searching step and searching area as Δc and $[C_L, C_U]$, respectively. Moreover, by exploiting the fact that $I_i(x)$ is non-decreasing with respect to x [8], the above exhausted searching algorithm can be improved. Notably, we would generate a sufficiently large Whittle's index table by executing this algorithm before the offline training procedure.

Derivation of a(t): During the offline training procedure, WIAC module first counts the number of agents willing to transmit data for monitoring devices, i.e., the agents satisfying $\bar{a}_m(t) = J + 1, m \in \{1, 2, \dots, M\}$. Particularly, this number at the t^{th} time slot is denoted as A(t) and defined as $A(t) \triangleq \sum_{m=1}^{M} \mathbbm{1}_{J+1}(\bar{a}_m(t))$. Next, WIAC module looks up the generated Whittle's index table and obtains the Whittle's indices for I monitoring devices according to their current AoII values $\boldsymbol{x}(t)$. Then, the index of the monitoring device with the l^{th} highest Whittle's index is denoted as $W_l(t) \in \{1, 2, \dots, I\}$ and A(t) monitoring devices with the top A(t) Whittle's indices are picked out. Finally, the A(t) agents satisfying $\bar{a}_m(t) = J + 1$, $m \in \{1, 2, \dots, M\}$ transmit data for the picked out A(t) monitoring devices over their channels and accordingly the action $\boldsymbol{a}(t) = (a_1(t), a_2(t), \dots, a_M(t))$ is computed as

$$a_m(t) = \begin{cases} I + \bar{a}_m(t) & \bar{a}_m(t) < J + 1 \\ W_{l_m}(t) & \bar{a}_m(t) = J + 1 \\ 0 & \bar{a}_m(t) = J + 2, \end{cases}$$
(17)

where l_m is the number of elements equaling J + 1 in $[\bar{a}_1(t), \bar{a}_2(t), \cdots, \bar{a}_m(t)]^T$.

B. Offline training

Based on the historical observed samples, we can easily approximate the values of $Pr\{G'|G\}$, p_i , and q_i , and then simulate an offline environment accordingly. Finally, we develop the offline training algorithm by interacting with it.

1) Offline environment simulation: To mimic the real environment, the offline environment needs to fulfill two functions:

- State evolution: Given $\hat{s}(t) = (x(t), G(t), b(t))$ and a(t), we first simulate x(t+1) based on (3), (4), and (5), and simulate G(t+1) based on (6) and the approximated $\Pr\{G'|G\}, p_i$, and q_i . Then, we directly compute b(t+1)based on (10). Thus, $\hat{s}(t+1)$ is obtained;
- Reward generation: Given $\hat{s}(t)$ and a(t), we compute $r_1(t), r_2(t), \dots, r_M(t)$ based on (11), (12) and the fact $r_m(t) = \hat{r}(t)$.

2) Offline training: As illustrated in the right part of Fig. 2, we alternatingly generate experiences by deploying the latest actor-critic network and update the actor-critic network according to the latest generated experiences. We specific these two steps as follows.

• Generation of experiences: First, obtain the $s_1(t)$, \cdots , $s_M(t)$ from the observed $\hat{s}(t)$ based on OD module. Next, by utilizing the actor networks for PPOs, obtain PPO action $\bar{a}(t)$ and the value of $\pi_{m,\theta_1}(\bar{a}_m(t)|s_m(t))$ for all $m \in \{1, 2, \ldots, M\}$ according to (15). Then, with the known $\bar{a}(t)$ and $\hat{s}(t)$, derive action a(t) by using the WIAC module. Finally, obtain the agent rewards $r_1(t), \cdots, r_M(t)$ by interacting with the offline environment. We pack the above information as experiences $e_1(t), e_2(t), \cdots, e_M(t)$, where $e_m(t)$ is defined as

$$e_m(t) \triangleq (\hat{\boldsymbol{s}}(t), \boldsymbol{s}_m(t), \bar{a}_m(t), \pi_{m,\boldsymbol{\theta}_1}(\bar{a}_m(t)|\boldsymbol{s}_m(t)), r_m(t)),$$
(18)

and then we store $e_m(t)$ in the experience buffer in PPO_m. Remarkably, we can continuously generate N_B experiences for each PPO before the update of the actorcritic networks and empty all the experience buffers after each update. • Updation of actor-critic networks: Each updation performs N_U epochs of optimization on the generated N_B experiences and each epoch would modify the parameters of the actor-critic networks for all PPOs. Particularly, at the beginning of each epoch, we first denote the actor and critic networks for PPO_m at this moment as π_{m,θ'_1} and V_{m,θ'_2} , respectively, Next, we estimate N_B value functions as

$$V_m(t, \hat{s}(t)) = r_m(t) + \alpha r_m(t+1) + \dots + \alpha^{N_B} r_m(N_B)$$
(19)

with $t \in \{1, 2, \dots, N_B\}$, where $\{V_m(t, \hat{s}(t))\}_{t=1}^{N_B}$ are the value functions of states $\{\hat{s}(t)\}_{t=1}^{N_B}$ for PPO_m [27]. Then, we utilize the current critic network V_{m,θ'_2} to estimate N_B advantage functions as

$$A_m(t) = V_m(t, \hat{s}(t)) - V_{m, \theta'_2}(\hat{s}(t)), t \in \{1, 2, \cdots, N_B\}, (20)$$

where $\{A_m(t)\}_{t=1}^{N_B}$ are the advantage functions for PPO_m [27], and utilize the current actor network π_{m,θ'_1} to derive N_B probability ratios as

$$R_m(t) = \begin{cases} 1 & b_m(t) > 0\\ \frac{\pi_{m,\theta_1'}(\bar{a}_m(t)|\boldsymbol{s}_m(t))}{\pi_{m,\theta_1}(\bar{a}_m(t)|\boldsymbol{s}_m(t))} & b_m(t) = 0, \end{cases}$$
(21)

for all $t \in \{1, 2, \dots, N_B\}$. Finally, the surrogate loss for PPO_m is computed as [27]

$$L_{m} = \sum_{t=1}^{N_{B}} \frac{1}{N_{B}} \Big(-\min(R_{m}(t)A_{m}(t), \\ \operatorname{clip}(R_{m}(t), 1 - \epsilon, 1 + \epsilon)A_{m}(t)) \Big)$$
(22)

$$+ c_1 (V_m(t, \hat{s}(t)) - V_{m, \theta'_2}(\hat{s}(t)))^2$$

$$-c_2H\left(\pi_{m,\boldsymbol{\theta}_1'}(\cdot|\boldsymbol{s}_m(t))\right)\right),$$

where $\operatorname{clip}(x, a, b) \triangleq \min(\max(x, a), b)$ clamps x into the area [a, b]; $H(\pi_{m, \theta'_1}(\cdot|s_m(t))$ is the entropy of the stochastic output generated by the current actor network π_{m, θ'_1} with $s_m(t)$ as input; ϵ , c_1 , and c_2 are some constants. Particularly, the first term in (22) is a pessimistic bound, which could improve the actor-critic networks for PPO in a considerably stable manner; the second MSE term is essential for the convergence of the training of the actor-critic network [27]; and the last term adopts an entropy bonus to ensure sufficient exploration. Remarkably, both the actor and critic networks backpropagate this surrogate loss to update their parameters θ'_1 and θ'_2 .

The details of the offline training algorithm are summarized in Algorithm 1.

Remark 3.2: Compared with the conventional training algorithm for MAPPO, the proposed training algorithm modifies the actor network for each PPO in (15) and the probability ratio derivation in (21). Such modifications have two advantages:

• The proposed training algorithm will not violate the timevarying constraints in (2) nor terminate the training procedure. The conventional training algorithm for MAPPO selects action by sampling from the action distribution generated by the actor network, and thus may select

Algorithm 1 Offline training algorithm for joint scheduling

- 1: Randomly initialize the actor-critic networks PPO_1 , PPO_2 , \cdots , and PPO_M ;
- 2: Initialize one experience buffer for each PPO;
- 3: Input the values of $\Pr\{G'|G\}, \{p_i\}_{i=1}^{I}, \{q_i\}_{i=1}^{I}, \{w_i\}_{i=1}^{I+J}, \{W_m\}_{m=1}^{M}, P/N, \alpha, \Delta c, C_L, C_U, N_B, N_U, \epsilon, c_1, c_2;$
- 4: Derive the Whittle's index table by executing the exhausted search algorithm specified in III-A3;
- 5: Generate the offline environment based on III-B1;
- 6: **for** episode $= 1, 2, \cdots$
- 7: Let $\boldsymbol{x}(1) = 0^{I \times 1}$ and $\boldsymbol{b}(1) = 0^{M \times 1}$. Let $\boldsymbol{G}(1)$ be any element in $\mathcal{G}^{J \times M}$;

8:
$$\hat{\boldsymbol{s}}(1) = (\boldsymbol{x}(1), \boldsymbol{G}(1), \boldsymbol{b}(1));$$

- 9: **for** $t = 1, 2, \dots, N_B$
- 10: Send $\hat{s}(t)$ to OD module and derive $\{s_m(t)\}_{m=1}^M$;
- 11: **for** $m = 1, 2, \dots, M$ 12: Send $s_m(t)$ to the
 - Send $s_m(t)$ to the actor Network for PPO_m and derive $\bar{a}_m(t)$, $\pi_{m,\theta_1}(\bar{a}_m(t)|s_m(t))$;

13: end for

14:

16:

17:

22:

23:

Send $\hat{s}(t)$ and $\bar{a}(t)$ to the WIAC module and derive a(t);

15: Send
$$\hat{s}(t)$$
 and $a(t)$ to the offline environment and derive $\hat{s}(t+1), r_1(t), \dots, r_M(t)$;

for
$$m = 1, 2, \dots, M$$

Pack experience $e_m(t)$ by (18) and store it into the experience buffer in PPO_m ;

18: **end for**

19: **end for**

- 20: for epoch = $1, 2, \cdots, N_U$
- 21: **for** $m = 1, 2, \cdots, M$
 - Load the experiences $\{e_m(t)\}_{t=1}^{N_B}$ from the experience buffer in PPO_m; Derive $\{V_m(t \ \hat{s}(t))\}_{m=1}^{N_B}$ $\{A_m(t)\}_{m=1}^{N_B}$

Derive
$$\{V_m(t, \mathbf{s}(t))\}_{t=1}^{T}, \{A_m(t)\}_{t=1}^{T}, \{P_m(t)\}_{t=1}^{NB}$$
 based on (10) (20) (21)

24: $\{R_m(t)\}_{t=1}^{N_B}$ based on (19), (20), (21); Derive L_m based on (22);

- 25: Update the actor-critic networks for PPO_m by backpropagating L_m ;
- 26: end for
- 27: end for
- 28: Empty the experience buffers for PPOs;
- 29: end for

infeasible ones. However, based on the modification in (15), the proposed algorithm forces the agent not to start new data transmission when its channel is occupied, by which the generated action never violates the constraints in (2);

• The proposed training algorithm perfectly extends the conventional training algorithm for MAPPO to solve Markov games with time-varying constraints. Particularly, the principal component for the surrogate loss of the conventional MAPPO utilizes both the advantage functions in (20) and the probability ratios in (21). The former ones are directly derived from the value functions, which are estimated based on the generated trajectory in (19). Now, with the modification in (15), the agent

action selected at the states satisfying $b_m(t) > 0$ in the generated trajectory is the optimal agent action since no other agent action is allowable at these states. Thus, based on the trajectory generated in this way, the value functions and advantage functions can be estimated faster and more accurately. The latter ones, i.e., the probability ratios, are modified by (21) and the reason for this modification is quite straightforward: when encountering the states satisfying $b_m(t) > 0$, the modified actor network by (15) always selects J + 2 as the action. Thus, the probability ratio at these states is equal to $\frac{1}{1} = 1$.

C. Online applying

The online applying algorithm is very similar to the offline one while only uses the trained actor networks for PPOs. Moreover, the values of x(t+1) and G(t+1) can only be derived from the online interactions with the real environment. Thus, we omit the details.

IV. NUMERICAL RESULTS

This section evaluates the performance of the proposed algorithm and compares it with the stat-of-the-art AoI-based algorithms. Specifically, we consider a IoT system with M =10 channels collecting data from I = 90 monitoring devices and J = 10 traditional devices, where the corresponding action space has a magnitude of 101¹⁰. Each monitoring device monitors one random process, where each random process has $|\mathcal{X}_i| = 10$ states. Meanwhile, the self-transition probabilities $\{p_i\}_{i=1}^{90}$ of these random processes satisfy $p_i=0.6,\ 1\leq$ $i \leq 60$ and $p_i = 0.9, \ 61 \leq i \leq 90$. The consumed time duration T_j for each data transmission of the j^{th} traditional device is uniformly picked from $\mathcal{T} \triangleq \{1, 2, \dots, 10\}$, i.e., $\boldsymbol{T} \triangleq [T_1, T_2, \cdots, T_J]^T \in \mathcal{T}^{J \times 1}$. The channel gain model refers [28], where each channel gain $g_{j,m}(t)$ takes value in $\{\bar{g}_{j,m}, \bar{g}_{j,m} + 1, \cdots, \bar{g}_{j,m} + 9\}, \bar{g}_{j,m}$ is uniformly picked from $\{0, 1, \dots, 40\}$, and $g_{i,m}(t)$ transits to the current value with probability 0.6 and to two adjacent values with equal probability 0.2. The bandwidths for all channels are set as $W_m = 1$ and the importance weights for all devices $\{w_i\}_{i=1}^{I+J}$ are uniformly picked from $\{1, 2\}$. Other parameters are set as P/N = 1, $\Delta c = 0.1$, $C_L = 0.1$, $C_U = 4000$, $\Delta C = 0.1$, $N_B = 4000$, $N_U = 80$, $\epsilon = 0.2$, $c_1 = 0.5$, and $c_2 = 0.01$. Moreover, both the actors and critics in WI-MAPPO utilize two hidden layers, each of which has 128 nodes. We compare the proposed WI-MAPPO with two AoI-based algorithms. The first one is AoI-based WI-MAPPO, which utilizes the same structure with WI-MAPPO, while the employed WIAC module is designed for minimizing AoI according to the method in [7]; The other algorithm, namely age-aware policy (AAP), utilizes the Lyapunov drift optimization and is currently the state-ofthe-art algorithm for joint schedules [13]. Remarkably, AAP cannot address the time-varying constraints issue and thus can only be applied in the scenario where the data transmissions for all traditional devices consume only 1 time slot, i.e., the scenario with $T = 1^{J \times 1}$, where $1^{J \times 1}$ is the J-by-1 vector with all entries as 1.



Figure 3: Performance comparisons between WI-MAPPOs with different discount factors.



Figure 4: Performance comparisons between WI-MAPPO and other MARL algorithms.

In Fig. 3, we investigate the performance of the proposed WI-MAPPO with different discount factors and approximate the value range of the discount factor satisfying the statement in Proposition 2.3. Specifically, it is observed that when α is no smaller than 0.8, WI-MAPPO achieves the same maximum on the average reward. This indicates that [0.8, 1] could be a proper value range as aforementioned. Meanwhile, it is observed that a too large discount factor, e.g., $\alpha = 0.95$, would slow down the convergence. Therefore, we select α as 0.9 for all the following experiences.

In Fig. 4, we validate the performance advantages of the WI-MAPPO over other MARL algorithms. The first algorithm we concerned is the conventional MAPPO. Compared with WI-MAPPO, MAPPO does not have the action space shrinkage provided by the WIAC module and thus suffers from a much larger output dimensionality on its actors, which equals 101. We test two MAPPO algorithms with separately 128 and 256 nodes on their hidden layers, and it is observed that both of them have far too slow convergence speeds, which also val-



Figure 5: Average throughput and average accuracy tradeoff for WI-MAPPOs and AAP in the scenario with 10 channels and 100 devices: (a) tradeoff comparisons between WI-MAPPOs and AAP in the non-constrained case; (b) tradeoff comparisons between WI-MAPPOs in the non-constrained and constrained cases.



Figure 6: Average throughput and average accuracy tradeoff for WI-MAPPOs and AAP in the scenario with 10 channels and 40 devices: (a) tradeoff comparisons between WI-MAPPOs and AAP in the non-constrained case; (b) tradeoff comparisons between WI-MAPPOs in the non-constrained and constrained cases.

idates Remark 3.1. The second algorithm simply extends the agent observation from $s_m(t)$ to $\hat{s}(t)$ and then designs a new WI-MAPPO accordingly. Certainly, this modified WI-MAPPO contains larger PPO actor networks and consumes more offline training and online applying computational resources than the original WI-MAPPO. Moreover, it is observed that it has bare extra gain over the original WI-MAPPO, which also validates Remark 2.2.

In Fig. 5(a), we simulate the scenario satisfying $T = 1^{T \times 1}$, and compare the performances of WI-MAPPO, AoI-based WI-MAPPO, and AAP. It is observed that when the throughput of the traditional devices is not important, the two AoIbased algorithms achieve almost the same average accuracy on predicting the monitored processes, which is around 0.37, and the original WI-MAPPO, which is AoII-based, performs much better and is around 0.42. This validates the advantage of AoII over AoI in pure task-oriented communications. Meanwhile, if we concern the throughput more, both AoII-based and AoI-based WI-MAPPOs gain much larger throughput than AAP when their achieved average accuracies are equal. And AoII-based WI-MAPPO greatly outperforms AoI-based WI-MAPPO. The reason for such performance advantages is also straightforward: AoI captures the aging of sampled status updates, while AoII further factors the semantic of status updates relative to prediction of real-time status at data source. Therefore, AoII-oriented scheduling reduces the required data traffic for monitoring devices to guarantee a certain level of prediction performance. In Fig. 5(b), we validate the ability of WI-MAPPO on handling the time-varying constraints by

comparing its performances in the non-constrained case, i.e., $T = 1^{J \times 1}$, and the constrained case, i.e., $T \in \mathcal{T}^{J \times 1}$. Remarkably, in the constrained case, the solution space is much smaller than the non-constrained case. Consequently, the optimal performance that any algorithm can achieve in the constrained case is also supposed to be worse than the non-constrained case. However, it is observed that both AoII-based and AoI-based WI-MAPPOs have almost the same performances in these two cases. This amazing result indicates that WI-MAPPO perfectly addresses the time-varying constraints issue for the considered Markov games.

In Fig. 6, we simulate a simple scenario with 30 monitoring devices and 10 traditional devices. Since there are only 40 devices requesting for data transmissions, both AoII-based and AoI-based algorithms could achieve the maximum average accuracy, which is around 0.685, when the throughput is not important. Moreover, WI-MAPPO algorithms again greatly outperform AAP. Remarkably, it is observed that AoII-based WI-MAPPO slightly outperforms the AoI-based WI-MAPPO in both the non-constrained and constrained cases. This indicates that the advantages of AoII-based algorithms over AoI-based ones would be more significant when the channel resources are not sufficient.

V. CONCLUSIONS

In this work, we study the joint schedule of task-oriented and data-oriented communications and formulate this problem as a challenging Markov decision problem. Insightful techniques and innovative algorithm are utilized to solve this problem as efficiently as possible. Specifically, to simplify this problem, we analysis its "stationary" feature and Blackwell policies and redesign the reward function based on the channel statistics, by which the solution space and state space are greatly shrunk in an equivalent manner and the optimality criteria is equivalently replaced by a discounted one. To overcome the large discrete action space issue, we convert this problem to an equivalent Markov decision game, where the original action for Markov decision problem is decomposed into low-dimension agent actions. Then, we validate the existence of Whittle's index and design a Whittle's index guided module to further shrink the action space. To overcome the time-varying action constraints issue, we modify the advantage function estimation kernel for MAPPO and extend the training algorithm to solve the constrained Markov games.

APPENDIX A PROOF OF PROPOSITION 2.1

It can be easily checked that: problem (**P1**) has infinite and countable states; and the reward of problem (**P1**) satisfies

$$r(t) \le \sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} W_m \log\left(1 + \frac{g_{|G|}P}{N}\right)$$
(23)

and thus it is upper bounded. Based on the above two properties, problem (**P1**) has stationary optimal policies if the following two conditions are satisfied [21, Proposition 5 & Theorem 1]:

1) Problem (**P1**) has a stationary policy which induces an ergodic Markov chain and has a finite average reward.

2) Define $V_{\alpha}(s) \triangleq \sup_{\theta} \lim_{T \to \infty} \mathbb{E}_{\pi_{\theta}, \Pr\{s' \mid s, a\}} \left[\sum_{t=1}^{T} \alpha^{t-1} r(t) \Big|_{s(1)=s} \right]$, where $\alpha \in (0, 1)$ is a discount factor; θ parameterizes policy π_{θ} ; the expectation is taken with respect to policy π_{θ} and state s(t); $V_{\alpha}(s)$ is the maximum total discounted reward that can be achieved by any policy starting from state s. Then, there exists a non-negative real number C such that $V_{\alpha}(s) - V_{\alpha}(s_0) \leq C$ holds for all state $s \in \delta$ and all $\alpha \in (0, 1)$, where s_0 is a reference state in δ .

To validate the first condition, we investigate the do-nothing policy π , where no device would be selected for data transmission over any channel or at any time slot. Apparently, this policy is stationary. In the following, we first show that the average reward under policy π is finite, and then show that the induced Markov chain is ergodic.

As for the average reward under policy π , it equals $\lim \inf_{T \to \infty} \mathbb{E}_{\pi,\Pr\{\boldsymbol{x}' \mid \boldsymbol{x}, \boldsymbol{a}\}} \left[\frac{1}{T} \sum_{t=1}^{T} \left(-\sum_{i=1}^{I} w_i x_i(t) \right) \right]$. Specifically, the process $\{x_i(t)\}$ under policy π forms an ergodic Markov chain, the transition of which is specified by equations (4), (5). And accordingly, the transition equations are $\mu_1 = (1 - p_i)\mu_0$ and $\mu_x = (1 - q_i)\mu_{x-1}, x = 2, 3, \cdots$, where μ_x is the state occurrence probability of the state x. Then, we can solve that $\mu_0 = \frac{q_i}{1+q_i-p_i}$ and $\mu_x = \frac{q_i(1-p_i)(1-q_i)^{x-1}}{1+q_i-p_i}, \forall x \in \mathbb{Z}^+$. And it follows

$$\liminf_{T \to \infty} \mathbb{E}_{\pi, \Pr\{\boldsymbol{x}' | \boldsymbol{x}, \boldsymbol{a}\}} \left[\frac{1}{T} \sum_{t=1}^{T} \left(-\sum_{i=1}^{I} w_i x_i(t) \right) \right]$$
$$= -\sum_{i=1}^{I} w_i \sum_{x=1}^{\infty} x \mu_x = -\sum_{i=1}^{I} w_i \frac{1-p_i}{(1+q_i-p_i)q_i} < \infty,$$

which implies that the average reward under policy π is finite. Now, we show that the process $\{(\boldsymbol{x}(t), \boldsymbol{G}(t), \boldsymbol{B}(t))\}$ under

policy π induces an ergodic Markov chain. This is pretty obvious, since $\{B(t)\} = \{\mathbf{0}_{J \times M}\}$ holds and processes $\{x_1(t)\}, \dots, \{x_I(t)\}, \{g_{1,1}(t)\}, \dots, \{g_{J,M}(t)\}$ are ergodic and independent.

To validate the second condition, we show that for all $\alpha \in (0, 1)$, $V_{\alpha}(s)$ is non-increasing with respect to x_i and $b_{j,m}$, where s = (x, G, B) is any state in \mathcal{S} and x_i and $b_{j,m}$ are the *i*th and (j,m)th entries of x and B, respectively. Consequently, by selecting s_0 as $(\mathbf{0}_{I\times 1}, \mathbf{G}_0, \mathbf{0}_{J\times M})$ with $\mathbf{G}_0 \triangleq \arg \max_{\mathbf{G} \in \mathcal{G}^{J\times M}} V_{\alpha}((\mathbf{0}_{I\times 1}, \mathbf{G}, \mathbf{0}_{J\times M})), V_{\alpha}(s) - V_{\alpha}(s_0) \leq 0$ holds for all $s \in \mathcal{S}$ and thus the second condition is verified.

Here, we only prove that $V_{\alpha}(s)$ is non-increasing with respect to x_i , and the proof for $b_{j,m}$ is very similar and thus omitted. First of all, based on the upper bound on reward in inequality (23), we compute the upper bound of $V_{\alpha}(s)$, denoted as $V_{\alpha,0}(s)$, as

$$V_{\alpha,0}(s) \triangleq \lim_{T \to \infty} \sum_{t=1}^{T} \alpha^{t-1} \sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} W_m \log\left(1 + \frac{g_{|G|}P}{N}\right)$$
$$= \frac{1}{1 - \alpha} \sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} W_m \log\left(1 + \frac{g_{|G|}P}{N}\right),$$

which is a constant and independent of s. Next, define $V_{\alpha,n}(s)$ for $n \in \mathbb{Z}^+$ by

$$V_{\alpha,n}(\boldsymbol{s}) \triangleq \max_{\boldsymbol{a}\in\mathcal{A}_{\boldsymbol{s}}} \Big\{ \alpha \sum_{\boldsymbol{s}'} \Pr\{\boldsymbol{s}'|\boldsymbol{s}, \boldsymbol{a}\} V_{\alpha,n-1}(\boldsymbol{s}') + r(t)|_{\boldsymbol{s}(t)=\boldsymbol{s},\boldsymbol{a}(t)=\boldsymbol{a}} \Big\}.$$

Then, we show that $V_{\alpha,n}(s)$ is non-increasing with respect to x_i for all $n \in \mathbb{Z}_{\geq 0}$ by induction:

- Rewrite s as (x_i, x_{-i}, G, B) , where x_{-i} consists of all entries of x except x_i , and construct s^* by $s^* \triangleq (x_i + 1, x_{-i}, G, B)$. Suppose $V_{\alpha,n}(s)$ is non-increasing with respect to x_i , i.e., $V_{\alpha,n}(s) \ge V_{\alpha,n}(s^*)$ holds, which is certainly true for n = 0.
- It can be easily verified that $\mathcal{A}_{s} = \mathcal{A}_{s^*}$ holds and

$$\alpha \sum_{\mathbf{s}'} \Pr\{\mathbf{s}' | \mathbf{s}, \mathbf{a}\} V_{\alpha, n}(\mathbf{s}') + r(t)|_{\mathbf{s}(t) = \mathbf{s}, \mathbf{a}(t) = \mathbf{a}}$$

$$\geq \alpha \sum_{\mathbf{s}'} \Pr\{\mathbf{s}' | \mathbf{s}^*, \mathbf{a}\} V_{\alpha, n}(\mathbf{s}') + r(t)|_{\mathbf{s}(t) = \mathbf{s}^*, \mathbf{a}(t) = \mathbf{a}}$$

holds for all $a \in \mathcal{A}_s$. Thus, $V_{\alpha,n+1}(s) \geq V_{\alpha,n+1}(s^*)$ holds.

Finally, since for all $\alpha \in (0,1)$, $V_{\alpha,n}(s)$ converges to $V_{\alpha}(s)$ as n goes to infinity [22, Proposition 7.3.1], we know that $V_{\alpha}(s) \geq V_{\alpha}(s^*)$, i.e., $V_{\alpha}(s)$ is non-increasing with respect to x_i .

APPENDIX B PROOF OF PROPOSITION 2.2

To begin with, we rephrase the optimal value of problem **(P1)** as

$$\max_{\{\boldsymbol{a}(t)\}} \liminf_{T \to \infty} \mathbb{E}_{\Pr\{\boldsymbol{x}' | \boldsymbol{x}, \boldsymbol{a}\}, \Pr\{\boldsymbol{G}' | \boldsymbol{G}\}} \left[\frac{1}{T} \sum_{t=1}^{T} r(t) \right]$$

$$\stackrel{(i)}{=} \max_{\pi \in \Pi^S} \lim_{T \to \infty} \mathbb{E}_{\pi, \Pr\{\boldsymbol{x}' | \boldsymbol{x}, \boldsymbol{a}\}, \Pr\{\boldsymbol{G}' | \boldsymbol{G}\}} \left[\frac{1}{T} \sum_{t=1}^{T} \left(-\sum_{i=1}^{I} w_i x_i(t) + \sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \left(\mathbbm{1}_{I+j} \left(a_m(t) \right) + \mathcal{G} \left(b_{j,m}(t) \right) \right) u_{j,m}(t) \right) \right]$$

$$\stackrel{(ii)}{=} \max_{\pi \in \Pi^S} \lim_{T \to \infty} \mathbb{E}_{\pi, \Pr\{\boldsymbol{x}' | \boldsymbol{x}, \boldsymbol{a}\}, \Pr\{\boldsymbol{G}' | \boldsymbol{G}\}} \left[\frac{1}{T} \sum_{t=1}^{T} \left(-\sum_{i=1}^{I} w_i x_i(t) + \sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbbm{1}_{I+j} \left(a_m(t) \right) \sum_{\tau=t}^{\tau+1} u_{j,m}(\tau) \right) \right],$$

$$(24)$$

where equality (i) holds due to Proposition 2.1 and equality (ii) can be easily derived by combining equalities in (2) and (7). In the following, we first show that for each stationary policy π for problem (P1), there exists another stationary policy $\hat{\pi}$ for problem (P2) such that their objective functions are equal, i.e., $(24)|_{\pi} = (14)|_{\hat{\pi}}$. Then, we show that the inverse holds, too. Based on these two results, problems (P1) and (P2) are obviously equivalent.

1) We first derive $(24)|_{\pi}$. Then, we develop another policy $\bar{\pi}$ for problem (P1), which is simpler than while equivalent

to policy π , i.e., $(24)|_{\pi}=(24)|_{\bar{\pi}}$. Finally, we introduce policy $\hat{\pi}$ for problem (**P2**) and show that $(24)|_{\bar{\pi}}=(14)|_{\hat{\pi}}$.

To derive $(24)|_{\pi}$, we list the countable states of problem (**P1**) as $s_1, s_2, \dots, s_l, \dots$, respectively, where $s_l = (x_l, G_l, B_l)$ is regarded as the l^{th} state in \mathcal{S} and the i^{th} entries of x_l is denoted as $x_{l,i}$. Then, denote the state occurrence distribution for problem (**P1**) under stationary policy π as $\mu^{\pi} \triangleq [\mu_1^{\pi}, \mu_2^{\pi}, \dots, \mu_l^{\pi}, \dots]^T$, where $\mu_l^{\pi} (\geq 0)$ is the state occurrence probability of the state s_l . And it follows

(24)

$$\begin{aligned} &(24)|_{\pi} \\ &\stackrel{(i)}{=} \sum_{l=1}^{\infty} \mu_{l}^{\pi} \mathbb{E}_{\boldsymbol{a} \sim \pi(\boldsymbol{s}_{l})} \left[\mathbb{E}_{\Pr\{\boldsymbol{x}'|\boldsymbol{x},\boldsymbol{a}\},\Pr\{\boldsymbol{G}'|\boldsymbol{G}\}} \left[\left(-\sum_{i=1}^{I} w_{i} \boldsymbol{x}_{l,i} + \sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbb{1}_{I+j} (a_{m}) \sum_{\tau=0}^{T_{j}-1} u_{j,m}(\tau) |_{g_{j,m}(0) = [\boldsymbol{G}_{l}]_{(j,m)}} \right) \right|_{\boldsymbol{s}_{l},\boldsymbol{a}_{l}} \right] \\ &= \sum_{l=1}^{\infty} \mu_{l}^{\pi} \left(-\sum_{i=1}^{I} w_{i} \boldsymbol{x}_{l,i} \Big|_{\boldsymbol{x}_{l}} + \mathbb{E}_{\boldsymbol{a} \sim \pi(\boldsymbol{s}_{l})} \mathbb{E}_{\Pr\{\boldsymbol{G}'|\boldsymbol{G}\}} \right] \\ &= \sum_{l=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbb{1}_{I+j} (a_{m}|_{\boldsymbol{a}}) \sum_{\tau=0}^{T_{j}-1} u_{j,m}(\tau) \Big|_{g_{j,m}(0) = [\boldsymbol{G}_{l}]_{(j,m)}} \right] \\ &\stackrel{(ii)}{=} \sum_{l=1}^{\infty} \mu_{l}^{\pi} \left(-\sum_{i=1}^{I} w_{i} \boldsymbol{x}_{l,i} \Big|_{\boldsymbol{x}_{l}} + \mathbb{E}_{\boldsymbol{a} \sim \pi(\boldsymbol{s}_{l})} \left[\sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbb{1}_{I+j} (a_{m}|_{\boldsymbol{a}}) \bar{u}_{j,m}(0) \Big|_{g_{j,m}(0) = [\boldsymbol{G}_{l}]_{(j,m)}} \right] \right), \end{aligned}$$

$$(25)$$

where $\boldsymbol{a} = [a_1, a_2, \dots, a_M]^T$ is the action in $\mathcal{A}_{\boldsymbol{s}_l}$; $\pi(\boldsymbol{s}_l)$ is the action distribution at state \boldsymbol{s}_l under policy π ; equality (*i*) holds since policy π is stationary; equality (*ii*) holds due to the definition in equality (12). An essential observation is that (25) does not involve \boldsymbol{B}_l . Thus, we realign it as

$$(25)$$

$$\stackrel{(i)}{=} \sum_{\hat{s} \in \hat{S}} \sum_{l \in \{l \mid s_{l} \in N(\hat{s})\}} \mu_{l}^{\pi} \left(-\sum_{i=1}^{I} w_{i} x_{i} \Big|_{x} + \mathbb{E}_{\boldsymbol{a} \sim \pi(\boldsymbol{s}_{l})} \right)$$

$$\left[\sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbb{1}_{I+j} \left(a_{m} \Big|_{\boldsymbol{a}} \right) \bar{u}_{j,m}(0) \Big|_{g_{j,m}(0) = [\boldsymbol{G}]_{(j,m)}} \right] \right)$$

$$(26)$$

$$\stackrel{(ii)}{=} \sum_{\hat{s} \in \hat{S}} \mu^{\pi}(N(\hat{s})) \left(-\sum_{i=1}^{I} w_{i} x_{i} \Big|_{x} \right) + \sum_{\hat{s} \in \hat{S}} \sum_{l \in \{l \mid s_{l} \in N(\hat{s})\}} \mu_{l}^{\pi} \mathbb{E}_{\boldsymbol{a} \sim \pi(\boldsymbol{s}_{l})}$$

$$\left[\sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbb{1}_{I+j} \left(a_{m}|_{a}\right) \bar{u}_{j,m}(0)\Big|_{g_{j,m}(0)=[G]_{(j,m)}}\right]$$
(27)

where in equality (i), \boldsymbol{x} and \boldsymbol{G} are the components of $\hat{\boldsymbol{s}}$, i.e., $\hat{\boldsymbol{s}} = (\boldsymbol{x}, \boldsymbol{G}, \boldsymbol{b})$ holds, and $N(\hat{\boldsymbol{s}})$ is the state set defined as $N(\hat{\boldsymbol{s}}) \triangleq \{\boldsymbol{s}_l = (\boldsymbol{x}_l, \boldsymbol{G}_l, \boldsymbol{B}_l) \in \mathcal{S} | \boldsymbol{x}_l = \boldsymbol{x}; \boldsymbol{G}_l =$ $\boldsymbol{G}; \sum_{j=1}^{J} [\boldsymbol{B}_l]_{(j,m)} = b_m, \forall m \in \{1, 2, \cdots, M\}\}$ with b_m being the m^{th} entry of \boldsymbol{b} ; in equality (ii), $\mu^{\pi}(N(\hat{\boldsymbol{s}}))$ is the summation of the state occurrence probabilities of all Now, we introduce a new policy $\bar{\pi}$ for problem (**P1**) and prove that (27) = (24) $|_{\bar{\pi}}$ holds. Specifically, denote $\bar{\pi}(s, a)$ and $\pi(s, a)$ as the probabilities of applying action a at state s = (x, G, B) under policies $\bar{\pi}$ and π , respectively. Then, construct policy $\bar{\pi}$ as

$$\bar{\pi}(\boldsymbol{s}, \boldsymbol{a}) = \sum_{l \in \{l | \boldsymbol{s}_l \in N(\hat{\boldsymbol{s}})\}} \mu_l^{\pi} \pi(\boldsymbol{s}_l, \boldsymbol{a}).$$
(28)

Here, we highlight again that $\hat{s} = (x, G, b)$ at the subscript of the RHS of (28) is induced from s by using $\boldsymbol{b} = [b_1, b_2, \dots, b_M]^T$ and $b_m = \sum_{j=1}^J [\boldsymbol{B}]_{(j,m)}$. Obviously, the action distribution that policy $\bar{\pi}$ follows at state s is actually the expected action distribution that policy π follows over the state set $N(\hat{s})$ and the input of policy $\bar{\pi}$ needs to know only $(\hat{s}, \boldsymbol{a})$ rather than (s, \boldsymbol{a}) . Therefore, policy $\bar{\pi}$ is simpler than policy π and accordingly we directly denote the action distribution at state s under policy $\bar{\pi}$ as $\bar{\pi}(\hat{s})$. Based on equality (28), (27) is equal to

$$\sum_{\hat{s}\in\hat{S}} \mu^{\pi}(N(\hat{s})) \left(-\sum_{i=1}^{I} w_{i}x_{i} \Big|_{x} \right) + \sum_{\hat{s}\in\hat{S}} \mu^{\pi}(N(\hat{s}))\mathbb{E}_{\boldsymbol{a}\sim\bar{\pi}}(\hat{s})$$

$$\left[\sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbb{1}_{I+j} \left(a_{m} \Big|_{\boldsymbol{a}} \right) \bar{u}_{j,m}(0) \Big|_{g_{j,m}(0)=[\boldsymbol{G}]_{(j,m)}} \right].$$
(29)

Now, we prove that $(29) = (24)|_{\bar{\pi}}$ holds, and an essential step is to show that $\mu^{\pi}(N(\hat{s})) = \mu^{\bar{\pi}}(N(\hat{s}))$ holds for all $\hat{s} \in \hat{S}$. Specifically, it can be easily checked that the transition probabilities of the state sets are equivalent in policies π and $\bar{\pi}$, i.e., $\Pr_{\pi}\{N(\hat{s}')|N(\hat{s})\} \triangleq \mathbb{E}_{\boldsymbol{a} \sim \pi(N(\hat{s}))}[\Pr\{N(\hat{s}')|N(\hat{s}),\boldsymbol{a}\}]$ equals $\Pr_{\bar{\pi}}\{N(\hat{s}')|N(\hat{s})\} \triangleq \mathbb{E}_{\boldsymbol{a} \sim \bar{\pi}(N(\hat{s}))}[\Pr\{N(\hat{s}')|N(\hat{s}),\boldsymbol{a}\}]$ for all $\hat{s}, \hat{s}' \in \hat{S}$. Consequently, the Markov chains for state sets under policies π and $\bar{\pi}$ are the same. And the occurrence probability of each state set is unique and can be derived by solving the Cerso limit of the Markov chain. Therefore, $\mu^{\pi}(N(\hat{s})) = \mu^{\bar{\pi}}(N(\hat{s}))$ holds for all $\hat{s} \in \hat{S}$. Based on this condition, it follows

$$(29) = \sum_{\hat{s} \in \hat{S}} \mu^{\bar{\pi}}(N(\hat{s})) \left(-\sum_{i=1}^{I} w_{i} x_{i} \Big|_{\boldsymbol{x}} \right) + \sum_{\hat{s} \in \hat{S}} \mu^{\bar{\pi}}(N(\hat{s})) \mathbb{E}_{\boldsymbol{a} \sim \bar{\pi}(\hat{s})} \\ \left[\sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbb{1}_{I+j} \left(a_{m} \Big|_{\boldsymbol{a}} \right) \bar{u}_{j,m}(0) \Big|_{g_{j,m}(0) = [\boldsymbol{G}]_{(j,m)}} \right]$$
(30)
=(24) $|_{\bar{\pi}}$.

Finally, we construct the policy $\hat{\pi}$ for problem (P2) by

$$\hat{\pi}(\hat{\boldsymbol{s}}, \boldsymbol{a}) = \sum_{l \in \{l | \boldsymbol{s}_l \in N(\hat{\boldsymbol{s}})\}} \mu_l^{\pi} \pi(\boldsymbol{s}_l, \boldsymbol{a}).$$
(31)

and show that $(24)|_{\bar{\pi}} = (14)|_{\hat{\pi}}$. Specifically, based on (28) and (31), the policy $\bar{\pi}$ for problem (**P1**) makes the same decision with the policy $\hat{\pi}$ for problem (**P2**) when their encountering states are \hat{s} and s, respectively. And based on (30), the instant

reward for problem (P1) under policy $\bar{\pi}$ can be equivalently regarded as

$$-\sum_{i=1}^{I} w_i x_i(t) + \sum_{j=1}^{J} w_{I+j} \sum_{m=1}^{M} \mathbb{1}_{I+j} (a_m(t)) \bar{u}_{j,m}(t),$$

which is exactly the instant reward $\hat{r}(t)$ for problem (P2). Consequently, the average rewards of problem (P1) under policy $\bar{\pi}$ is equals to that of problem (P2) under policy $\hat{\pi}$, i.e., $(24)|_{\bar{\pi}} = (14)|_{\hat{\pi}}$ holds.

2) As for the proof for the inverse, i.e., for each stationary policy $\hat{\pi}$ for problem (**P2**), there exists another stationary policy π for problem (**P1**) such that $(24)|_{\pi} = (14)|_{\hat{\pi}}$ holds, the utilized techniques are similar to the proof in *I*) and thus omitted. We only highlight that the policy π for problem (**P1**) is constructed by $\pi(s, a) = \hat{\pi}(\hat{s}, a)$.

APPENDIX C Sketch Proof of Proposition 2.3

Define $\hat{V}_{\alpha}(\hat{s})$ by $\hat{V}_{\alpha}(\hat{s}) \triangleq \sup_{\theta} \lim_{T \to \infty} \mathbb{E}_{\pi_{\theta}, \Pr\{\hat{s}' | \hat{s}, a\}} \left[\sum_{t=1}^{T} \alpha^{t-1} \hat{r}(t) |_{\hat{s}(1)=\hat{s}} \right]$, where \hat{s} can be any state in \hat{S} ; $\alpha \in (0, 1)$ is a discount factor; θ parameterizes the policy π_{θ} . Similar to the proof in Appendix A, we can show that $\hat{V}_{\alpha}(\hat{s})$ is non-increasing with respect to x_i and consequently, the optimal policies for problem (**P3**) are of threshold type with respect to x_i : if the optimal policy is to transmit data for the i^{th} monitoring device at state $(x_i, \mathbf{x}_{-i}, \mathbf{G}, \mathbf{b})$, it also transmits data for the i^{th} monitoring device at state $(x_i + 1, \mathbf{x}_{-i}, \mathbf{G}, \mathbf{b})$. Since the optimal threshold cannot be infinitely large, there are finite number of threshold-type policies possibly being the optimal policies. Accordingly, based on [23, Proposition 4.1.3], there exist Blackwell optimal policies for problem (**P3**), and based on [23, Proposition 4.1.7], these policies optimize problem (**P2**). Thus, Proposition 2.3 holds.

APPENDIX D PROOF OF PROPOSITION 3.1

We first introduce the explicit formulation of the decoupled sub-problems. Next, we validate that the optimal policies for these problems are of threshold type. Then, we derive the optimal policies for these sub-problems. Finally, we validate the existence of the Whittle's index for problem (**P5**).

1) Decoupled sub-problem: To decouple the problem (P5), we let all monitoring devices be selfish such that each of them aims to minimize its own average weighted AoII. Moreover, each monitoring device is allowed to transmit data at any time slot as long as it pays an additional cost C for each transmission. The goal of each sub-problem is to find the optimal scheduling policy which strikes the balance between the average additional costs and the average weighted AoII for each monitoring device. We formulate the i^{th} sub-problem as the following Markov decision problem.

- state: $x_i(t) \in \mathbb{Z}_{\geq 0}$;
- action: a(t) ∈ {0,1}, where a(t) = 1 means to transmit data for the ith monitoring device at the tth time slot and a(t) = 1 means not;

- transitions: $\Pr\{x_i(t+1) = 0|a(t) = 1\} = p_i$; $\Pr\{x_i(t+1) = 1\} = 1 p_i$; $\Pr\{x_i(t+1) = 0|a(t) = 0, x_i(t) = 0\} = p_i$; $\Pr\{x_i(t+1) = 1|a(t) = 0, x_i(t) = 0\} = 1 p_i$; $\Pr\{x_i(t+1) = 0|a(t) = 0, x_i(t) > 0\} = q_i$; $\Pr\{x_i(t+1) = x_i(t) + 1|a(t) = 0, x_i(t) > 0\} = 1 q_i$;
- cost: $c(t) = x_i(t) + a(t)C;$
- optimality criteria: *lim* average optimality criteria.

2) The structure of the optimal policy: We first study the decoupled sub-problem with total discounted cost criteria and analyze the corresponding optimal policies. Specifically, define $V_{\alpha}(x)$ by $V_{\alpha}(x) \triangleq \inf_{\theta} \lim_{T \to \infty} \mathbb{E}_{\pi_{\theta}, \Pr\{x'_i | x_i, a\}} \left[\sum_{t=1}^{T} \alpha^{t-1} \right]$ $c(t)|_{x_i(1)=x}$, where $\alpha \in (0,1)$ is a discount factor; θ parameterizes the policy π_{θ} . Similar to the proof in Appendix A, we can show that $V_{\alpha}(x)$ is non-decreasing with respect to x, i.e., $V_{\alpha}(x) \leq V_{\alpha}(x+1)$ holds for all $x \in \mathbb{Z}_{\geq 0}$. Now, we show that the optimal policy for the discounted version of the decoupled sub-problem is of threshold type: 1) consider the optimal policy is to transmit data for the i^{th} monitoring device at the state $x \in \mathbb{Z}^+$; 2) then, based on Bellman's optimality equation (Prop 7.3.1 in [22]), $x + C + \alpha(p_i V_{\alpha}(0) + (1 - 1))$ $p_i V_{\alpha}(x+1) \leq x + \alpha (q_i V_{\alpha}(0) + (1-q_i) V_{\alpha}(x+1))$ holds; 3) since $V_{\alpha}(x+1) \leq V_{\alpha}(x+2)$ holds, $x+1+C+\alpha(p_iV_{\alpha}(0)+$ $(1 - p_i)V_{\alpha}(x + 2)) \le x + 1 + \alpha(q_i V_{\alpha}(0) + (1 - q_i)V_{\alpha}(x + 2))$ holds, too; 4) thus, the optimal policy will also transmit data for the i^{th} monitoring device at the state x+1, which completes the proof.

Similar to the proof in Appendix C, there are finite number of threshold-type policies possibly being the optimal policies for the discounted decoupled sub-problems. Then, based on Prop 4.1.3 and Prop 4.1.7 in [23], there exist threshold-type optimal policies for the original decoupled sub-problem with *lim* average optimality criteria.

3) The derivation of the optimal policy: To derive the optimal policy, we randomly investigate a threshold-type policy π_{x_0} : if $x \ge x_0$, the *i*th monitoring device transmits data; if $x < x_0$, the *i*th monitoring device does not transmit data. By solving the transition equations, we derive that

$$\mu_{i,x} = \begin{cases} \frac{1}{1 + \frac{1-p_i}{q_i} - (1-p_i)\left(\frac{1}{q_i} - \frac{1}{p_i}\right)(1-q_i)^{x_0-1}} & x=0\\ (1-p_i)(1-q_i)^{x-1}\mu_{i,0} & x=1,2,\cdots,x_0\\ (1-p_i)^{x-x_0+1}(1-q_i)^{x_0-1}\mu_{i,0} & x=x_0+1,x_0+2,\cdots,\end{cases}$$
(32)

where $\mu_{i,x}$ is the state occurrence probability of the state x in the *i*th sub-problem under policy π_{x_0} . And the average cost equals

$$f_i(x_0, C) \triangleq \sum_{x=0}^{x_0-1} w_i x \mu_{i,x} + \sum_{x=x_0}^{\infty} (w_i x + C) \mu_{i,x}$$
$$= \frac{\beta_1 + (\beta_2 + \beta_3 x_0)(1 - q_i)^{x_0}}{\beta_4 - \beta_5 (1 - q_i)^{x_0-1}},$$
(33)

where

$$\begin{split} \beta_1 &= w_i \frac{1-p_i}{q_i^2} > 0; \\ \beta_2 &= w_i \frac{(1-p_i)^2}{p_i^2(1-q_i)} - w_i \frac{1-p_i}{q_i^2} + \frac{1-p_i}{p_i(1-q_i)}C; \\ \beta_3 &= w_i \frac{1-p_i}{1-q_i} \left(\frac{1}{p_i} - \frac{1}{q_i}\right) < 0; \ \beta_4 = 1 + \frac{1-p_i}{q_i} > 0; \\ \beta_5 &= (1-p_i) \left(\frac{1}{q_i} - \frac{1}{p_i}\right) > 0. \end{split}$$

Thus, we can derive the optimal policy by finding the optimal threshold $x_i(C) \in \mathbb{Z}^+$, which is defined by $x_i(C) = \arg \min_{x_0 \in \mathbb{Z}^+} f_i(x_0, C)$. Notably, it is not easy to derive the exact value of $x_i(C)$ and neither will we derive it. Instead, the analyses here are used to prove the existence of the Whittle's index for problem (**P5**).

4) The indexability for the decoupled sub-problem: Based on [8], the existence of the Whittle's index is guaranteed if all the sub-problems are indexable. Specifically, we give the explicit definition of the indexability as follows.

Definition D.1(indexability): Define $Z_i(C) = \{x \in \mathbb{Z}_{\geq 0} | x < x_i(C)\}$ as the set of states where the optimal policy is not to transmit data for the *i*th monitoring device. Then, the *i*th decoupled sub-problem is said to be indexable if it follows

$$C' \ge C \Rightarrow Z_i(C') \supseteq Z_i(C)$$

Apparently, it is difficult to directly validate the indexability for the decoupled sub-problem based on the above definition. Instead, we refer [29, Proposition 2.2], according to which, the i^{th} decoupled sub-problem is indexable as long as $\sum_{x=x_0}^{\infty} \mu_{i,x}$ is decreasing with respect to x_0 . Based on (32), it follows

$$\sum_{x=x_0}^{\infty} \mu_{i,x} = \frac{1-p_i}{q_i \left(\frac{1+\frac{1-p_i}{q_i}}{(1-q_i)^{x_0-1}} - (1-p_i)(\frac{1}{q_i} - \frac{1}{p_i})\right)}$$

and apparently, $\sum_{x=x_0}^{\infty} \mu_{i,x}$ is decreasing with respect to x_0 . This completes the proof.

REFERENCES

- L. Chettri and R. Bera, "A comprehensive survey on internet of things (IoT) toward 5G wireless systems," *IEEE Internet Things J.*, vol. 7, no. 1, pp. 16-32, Jan. 2020.
- [2] M. A. Al Mamun and M. R. Yuce, "Sensors and systems for wearable environmental monitoring toward IoT-enabled applications: A review," *IEEE Sens. J.*, vol. 19, no. 18, pp. 7771-7788, Sept. 2019.
- [3] W. Feng, J. Tang, Y. Yu, J. Song, N. Zhao, G. Chen, K.-K. Wong, and J. Chambers, "UAV-enabled SWIPT in IoT networks for emergency communications," *IEEE Wireless Commun.*, vol. 27, no. 5, pp. 140-147, Oct. 2020.
- [4] H. Elayan, M. Aloqaily, and M. Guizani, "Digital twin for intelligent context-aware IoT healthcare systems," *IEEE Internet Things J.*, vol. 8, no. 23, pp. 16749-16757, Dec. 2021.
- [5] S. Kaul, R. Yates, and M. Gruteser, "Real-time status: how often should one update?" in *Proc. IEEE INFOCOM*, Orlando, FL, USA, Mar. 2012.
- [6] J. Sun, Z. Jiang, B. Krishnamachari, S. Zhou, and Z. Niu, "Closed-form Whittle's index-enabled random access for timely status update," *IEEE Trans. Commun.*, vol. 68, no. 3, pp. 1538-1551, Mar. 2019.
- [7] A. Maatouk, S. Kriouile, M. Assad, and A. Ephremides, "On the optimality of the Whittle's index policy for minimizing the age of information," *IEEE Trans. Wireless Commun.*, vol. 20, no. 2, pp. 1263-1277, Feb. 2021.
- [8] P. Whittle, "Restless bandits: activity allocation in a changing world," J. Appl. Probab., vol. 25, no. A, pp. 287-298, 1988.

- [9] 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.
- [10] M. J. Neely, "Stochastic network optimization with application to communication and queueing systems," *Synth. Lect. Commun. Netw.*, vol. 3, no. 1, pp. 1-211, Sept. 2010.
- [11] I. Kadota, A. Sinha, and E. Modiano, "Scheduling algorithms for optimizing age of information in wireless networks with throughput constraints," *IEEE/ACM Trans. Netw.*, vol. 27, no. 4, pp. 1359-1372, Aug. 2019.
- [12] A. Kosta, N. Pappas, A. Ephremides, and V. Angelakis, "Age of information and throughput in a shared access network with heterogeneous traffic," in *IEEE Proc. GLOBECOM*, Abu Dhabi, United Arab Emirates, Feb. 2018, pp. 1-6.
- [13] J. Sun, L. Wang, Z. Jiang, S. Zhou, and Z. Niu, "Age-optimal scheduling for heterogeneous traffic with timely throughput constraints," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 5, pp. 1485-1498, May 2021.
- [14] A. Maatouk, S. Kriouile, M. Assaad, and A. Ephremides, "The age of incorrect information: a new performance metric for status updates," *IEEE/ACM Trans, Netw.*, vol. 28, no. 5, pp. 2215-2228, Oct. 2020.
- [15] A. Maatouk, M. Assaad, and A. Ephremides, "The age of incorrect information: and enabler of semantics-empowered communication," arXiv preprint arXiv:2012.13214, Dec. 2020.
- [16] C. Kam, S. Kompella, and A. Ephremides, "Age of incorrect information for remote estimation of a binary Markov source," in *Proc. IEEE INFOCOM WKSHPS*, Toronto, ON, Canada, Jul. 2020, pp. 1-6.
- [17] Y. Chen and A. Ephremides, "Scheduling to minimize age of incorrect information with imperfect channels state information," *Entropy*, vol. 23, no. 12, pp. 1572, Nov. 2021.
- [18] Y. Chen and A. Ephremides, "Minimizing age of incorrect information for unreliable channel with power constraint," arXiv preprint arXiv:2101.08908, Jan. 2021.
- [19] S. Kriouile and M. Assaad, "When to pull data from sensors for minimum distance-based age of incorrect information metric," *arXiv* preprint arXiv:2202.02878, Feb. 2022.
- [20] M. L. Puterman, Markov decision processes: discrete stochastic dynamic programming, John Wiley & Sons, Hoboken, New Jersey, 2014.
- [21] L. I. Sennott, "Average cost optimal stationary policies in infinite state Markov decision processes with unbounded costs," *Math. Oper, Res.*, vol. 37, no. 4, pp. 626-633, Aug. 1989.
- [22] D. P. Bertsekas, Dynamic programming and optimal control: volume I, Athena scientific, Belmont, MA, 2012.
- [23] D. P. Bertsekas, Dynamic programming and optimal control: volume II, Athena scientific, Belmont, MA, 2011.
- [24] R. S. Sutton and A. G. Barto, *Reinforcement learning: an introduction*, MIT press, Cambridge, MA, 2018.
- [25] G. Dulac-Arnold, R. Evans, H. Van Hasselt, P. Sunehag, T. Lillicrap, J. Hunt, T. Mann, T. Weber, T. Degris, and B. Coppin, "Deep reinforcement learning in large discrete action spaces," *arXiv preprint arXiv:1512.07679*, Dec. 2015.
- [26] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," in *ICLR*, San Juan, Puerto Rico, USA, May 2016.
- [27] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, Jul. 2017.
- [28] H. S. Wang and N. Moayeri, "Finite-state Markov channel-a useful model for radio communication channels," *IEEE Trans. Veh. Tech.*, vol. 44, no. 1, pp. 163-171, Feb. 1995.
- [29] M. Larrañaga, "Dynamic control of stochastic and fluid resource-sharing systems," Ph.D. dissertation, Institut National Polytechnique de Toulouse, Toulouse, French Republic, 2015.