

# Joint Status Sampling and Updating for Minimizing Age of Information in the Internet of Things

Bo Zhou, *Member, IEEE* and Walid Saad, *Fellow, IEEE*

**Abstract**—The effective operation of time-critical Internet of things (IoT) applications requires real-time reporting of fresh status information of underlying physical processes. In this paper, a real-time IoT monitoring system is considered, in which the IoT devices sample a physical process with a sampling cost and send the status packet to a given destination with an updating cost. This joint status sampling and updating process is designed to minimize the average age of information (AoI) at the destination node under an average energy cost constraint at each device. This stochastic problem is formulated as an infinite horizon average cost constrained Markov decision process (CMDP) and transformed into an unconstrained Markov decision process (MDP) using a Lagrangian method. For the single IoT device case, the optimal policy for the CMDP is shown to be a randomized mixture of two deterministic policies for the unconstrained MDP, which is of threshold type. This reveals a fundamental tradeoff between the average AoI at the destination and the sampling and updating costs. Then, a structure-aware optimal algorithm to obtain the optimal policy of the CMDP is proposed and the impact of the wireless channel dynamics is studied while demonstrating that channels having a larger mean channel gain and less scattering can achieve better AoI performance. For the case of multiple IoT devices, a low-complexity semi-distributed suboptimal policy is proposed with the updating control at the destination and the sampling control at each IoT device. Then, an online learning algorithm is developed to obtain this policy, which can be implemented at each IoT device and requires only the local knowledge and small signaling from the destination. The proposed learning algorithm is shown to converge almost surely to the suboptimal policy. Simulation results show the structural properties of the optimal policy for the single IoT device case; and show that the proposed policy for multiple IoT devices outperforms a zero-wait baseline policy, with average AoI reductions reaching up to 33%.

**Index Terms**—Internet of things, status update, age of information, Markov decision processes, structural analysis, distributed stochastic learning.

## I. INTRODUCTION

With the rapid proliferation of the Internet of Thing (IoT) devices, delivering timely status information of the underlying physical processes has become increasingly critical for many real-world IoT and cyber-physical system applications [2], [3], such as environment monitoring in sensor networks and vehicle tracking in smart transportation systems. Given the criticality of IoT applications, it is imperative to maintain the status information of the physical process at the destination nodes as fresh as possible, for effective monitoring and control.

This research was supported by the U.S. National Science Foundation under Grants IIS-1633363, CNS-1836802, and CNS-1460316. A preliminary version of this work has been presented at IEEE GLOBECOM 2018 [1].

B. Zhou and W. Saad are with Wireless@VT, Bradley Department of Electrical and Computer Engineering, Virginia Tech, Blacksburg, VA 24061, USA. Email: {ecebo, walids}@vt.edu.

To quantify the freshness of the status information of the physical process, the concept of *age of information* (AoI) has been proposed as a key performance metric [4] that quantifies the time elapsed since the generation of the most recent IoT device status packet received at a given destination. In contrast to conventional delay metrics, which measure the time interval between the generation and the delivery of each individual packet, the AoI considers the packet delay and the generation time of each packet, and, hence, characterizes the freshness of the status information from the perspective of the destination. Therefore, optimizing the AoI in an IoT would lead to distinctively different system designs from those used for conventional delay optimization. For example, it has been shown that the last-come-first-served (LCFS) principle achieves a lower AoI than the conventional first-come-first-served (FCFS) principle [5].

The problem of minimizing the AoI has attracted significant recent attention [4]–[17]. Generally, these works can be classified into two broad groups based on the model of the generation process of the status packets. The first group [4]–[10] models the generation process of the status packets as a queueing system in which the status packets arrive at the source node stochastically and are queued before being forwarded to the destination. Queueing theory has also been used to analyze and optimize the average AoI for FCFS [4] and LCFS systems [5]. The works in [6]–[8] propose scheduling schemes that seek to minimize the average AoI in wireless broadcast networks. In [9] and [10], the authors study the problem of AoI minimization in wireless multiaccess networks and propose decentralized scheduling policies with near-optimal performance. In the second group of works [11]–[16], the status packets can be generated at any time by the source node. The authors in [11] and [12] propose optimal updating policies to minimize the average AoI for status update systems, with a single source and multiple sources, respectively. In [13]–[15], the authors propose optimal status updating schemes for an energy harvesting source to minimize the average AoI. The authors in [16] introduce an optimal status updating scheme to minimize the average AoI under resource constraints. Motivated by recent research on AoI, the authors in [17] study the remote estimation problem for a Wiener process and propose an optimal sampling policy to minimize the estimation error.

In the existing literature, e.g., [4]–[17], the source node is usually required to perform simple monitoring tasks, such as reading a temperature sensor, and, hence, the cost for generating status packets is assumed to be negligible. However, next-generation IoT devices can now perform more complex

tasks<sup>1</sup>, such as initial feature extraction and classification for computer vision applications, by using neural networks and on-device artificial intelligence [19], [20]. For such applications, generating the status update packets incurs energy cost for the IoT devices. Moreover, compared to the status packets generated for simple monitoring tasks (e.g., a temperature reading), a generated status packet for sophisticated artificial intelligence tasks carries richer information on the underlying physical systems (e.g., objects detected in an image or video sequence). Therefore, there will also incur some energy cost and time delay for transmitting those status packets with relatively large size to the destination node. In presence of the energy cost pertaining to the sampling and updating processes, a key open problem is to study how to intelligently sample the underlying physical systems and send status packets to the destination, in order to minimize the AoI.

The main contribution of this paper is, thus, to jointly design the status sampling and updating processes that can minimize the average AoI at the destination under an average energy cost constraint for each IoT device, by taking into account the energy cost for generating and updating status packets. In particular, our key contributions include:

- ≤ For the single IoT device case, we formulate this stochastic control problem as an infinite horizon average cost constrained Markov decision process (CMDP) and transform the CMDP into a parameterized unconstrained Markov decision process (MDP) using a Lagrangian method. We show that the optimal policy for the CMDP is a randomized mixture of two deterministic policies for the unconstrained MDP. By using the special properties of the AoI dynamics, we derive key properties of the value function for the unconstrained MDP. Based on these properties, we show that the optimal sampling and updating process for the unconstrained MDP is threshold-based with the AoI state at the device and the AoI state at the destination. This reveals a fundamental tradeoff between the average AoI at the destination and the sampling and updating costs. Then, we propose a structure-aware optimal algorithm to obtain the optimal policy for the CMDP. We also study the influence of the wireless channel fading distribution on the optimal average AoI at the destination. By using the concept of stochastic dominance, we show that channels having a larger mean channel gain and less scattering can achieve better AoI performance.
- ≤ For the case of multiple IoT devices, to obtain the optimal sampling and updating policy, we also formulate a CMDP and convert it to an unconstrained MDP. We show that the optimal sampling and updating policy, which adapts to the AoI and channels states of all IoT devices, is a function of the Q-factors of the unconstrained MDP. To overcome the curse of dimensionality and to distribute the system's controls, we propose a low-complexity semi-distributed suboptimal policy by approximating the opti-

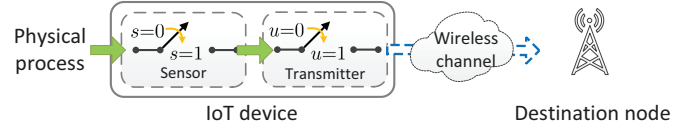


Fig. 1: Illustration of a real-time monitoring system with a single IoT device.

mal Q-factors into the sum of per-device Q-factors, based on approximate dynamic programming. Then, we propose an online learning algorithm that allows each device to learn its per-device Q-factor, which requires only the knowledge of the local AoI and channel states, as well as small signaling from the destination. The proposed semi-distributed online learning algorithm is shown to converge almost surely to the proposed suboptimal policy.

- ≤ We provide extensive simulations to illustrate additional structural properties of the optimal policy for the single device case. We show that the optimal thresholds for sampling and updating are non-decreasing with the sampling cost and the updating cost, respectively, and the optimal action is also threshold-based with respect to the channel state. For the case of multiple devices, numerical results show that the proposed semi-distributed policy outperforms a zero-wait baseline policy (i.e., sampling immediately after updating), with average AoI reductions reaching up to 33%. In summary, the derived results provide novel and holistic insights on the design of AoI-aware sampling and updating in practical IoT systems.

The rest of this paper is organized as follows. In Section II, we present the single IoT device model and analyze its properties. In Section III, we present the analysis for the case of multiple IoT devices using online learning. Section IV presents and analyzes numerical results. Finally, conclusions are drawn in Section V.

## II. OPTIMAL SAMPLING AND UPDATING CONTROL FOR A SINGLE IoT DEVICE

### A. System Model

Consider a real-time IoT monitoring system composed of a single IoT device and a destination node (e.g., a base station or control center), as illustrated in Fig. 1. The IoT device encompasses a sensor which can monitor the real-time status of a physical process (referred to hereinafter as *status sampling*) and a transmitter which can send status information packets to the destination through a wireless channel (referred to hereinafter as *status updating*). For the status sampling process, different from the existing literature where the device is usually assumed to perform simple sampling tasks [4]–[17], e.g., temperature and humidity monitoring, here, we consider that the IoT device can perform more sophisticated tasks, e.g., initial feature extraction and pre-classification using machine learning and neural network tools. Hence, the time for status sampling and updating is not negligible and there will be some associated energy expenditures, which constrain the operation of the IoT device.

<sup>1</sup>One practical example is the Nest Cam IQ indoor security camera, which uses on-device vision processing to watch for motion, distinguish family members, and send alerts if someone is not recognized [18].

We consider a time-slotted system with unit slot length (without loss of generality) that is indexed by  $t = 1, 2, \dots$ . Let  $h[t] \in \mathcal{H}$  be the channel state, representing the channel gain at slot  $t$ , where  $\mathcal{H}$  is the finite channel state space. We assume a block fading wireless channel over all time slots and we consider an i.i.d. channel state process  $\{h[t]\}$  that is distributed according to a general distribution  $p_{\mathcal{H}}(h)$ . Note that, the analytical framework and results can be readily extended to the Markovian fading channels.

1) *Monitoring Model*: In each slot, the IoT device must decide whether to generate a status packet and whether to send to the status packet to the destination. Let  $s[t] \in \{0, 1\}$  be the sampling action of the device at slot  $t$ , where  $s[t] = 1$  indicates that the device samples the physical process and generates a status packet at slot  $t$ , and  $s[t] = 0$ , otherwise. We consider that a newly generated status packet will replace the older one at the device, as the destination will not benefit from receiving an outdated status update. This is similar to the LCFS principle explored in [5]. Let  $C_s$  be the sampling cost for generating the status packet. This cost captures the computational cost needed for running some pre-classification algorithms using neural network models. We assume that the status sampling process takes one time slot. Let  $u[t] \in \{0, 1\}$  be the update action of the device at slot  $t$ , where  $u[t] = 1$  indicates that the device sends the status packet to the destination at slot  $t$  and  $u[t] = 0$ , otherwise. The IoT device can only send the status packet available locally. We denote by  $C_u(h)$  the minimum transmission power required by the IoT device for successfully updating a status packet to the destination within a slot when the channel state is  $h$ . Without loss of generality, we assume that  $C_u(h)$  is decreasing with  $h$ .

Let  $\mathbf{w}[t] \triangleq (s[t], u[t]) \in \mathcal{W} \triangleq \{0, 1\} \times \{0, 1\}$  be the control action vector of the IoT device at  $t$ . Then, the energy cost at the device associated with action  $\mathbf{w}[t]$  is given by:  $C(\mathbf{w})[t] \triangleq s[t]C_s + u[t]C_u(h)[t]$ .

2) *Age of Information Model*: We adopt the AoI as the key performance metric to quantify the freshness of the status information packet [4]. The AoI is essentially defined as the time elapsed since the generation of the last status update of the physical process. Let  $A_r[t]$  be the AoI at the destination at the beginning of slot  $t$ . Then, we have  $A_r[t] = t - \delta[t]$ , where  $\delta[t]$  is the time slot during which the most up-to-date status packet received by the destination was generated. Note that, the IoT device can only send its currently available status packet to the destination. Thus, the AoI at the destination depends on the AoI at the device, i.e., the age of the status packet at the device. Let  $A_l[t]$  be the AoI at the device at the beginning of slot  $t$ . The AoI at the device and the AoI at the destination are maintained by the device and can be implemented using counters. Let  $\hat{A}_l$  and  $\hat{A}_r$  be, respectively, the upper limits of the corresponding counters for the AoI at the device and the AoI at the destination. We assume that  $\hat{A}_l$  and  $\hat{A}_r$  are finite. This is due to that, for time-critical IoT applications, it is not meaningful for the destination node to receive a status information with an infinite age. Such highly outdated status information will not be of any use to the system or underlying application. Note that, the obtained results hold for arbitrarily finite  $\hat{A}_l$  and  $\hat{A}_r$ , no matter how

small or large  $\hat{A}_l$  and  $\hat{A}_r$  are. We denote by  $\mathcal{A}_l \triangleq \{1, 2, \dots, \hat{A}_l\}$  and  $\mathcal{A}_r \triangleq \{1, 2, \dots, \hat{A}_r\}$  the state space for the AoI at the device and the AoI at the destination, respectively. We also define  $\mathbf{A}[t] \triangleq (A_l[t], A_r[t]) \in \mathcal{A}$  as the system AoI state at the beginning of slot  $t$ , where  $\mathcal{A} \triangleq \mathcal{A}_l \times \mathcal{A}_r$  is the system AoI state space.

For the AoI at the device, if the device samples the physical process at slot  $t$  (i.e.,  $s[t] = 1$ ), then the AoI decreases to one (due to one slot used for status sampling), otherwise, the AoI increases by one. Then, the dynamics of the AoI at the device will be given by:

$$A_l[t+1] = \begin{cases} 1, & \text{if } s[t] = 1, \\ \min\{A_l[t] + 1, \hat{A}_l\}, & \text{otherwise.} \end{cases} \quad (1)$$

For the AoI at the destination, if the device sends the status packet to the destination at slot  $t$  (i.e.,  $u[t] = 1$ ), then the AoI decreases to the AoI at the device at slot  $t$  plus one (due to one slot used for status packet transmission), otherwise, the AoI increases by one. Then, the dynamics of the AoI at the destination will be given by:

$$A_r[t+1] = \begin{cases} \min\{A_l[t] + 1, \hat{A}_r\}, & \text{if } u[t] = 1, \\ \min\{A_r[t] + 1, \hat{A}_r\}, & \text{otherwise.} \end{cases} \quad (2)$$

Note that, the analytical framework can be extended to the scenario in which more than one slot are needed to generate or send a status packet, by modifying the AoI dynamics in (1) and (2), accordingly. Clearly, it may not be optimal for the device to sample the physical process immediately after updating the status. The reason is that the newly generated status packet, if not transmitted to the destination immediately (due to a possibly poor channel state), can become stale and less useful for the destination, yielding energy waste for sampling. Therefore, we are motivated to investigate how to jointly control the sampling and updating processes so as to minimize the AoI at the destination, under the stringent energy constraint at the IoT device.

## B. CMDP Formulation and Optimality Equation

1) *CMDP Formulation*: Given an observed system AoI state  $\mathbf{A}$  and channel state  $h$ , the IoT device determines the sampling and updating action  $\mathbf{w}$  according to the following policy.<sup>2</sup>

**Definition 1:** A stationary sampling and updating policy  $\pi$  is defined as a mapping from the system AoI and the channel states  $\mathcal{A} \times \mathcal{H}$  to the control action of the device  $\mathcal{W}$ , where  $\pi(\mathbf{A}, h) = \mathbf{w}$ .

Under the i.i.d. assumption for the channel state process and the AoI dynamics in (1) and (2), the induced random process  $\{\mathbf{A}[t], h[t]\}$  is a controlled Markov chain. Hereinafter, as is commonly done (e.g., see [21] and [16]), we restrict our attention to stationary unichain policies to guarantee the existence of the stationary optimal policies. For a given stationary

<sup>2</sup>Here, we consider the entire AoI state space of  $\mathbf{A}$ . However, in practice, one may only consider the AoI states  $\mathbf{A}$  such that  $A_r \geq A_l$  without sacrificing optimality.

unichain policy  $\pi$ , the average AoI at the destination and the average energy cost will be:

$$\bar{A}_r(\pi) \triangleq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[A_r]t, \quad (3)$$

$$\bar{C}(\pi) \triangleq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[C(\mathbf{w})]t, \quad (4)$$

where the expectation is taken with respect to the measure induced by the policy  $\pi$ .

We seek to find the optimal sampling and updating policy that minimizes the average AoI at the destination under an average energy cost constraint at the device, as follows:

$$\bar{A}_r^* \triangleq \min_{\pi} \bar{A}_r(\pi), \quad (5a)$$

$$\text{s.t. } \bar{C}(\pi) \leq C^{\max}. \quad (5b)$$

Here  $\pi$  is a stationary unichain policy and  $\bar{A}_r^*$  denotes the minimum average AoI at the destination achieved by the optimal policy  $\pi^*$  under the constraint in (5b). The problem in (5) is an infinite horizon average cost CMDP, which is known to be challenging due to the curse of dimensionality.

2) *Optimality Equation*: To obtain the optimal policy  $\pi^*$  for the CMDP in (5), we reformulate the CMDP into a parameterized unconstrained MDP using the Lagrangian approach [22]. For a given Lagrange multiplier  $\lambda$ , we define the Lagrange cost at slot  $t$  as

$$L(\mathbf{A})t, h)t, \mathbf{w})t; \lambda \triangleq A_r)t + \lambda C(\mathbf{w})t. \quad (6)$$

Then, the average Lagrange cost under policy  $\pi$  is given by:

$$\bar{L}(\pi); \lambda \triangleq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[L(\mathbf{A})t, h)t, \mathbf{w})t; \lambda]. \quad (7)$$

Now, we have an unconstrained MDP whose goal is to minimize the average Lagrange cost:

$$\bar{L}^*(\lambda) \triangleq \min_{\pi} \bar{L}(\pi); \lambda, \quad (8)$$

where  $\bar{L}^*(\lambda)$  is the minimum average Lagrange cost achieved by the optimal policy  $\pi_{\lambda}^*$  for a given  $\lambda$ . According to [21, Theorem 1] and [23, Theorem 4.4], we have the following relation between the optimal solutions of the problems in (5) and (8).

**Lemma 1:** The optimal average AoI cost in (5a) and the optimal average Lagrange cost in (8) satisfy:

$$\bar{A}_r^* = \max_{\lambda \geq 0} \bar{L}^*(\lambda) - \lambda C^{\max}. \quad (9)$$

The optimal policy  $\pi^*$  of the CMDP in (5) is a randomized mixture of two deterministic stationary policies  $\pi_{\lambda_1}^*$  and  $\pi_{\lambda_2}^*$ , in the form of

$$\pi^* = \alpha \pi_{\lambda_1}^* + (1 - \alpha) \pi_{\lambda_2}^*, \quad (10)$$

where  $\alpha \in [0, 1]$  is the randomization parameter, and  $\pi_{\lambda_1}^*$  and  $\pi_{\lambda_2}^*$  are the optimal policies of the unconstrained MDP in (8) under the Lagrange multipliers  $\lambda_1$  and  $\lambda_2$ , respectively.

To obtain the optimal policy  $\pi^*$  of the CMDP, according to [24, Propositions 4.2.1, 4.2.3, and 4.2.5] (these propositions are restated in Appendix H), we first obtain the optimal policy

$\pi_{\lambda}^*$  for a given  $\lambda$  of the unconstrained MDP by solving the following Bellman equation.

**Lemma 2:** For any  $\lambda$ , there exists  $\theta_{\lambda}, \{V(\mathbf{A}, h); \lambda\}$  satisfying:

$$\begin{aligned} \theta_{\lambda} + V(\mathbf{A}, h); \lambda &= \min_{\mathbf{w} \in \mathcal{W}} \{L(\mathbf{A}, h, \mathbf{w}); \lambda \\ &+ p_{\mathcal{H}}(\mathbf{h})h[V(\mathbf{A}, h); \lambda] + \sum_{\mathbf{A} \in \mathcal{A}} p_{\mathcal{A}}(\mathbf{A})V(\mathbf{A}, h); \lambda\}, \end{aligned} \quad (11)$$

where  $\mathbf{A}^{\infty}$  satisfies the AoI dynamics in (1) and (2),  $\theta_{\lambda} = \bar{L}^*(\lambda)$  is the optimal value to (8) for all initial state  $(\mathbf{A})1, h)1$ , and  $V(\mathbf{A}, h)$  is the value function which is a mapping from  $(\mathbf{A}, h)$  to real values. Moreover, for a given  $\lambda$ , the optimal policy achieving  $\bar{L}^*(\lambda)$  will be

$$\begin{aligned} \pi_{\lambda}^*(\mathbf{A}, h) &= \arg \min_{\mathbf{w} \in \mathcal{W}} \{L(\mathbf{A}, h, \mathbf{w}); \lambda \\ &+ p_{\mathcal{H}}(\mathbf{h})h[V(\mathbf{A}, h); \lambda] + \sum_{\mathbf{A} \in \mathcal{A}} p_{\mathcal{A}}(\mathbf{A})V(\mathbf{A}, h); \lambda\}. \end{aligned} \quad (12)$$

From Lemma 2, we can see that  $\pi_{\lambda}^*$ , which is given by (12), depends on  $(\mathbf{A}, h)$  through the value function  $V(\mathbf{A}, h)$ . Determining  $V(\mathbf{A}, h)$  involves solving the Bellman equation in (11), for which there is no closed-form solution in general [24]. Numerical algorithms such as value iteration and policy iteration are usually computationally impractical to implement for an IoT due to the curse of dimensionality and they do not typically yield many design insights. Therefore, it is desirable to analyze the structural properties of  $\pi_{\lambda}^*$ , as we do next.

### C. Structural Analysis and Algorithm Design

First, we characterize the structural properties of  $\pi_{\lambda}^*$  for the unconstrained MDP in (8). Then, we propose a novel structure-aware optimal algorithm to obtain the optimal policy  $\pi^*$  for the CMDP in (5). Finally, we study the effects of the wireless channel fading.

1) *Optimality Properties*: By using the relative value iteration algorithm and the special structures of the AoI dynamics in (1) and (2), we can prove the following property.

**Lemma 3:** Given  $\lambda \geq 0$ ,  $V(\mathbf{A}, h); \lambda$  is non-decreasing with  $A_l$  and  $A_r$  for any  $h \in \mathcal{H}$ .

*Proof:* See Appendix A. ■

Then, we introduce the state-action Lagrange cost function, which is related to the right-hand side of the Bellman equation in (11) and is given by:

$$J(\mathbf{A}, h, \mathbf{w}); \lambda \triangleq L(\mathbf{A}, h, \mathbf{w}); \lambda + p_{\mathcal{H}}(\mathbf{h})h[V(\mathbf{A}, h); \lambda]. \quad (13)$$

We now define  $\Delta J_{\mathbf{w}, \mathbf{w}'}(\mathbf{A}, h); \lambda \triangleq J(\mathbf{A}, h, \mathbf{w}); \lambda - J(\mathbf{A}, h, \mathbf{w}'); \lambda$ . If  $\Delta J_{\mathbf{w}, \mathbf{w}'}(\mathbf{A}, h); \lambda \geq 0$ , we say that action  $\mathbf{w}$  *dominates* action  $\mathbf{w}'$  at state  $(\mathbf{A}, h)$  for a given  $\lambda$ . By Lemma 3, if  $\mathbf{w}$  dominates all other actions at state  $(\mathbf{A}, h)$  for a given  $\lambda$ , then we have  $\pi_{\lambda}^*(\mathbf{A}, h) = \mathbf{w}$ . Based on Lemma 3, we can obtain the following properties of  $\Delta J_{\mathbf{w}, \mathbf{w}'}(\mathbf{A}, h); \lambda$ .

**Lemma 4:** Given  $\lambda \geq 0$ , for any  $\mathbf{A} \in \mathcal{A}$ ,  $h \in \mathcal{H}$ , and  $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$ ,  $\Delta J_{\mathbf{w}, \mathbf{w}'}(\mathbf{A}, h); \lambda$  has the following properties:

- A) If  $\mathbf{w} = (0, 0)$ , then  $\Delta J_{\mathbf{w}, \mathbf{w}'}(\mathbf{A}, h); \lambda$  is non-decreasing with  $A_l$  for  $\mathbf{w}' = (1, 0)$  and non-decreasing with  $A_r$  for  $\mathbf{w}' = (0, 1)$  or  $(1, 1)$ .

- B) If  $w = ]0, 1[$  or  $]1, 1[$ , then  $\Delta J_{w, w^\infty} A, h; \lambda[$  is non-increasing with  $A_r$  for any  $w^\infty \neq w$ .
- C) If  $w = ]1, 0[$ , then  $\Delta J_{w, w^\infty} A, h; \lambda[$  is non-increasing with  $A_l$  for any  $w^\infty \neq w$ .

*Proof:* See Appendix B. ■

Lemma 4 follows from the special properties of the AoI dynamics and is essential for the characterization of the structural properties of  $\pi_\lambda^*$ . The property shown in Lemma 4 is similar to the diminishing-return property of multimodularity functions [25]. From Lemma 4, we can see that for a given  $\lambda$  and  $h$ , if action  $w$  dominates action  $w^\infty$  for some AoI state  $A$ , then  $w$  still dominates  $w^\infty$  for another AoI  $A^\infty$ , provided that  $A$  and  $A^\infty$  satisfy certain conditions such that  $\Delta J_{w, w^\infty} A^\infty, h; \lambda[ \geq \Delta J_{w, w^\infty} A, h; \lambda[ \geq 0$ . Before presenting the structure of  $\pi_\lambda^*$  in Theorem 1, we make the following definitions:

$$\Phi_w)A_r, h; \lambda[ \triangleq \{A_l \nmid A_l \in \mathcal{A}_l \text{ and } \Delta J_{w, w^\infty} A, h; \lambda[ \geq 0 \\ \forall w^\infty \in \mathcal{W} \text{ and } w^\infty \neq w\}, \quad (14)$$

$$\Psi_w)A_l, h; \lambda[ \triangleq \{A_r \nmid A_r \in \mathcal{A}_r \text{ and } \Delta J_{w, w^\infty} A, h; \lambda[ \geq 0 \\ \forall w^\infty \in \mathcal{W} \text{ and } w^\infty \neq w\}. \quad (15)$$

Then, we define:

$$\phi_w^+)A_r, h; \lambda[ \triangleq \begin{cases} \max \Phi_w)A_r, h; \lambda[, & \text{if } \Phi_w)A_r, h; \lambda[ \neq \emptyset, \\ \in, & \text{otherwise,} \end{cases} \quad (16)$$

$$\phi_w)A_r, h; \lambda[ \triangleq \begin{cases} \min \Phi_w)A_r, h; \lambda[, & \text{if } \Phi_w)A_r, h; \lambda[ \neq \emptyset, \\ +\infty, & \text{otherwise,} \end{cases} \quad (17)$$

$$\psi_w^+)A_l, h; \lambda[ \triangleq \begin{cases} \max \Psi_w)A_l, h; \lambda[, & \text{if } \Psi_w)A_l, h; \lambda[ \neq \emptyset, \\ \in, & \text{otherwise,} \end{cases} \quad (18)$$

$$\psi_w)A_l, h; \lambda[ \triangleq \begin{cases} \min \Psi_w)A_l, h; \lambda[, & \text{if } \Psi_w)A_l, h; \lambda[ \neq \emptyset, \\ +\infty, & \text{otherwise.} \end{cases} \quad (19)$$

**Theorem 1:** Given  $\lambda$ , for any  $A \in \mathcal{A}$  and  $h \in \mathcal{H}$ , there exists an optimal policy satisfying the following structural properties:

- A)  $\pi_\lambda^*)A, h[ = ]0, 0[$ , for all  $A \in \mathcal{A}_0; h; \lambda[ \triangleq \{A \nmid A_l \geq \phi_{0,0}^+)A_r, h; \lambda[, A_r \geq \psi_{0,0}^+)A_l, h; \lambda[$ .
- B)  $\pi_\lambda^*)A, h[ = ]0, 1[$  if  $A_r \geq \psi_{0,1}^+)A_l, h; \lambda[$ .
- C)  $\pi_\lambda^*)A, h[ = ]1, 0[$  if  $A_l \geq \phi_{1,0}^+)A_r, h; \lambda[$ .
- D)  $\pi_\lambda^*)A, h[ = ]1, 1[$  if  $A_r \geq \psi_{1,1}^+)A_l, h; \lambda[$ .

Theorem 1 characterizes the structural properties of the optimal policy  $\pi_\lambda^*$  of the unconstrained MDP in (8) for a given  $\lambda$ . Fig. 2 illustrates the analytical results of Theorem 1, where the optimal policy is computed numerically using policy iteration [26, Chapter 8.6]. Fig. 2 shows that, if the AoI state falls into the region of the black squares (i.e.,  $\mathcal{A}_0; h; \lambda[$ ), the device will remain idle and will not sample the physical process nor send the status update. Thus,  $\mathcal{A}_0; h; \lambda[$  is referred to as the *idle region*. For given  $A_l$ ,  $h$ , and  $\lambda$ , the scheduling of  $]0, 1[$  (or  $]1, 1[$ ) is threshold-based with respect to  $A_r$ . In other words, when  $A_r$  is small, it is not efficient to send a new status update to the destination, as a higher updating

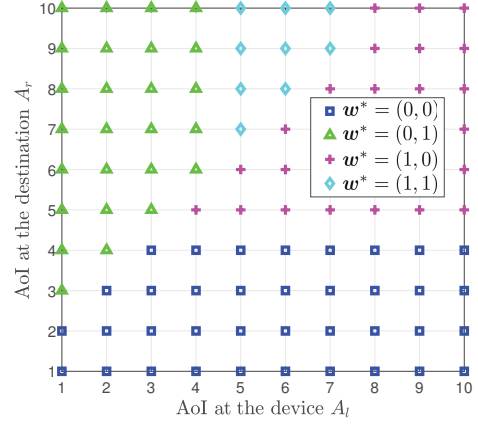


Fig. 2: Structure of the optimal policy  $\pi_\lambda^*$  for a given Lagrange multiplier  $\lambda$  and channel state  $h$ .  $\hat{A}_l = \hat{A}_r = 10$ .  $C_s = 2$ ,  $C_u)h[ = 3.5(h$ , where  $h = 1, 2$ ,  $C^{\max} = 3$ .

cost per age is consumed. Meanwhile, when  $A_r$  is large, it is more efficient to update the status, as the status packet at the destination becomes more outdated. For given  $A_r$ ,  $h$ , and  $\lambda$ , the scheduling of  $]1, 0[$  is threshold-based with respect to  $A_l$ . Hence when  $A_l$  is small, it is not efficient to sample the physical process, as a higher sampling cost per age is incurred. In contrast, when  $A_l$  is large, it is more efficient to generate a new status packet, as the status packet at the device becomes more outdated and less useful for the destination. These observations indicate that the zero-wait policy (i.e., transmit the status packet immediately after sampling) may be detrimental to the minimization of the AoI, because of the energy cost constraint. This is reminiscent of the result in [11], however, the work in [11] obtained this outcome because of the considered random service times. *These threshold-based properties reveal a fundamental tradeoff between the AoI at the destination and the sampling and updating costs.* Such structural properties provide valuable insights for the design of the sampling and updating processes in practical IoT systems. Here, we would like to emphasize that, although the threshold-based structures may look intuitive, it is challenging to prove these structures rigorously. This is due to the coupled two AoI states and the special AoI dynamics. Moreover, it is not always possible to fully characterize the structural properties of the optimal policy, e.g., the structure with respect to  $h$  and the structures of the thresholds, as the (generally required) key property of the value function, i.e., the multimodularity [25], does not hold for our value function.

2) *Algorithm Design:* By exploiting the results of Theorem 1, we first propose a structure-aware algorithm to compute the optimal policy  $\pi_\lambda^*$  for a given  $\lambda$ , and, then, we describe how to update  $\lambda$  and obtain the optimal policy  $\pi^*$ . Note that, although the exact values of the thresholds in Theorem 1 rely on the exact values of  $V)A, h; \lambda[$ , the threshold-based structure only relies on the properties of  $V)A, h; \lambda[$  and  $J)A, h, w; \lambda[$ . These properties can be exploited to reduce the computational complexity for obtaining the optimal policy, without knowing the exact values of the thresholds. In particular, by properties B)-D) in Theorem 1, we know that the optimal action for

---

**Algorithm 1** Structure-aware Policy Iteration Algorithm
 

---

- 1: Set  $\pi_{\lambda,0}^* \rceil A, h[ = \rceil 0, 0[$  for all  $\rceil A, h[ \in \mathcal{A} * \mathcal{H}$ , select reference state  $\rceil A', h'[$ , and set  $m = 0$ .
- 2: (Policy Evaluation) Given policy  $\pi_{\lambda,m}^*$ , compute the value  $\theta_{\lambda,m}$  and value function  $V_m \rceil A, h[$  from the linear system of equations<sup>3</sup>

$$\begin{cases} \theta_{\lambda,m} + V_m \rceil A, h; \lambda[ = L \rceil A, h, \pi_{\lambda,m}^* \rceil A, h[; \lambda[ \\ \quad + \prod_{h' \in \mathcal{H}} p_{\mathcal{H}}(h'|h) V_m \rceil A, h'; \lambda[, \forall \rceil A, h[ \\ V_m \rceil A', h'; \lambda[ = 0 \end{cases} \quad (21)$$

where  $A^\infty$  satisfies the AoI dynamics in (1) and (2) under the action  $\pi_{\lambda,m}^* \rceil A, h[$ .

- 3: (Structured Policy Improvement) Find a new policy  $\pi_{\lambda,m+1}^*$  for each  $\rceil A \in \mathcal{A}$  and  $h \in \mathcal{H}$ , where for each  $\rceil A, h[ \in \mathcal{A} * \mathcal{H}$ ,  $\pi_{\lambda,m+1}^* \rceil A, h; \lambda[$  is such that:
  - if  $\pi_{\lambda,m+1}^* \rceil A_l, A_r \rceil 1, h; \lambda[ = \rceil 0, 1[$ , **then**  $\pi_{\lambda,m+1}^* \rceil A, h; \lambda[ = \rceil 0, 1[$ .
  - else if  $\pi_{\lambda,m+1}^* \rceil A_l \rceil 1, A_r \rceil h; \lambda[ = \rceil 1, 0[$ , **then**  $\pi_{\lambda,m+1}^* \rceil A, h; \lambda[ = \rceil 1, 0[$ .
  - else if  $\pi_{\lambda,m+1}^* \rceil A_l, A_r \rceil 1, h; \lambda[ = \rceil 1, 1[$ , **then**  $\pi_{\lambda,m+1}^* \rceil A, h; \lambda[ = \rceil 1, 1[$ .
  - else

$$\begin{aligned} \pi_{\lambda,m+1}^* \rceil A, h; \lambda[ = \arg \min_{w \in \mathcal{W}} \Big\{ & L \rceil A, h, w; \lambda[ \\ & + \prod_{h' \in \mathcal{H}} p_{\mathcal{H}}(h'|h) V_m \rceil A, h'; \lambda[ \Big\} \quad (22) \end{aligned}$$

- 4: Go to Step 2 until  $\mu_{l+1}^* = \mu_l^*$ .
- 

a certain system state is still optimal for some other system states. In particular, we can see that, for all  $A$  and  $h$ ,

$$\begin{cases} \pi_{\lambda}^* \rceil A_l, A_r, h[ = \rceil 0, 1[ \quad \text{if} \quad \pi_{\lambda}^* \rceil A_l, A_r + 1, h[ = \rceil 0, 1[, \\ \pi_{\lambda}^* \rceil A_l, A_r, h[ = \rceil 1, 0[ \quad \text{if} \quad \pi_{\lambda}^* \rceil A_l + 1, A_r, h[ = \rceil 1, 0[, \\ \pi_{\lambda}^* \rceil A_l, A_r, h[ = \rceil 1, 1[ \quad \text{if} \quad \pi_{\lambda}^* \rceil A_l, A_r + 1, h[ = \rceil 1, 1[. \end{cases} \quad (20)$$

Therefore, to find  $\pi_{\lambda}^*$ , we only need to minimize in the right-hand side of (12) for some  $A$ , rather than for all  $A$ , which reduces the computational complexity. By incorporating (20) into a standard policy iteration algorithm, we can develop a structure-aware policy iteration algorithm, as shown in Algorithm 1. It can be seen that Algorithm 1 is a monotone policy iteration algorithm (see an example in [26, Chapter 8.11.2]), and thus, converges to the optimal policy  $\pi_{\lambda}^*$  [26, Theorem 8.6.6] (restated in Appendix H). Note that, when one of the “if” conditions in Step 3 of Algorithm 1 is satisfied for a certain system state, we can determine the optimal action immediately, without performing the minimization in (22). The computational complexity saving for each iteration in Algorithm 1 is  $O(\mathfrak{M} \mathfrak{H} \mathfrak{A} \mathfrak{H}^2)$  [27]. This is reasonable since the complexity saving grows exponentially with the state space.

From (10), we know that obtaining  $\pi^*$  requires computing the two Lagrange multipliers  $\lambda_1$  and  $\lambda_2$ , and the randomization

parameter  $\alpha$ . As in [23] and [28], we set  $\lambda_1 = \lambda^* - \eta$  and  $\lambda_2 = \lambda^* + \eta$ , where the perturbation parameter  $\eta$  is some small constant and  $\lambda^*$  is the optimal Lagrange multiplier satisfying  $\lambda^* = \min \{ \lambda : \bar{C} \pi_{\lambda}^* [ \geq C^{\max} \}$ . By using the Robbins-Monro algorithm [29], which is a stochastic gradient-based algorithm, the Lagrange multiplier is updated according to  $\lambda_{m+1} = \lambda_m + \epsilon_m \left( \bar{C} \pi_{\lambda_m}^* [ - C^{\max} \right)$ , where the step  $\epsilon_m = \frac{1}{m}$  and  $\lambda_1$  is initialized with a sufficiently large number. The generated sequence  $\{ \lambda_m \}$  converges to the optimal Lagrange multiplier  $\lambda^*$  [29]. Then, the randomization parameter  $\alpha$  is given by:  $\alpha = \bar{C}^{\max} \bar{C} \pi_{\lambda_2}^* [ / ((\bar{C} \pi_{\lambda_1}^* [ - \bar{C} \pi_{\lambda_2}^* [))$ . Here,  $\alpha$  is chosen such that  $\alpha \bar{C} \pi_{\lambda_1}^* [ + (1 - \alpha) \bar{C} \pi_{\lambda_2}^* [ = C^{\max}$ . Then, for some perturbation parameter  $\eta$ , by [23, Theorem 4.3], we have that  $\pi^* = \alpha \pi_{\lambda_1}^* + (1 - \alpha) \pi_{\lambda_2}^*$  is the optimal policy for the CMDP. Note that,  $\alpha$  is guaranteed to lie in  $\rceil 0, 1[$  since  $\bar{C} \pi_{\lambda}^* [$  is non-increasing with  $\lambda$  [23]. So far, we have characterized the structural properties of the optimal policy  $\pi_{\lambda}^*$  and developed a structure-aware optimal algorithm.

3) *Effects of Wireless Channel Dynamics*: Next, we study the influence of the wireless channel fading distribution on the optimal average AoI at the destination. The results are established by using the stochastic dominance relations of random variables. From [21] and [30], we present the following definition.

**Definition 2:** Let  $x) \gamma[$  be a random variable with the support on the set  $\mathcal{X}$  according to a probability measure  $\mu) \gamma[$  parameterized by some  $\gamma$ .  $x) \gamma_1[$  is said to *stochastically dominate*  $x) \gamma_2[$  on the set of functions  $\mathcal{F}$ , or  $x) \gamma_1[ \prec_{\mathcal{F}} x) \gamma_2[$ , if  $\mathbb{E}[f(x) \gamma_1][ \geq \mathbb{E}[f(x) \gamma_2][$ , for all functions  $f \in \mathcal{F}$ . If  $\mathcal{F}$  is the set of increasing functions, then  $\prec_{\mathcal{F}}$  corresponds to the *first-order stochastic dominance*. If  $\mathcal{F}$  is the set of increasing and concave functions, then  $\prec_{\mathcal{F}}$  corresponds to the *second-order stochastic dominance*.

Consider two channels  $I$  and  $J$ . Let  $h^I \in \mathcal{H}$  and  $h^J \in \mathcal{H}$  be random variables with the fading distributions  $p_{\mathcal{H}}^I(h)$  and  $p_{\mathcal{H}}^J(h)$  for channels  $I$  and  $J$ , respectively.

**Theorem 2:** If  $h^I$  first-order stochastically dominates  $h^J$ , then we have

$$\bar{A}_r^{I*} \geq \bar{A}_r^{J*}, \quad (23)$$

where  $\bar{A}_r^{I*}$  and  $\bar{A}_r^{J*}$  are the optimal average AoI at the destination under channels  $I$  and  $J$ , respectively.

*Proof:* See Appendix D. ■

Theorem 2 demonstrates that channels with larger mean channel gain can achieve smaller AoI at the destination under the same resource constraint. Following the proof of Theorem 2, we have the following corollary for second-order stochastically dominating channels.

**Corollary 1:** If  $h^I$  second-order stochastically dominates  $h^J$  and  $C_u) h[$  is decreasing and convex with  $h$ , then the optimal average AoI at the destination under channel  $I$  is smaller than that under channel  $J$ , i.e.,  $\bar{A}_r^{I*} \geq \bar{A}_r^{J*}$ .

Note that, for the transmission cost  $C_u) h[$  defined according to the Shannon's formula (e.g., in [21]), it can be easily seen that  $C_u) h[$  satisfies the conditions of Corollary 1. If  $h^I$  has the same mean as  $h^J$ , by Definition 2, the second-order stochastic dominance of  $h^J$  over  $C_u) h[$  indicates that  $h^I$  has smaller variance (i.e., less scattering) than  $h^J$ . Therefore,

<sup>3</sup>The solution to (21) can be derived using Gaussian elimination or the relative value iteration method [24].



Corollary 1 reveals that channels with less scattering and the same mean channel gain can achieve a smaller AoI at the destination under the same resource constraint. The results obtained in Theorem 2 and Corollary 1 reveal the fundamental monotone dependency of the optimal AoI at the destination on the transmission probability distribution of the CMDP in (5).

Thus far, we have analyzed the optimality properties for the case of a single IoT device so that to gain a deep understanding of the behavior of the optimal sampling and updating policy for the real-time monitoring system. Next, we consider a more general scenario in which there are multiple IoT devices. For such a scenario, the system state space is much larger than that for the case of a single IoT device, as it grows exponentially with the number of the devices. This hinders the structural analysis of the optimal policy and the design of an optimal algorithm with low-complexity. Therefore, we will focus on the design of a low-complexity suboptimal solution for the case of multiple IoT devices.

### III. SEMI-DISTRIBUTED SUBOPTIMAL SAMPLING AND UPDATING CONTROL FOR MULTIPLE IOT DEVICES

#### A. System Model and Problem Formulation

We now extend the real-time monitoring system in Section II to a more general scenario, in which a set  $\mathcal{K}$  of  $K$  IoT devices sample the associated physical processes and update the status packets to a common destination. Hereinafter, with some notation abuse, for each IoT device  $k \in \mathcal{K}$ , we denote by  $\mathbf{A}_k(t) \triangleq (A_{l,k}(t), A_{r,k}(t)) \in \mathcal{A}_k$ ,  $\mathbf{h}_k(t) \triangleq (s_k(t), u_k(t)) \in \mathcal{W}_k$  the AoI state, the channel state, and the control action vector at slot  $t$ , respectively. Under action  $\mathbf{w}_k(t)$ , the AoI state  $\mathbf{A}_k(t)$  for each IoT device  $k$  is updated in the same manner of (1) and (2). We define  $\mathbf{A}(t) \triangleq (\mathbf{A}_k(t))_{k \in \mathcal{K}} \in \mathcal{A} \triangleq \prod_{k \in \mathcal{K}} \mathcal{A}_k$ ,  $\mathbf{h}(t) \triangleq (\mathbf{h}_k(t))_{k \in \mathcal{K}} \in \mathcal{H} \triangleq \prod_{k \in \mathcal{K}} \mathcal{H}_k$ , and  $\mathbf{w}(t) \triangleq (\mathbf{w}_k(t))_{k \in \mathcal{K}} \in \mathcal{W}$  as the system AoI state, the system channel state, and the system control action at slot  $t$ , respectively. Let  $C_{s,k}$  and  $C_{u,k}h_k$  be the sampling cost and the updating cost under channel state  $h_k$  of IoT device  $k$ , respectively. We assume that, the channel state processes  $\{h_k(t)\}_{k \in \mathcal{K}}$  at the devices are mutually independent. As in [9], we consider that, in each slot, the multiple IoT devices cannot update their status packets concurrently; otherwise collisions occur and no status packets will be transmitted to the destination successfully. Thus, *different from the case of a single IoT device, the updating process of the multiple IoT devices should be carefully scheduled to avoid such collisions.* Mathematically, we have  $\prod_{k \in \mathcal{K}} u_k(t) \geq 1$ , for all  $t$ . Then, we define  $\mathcal{W} \triangleq \mathcal{S} * \mathcal{U}$  as the feasible system control action space, where  $\mathcal{S} \triangleq \{0, 1\}^K$  and  $\mathcal{U} \triangleq \{u_k \in [0, 1] \mid \forall k \in \mathcal{K} \text{ and } \prod_{k \in \mathcal{K}} u_k \geq 1\}$ . Note that, the proposed analytical framework and algorithm design can be readily extended to support the orthogonal frequency division multiple access (OFDMA) mode, in which multiple IoT devices can update their status at the same time without collisions over different non-overlapping channels [31].

Similar to the single device case, given an observed system AoI state  $\mathbf{A}$  and system channel state  $\mathbf{h}$ , the system control action  $\mathbf{w}$  is derived as per the following policy.

**Definition 3:** A *feasible stationary sampling and updating policy*  $\pi = (\pi_s, \pi_u)$  is defined as a mapping from the system AoI state and the system channel state  $(\mathbf{A}, \mathbf{h}) \in \mathcal{A} * \mathcal{H}$  to the feasible system control action of the IoT devices  $\mathbf{w} \in \mathcal{W}$ , where  $\pi_s(\mathbf{A}, \mathbf{h}) = \mathbf{s}$  and  $\pi_u(\mathbf{A}, \mathbf{h}) = \mathbf{u}$ .

Under a given stationary unichain policy  $\pi$ , the average AoI at the destination and the average energy cost for each IoT device  $k$  are respectively given by:

$$\bar{A}_r(\pi) \triangleq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E} [A_{r,k}(t)] \quad (24)$$

$$\bar{C}_k(\pi) \triangleq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E} [C_k(\mathbf{w}_k(t))] \quad \forall k \in \mathcal{K}, \quad (25)$$

where  $C_k(\mathbf{w}_k(t)) \triangleq s_k(t)C_{s,k} + u_k(t)C_{u,k}h_k(t)$  and the expectation is taken with respect to the measure induced by the policy  $\pi$ .

We want to find the optimal *feasible* sampling and updating policy that minimizes the average AoI at the destination, under an average energy cost constraint for each IoT device, as follows:

$$\bar{A}_r^* \triangleq \min_{\pi} \bar{A}_r(\pi), \quad (26a)$$

$$\text{s.t. } \bar{C}_k(\pi) \leq C_k^{\max}, \forall k \in \mathcal{K}. \quad (26b)$$

To obtain the optimal policy  $\pi^*$  in (26), we again introduce the Lagrangian for a given vector of Lagrange multipliers  $\lambda \triangleq (\lambda_k)_{k \in \mathcal{K}}$ , given by:

$$\bar{L}(\pi; \lambda) \triangleq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E} [L(\mathbf{A}(t), \mathbf{h}(t), \mathbf{w}(t); \lambda)], \quad (27)$$

where  $L(\mathbf{A}(t), \mathbf{h}(t), \mathbf{w}(t); \lambda) \triangleq \sum_{k=1}^K A_{r,k}(t) + \lambda_k C_k(\mathbf{w}_k(t))$  and  $C_k^{\max}$  is the Lagrange cost at slot  $t$ . Then, the corresponding unconstrained MDP for a given  $\lambda$  will be:

$$\bar{L}^*(\lambda) \triangleq \min_{\pi} \bar{L}(\pi; \lambda), \quad (28)$$

where  $\bar{L}^*(\lambda)$  is the minimum average Lagrange cost achieved by the optimal policy  $\pi_{\lambda}^*$  for a given  $\lambda$ . The optimal average AoI at the destination in (26a) is given by  $\bar{A}_r^* = \max_{\lambda} \bar{L}^*(\lambda)$ . In the following lemma, we summarize the solution to the unconstrained MDP in (28).

**Lemma 5:** For any  $\lambda$ , there exists  $(\theta_{\lambda}, Q)A, \mathbf{h}, \mathbf{u}; \lambda$  satisfying:

$$\begin{aligned} \theta_{\lambda} + Q)A, \mathbf{h}, \mathbf{u}; \lambda &= \min_{s \in \mathcal{S}} \left\{ L(\mathbf{A}, \mathbf{h}, \mathbf{w}; \lambda) \right. \\ &\quad \left. + \min_{\mathbf{h} \in \mathcal{H}} p_{\mathcal{H}}(\mathbf{h}) \min_{\mathbf{u} \in \mathcal{U}} Q)A^{\infty}, \mathbf{h}^{\infty}, \mathbf{u}^{\infty}; \lambda \right\}, \quad (29) \end{aligned}$$

where  $A^{\infty}$  satisfies the AoI dynamics in (1) and (2) for each IoT device,  $\theta_{\lambda} = \bar{L}^*(\lambda)$  is the optimal value to (28) for all initial state  $(\mathbf{A}(1), \mathbf{h}(1), \mathbf{u}(1))$ , and  $Q)X$  is the Q-factor which is a mapping from  $(\mathbf{A}, \mathbf{h}, \mathbf{u})$  to real values. Moreover, for a given  $\lambda$ , the optimal policy achieving the optimal value  $\bar{L}^*(\lambda)$  is given by  $\pi_{\lambda}^*(\mathbf{A}, \mathbf{h}) = (\pi_{\lambda,s}^*, \pi_{\lambda,u}^*)A, \mathbf{h}$ , where  $\pi_{\lambda,s}^*(\mathbf{A}, \mathbf{h})$  attains the minimum of the right-hand side of (29) and  $\pi_{\lambda,u}^*(\mathbf{A}, \mathbf{h}) = \arg \min_{\mathbf{u} \in \mathcal{U}} Q)A, \mathbf{h}, \mathbf{u}; \lambda$ .

*Proof:* See Appendix E. ■

From Lemma 5, we can see that the optimal sampling and updating action depends on the Q-factor  $Q_k)A, h, u; \lambda[$  and the  $K$  Lagrange multipliers. For a given  $\lambda$ , obtaining the Q-factor  $Q_k)$  requires solving the Bellman equation in (29), which suffers from the curse of the dimensionality due to the exponential growth of the cardinality of the system state space ( $\mathcal{A} * \mathcal{H} = \sum_{k=1}^K A_{l,k}^{\max} A_{r,k}^{\max} \mathcal{H}_k$ ). Even if we could obtain the optimal Q-factors by solving (29), the derived control will be centralized thus requiring a knowledge of the system AoI states and channel states at each slot by the destination node, which is highly undesirable. Moreover, the optimal policy of the CMDP in (26) is a randomized stationary policy with a degree of randomization no greater than  $K$  [32], and, thus, may not be very suitable for practical implementations. Note that, since we need to jointly control the sampling and updating processes, our problem cannot be cast into a restless multi-armed bandit problem (RMAB) [33] as is often done in the literature<sup>4</sup>, thus rendering the existing low-complexity solutions (e.g., [7], [8], and [9]) not applicable. Therefore, we next introduce a novel semi-distributed low-complexity algorithm to obtain a deterministic suboptimal sampling and updating policy.

### B. Algorithm Design

In this subsection, we first approximate the Q-factor  $Q_k)A, h, u; \lambda[$  by the sum of the per-device Q-factor  $Q_k)A_k, h_k, u_k; \lambda_k[$ . Based on the approximated Q-factor, we propose a semi-distributed sampling and updating policy, inspired by [34]. Then, we develop an online learning algorithm that enables each device to determine its per-device Q-factor and the associated Lagrange multiplier based on the observation of its AoI and channel states. Finally, we show that the proposed semi-distributed learning algorithm converges to the proposed suboptimal policy.

1) *Semi-Distributed Sampling and Updating Control*: To reduce the complexity for obtaining the optimal Q-factor, we adopt the linear approximation architecture [24] to approximate the Q-factor in (29) by the sum of the per-device Q-factor  $Q_k)A_k, h_k, u_k; \lambda_k[$ :

$$Q_k)A, h, u; \lambda[ \approx \sum_{k=1}^K Q_k)A_k, h_k, u_k; \lambda_k[, \quad (30)$$

where  $Q_k)A_k, h_k, u_k; \lambda_k[$  satisfies the following per-device Q-factor fixed point equation of each IoT device  $k$  for each given  $\lambda_k$ :

$$\begin{aligned} \theta_k + Q_k)A_k, h_k, u_k; \lambda_k[ = & \min_{s_k \in \{0,1\}} \left\{ L_k)A_k, h_k, s_k, u_k; \lambda_k[ \right. \\ & + \left. \sum_{h_k^\infty \in \mathcal{H}_k} p_{\mathcal{H}_k}(h_k^\infty) \min_{u_k^\infty \in \{0,1\}} Q_k)A_k^\infty, h_k^\infty, u_k^\infty; \lambda_k[ \right. \\ & \left. \forall A_k, h_k, u_k \in \{0,1\} * \mathcal{H}_k * \{0,1\} \right\}. \quad (31) \end{aligned}$$

Here,  $L_k)A_k, h_k, s_k, u_k; \lambda_k[ = A_{r,k} + \lambda_k)C_k)w_k[$  is the per-device Lagrange cost for IoT device  $k$ . Then, according to

Lemma 5, the destination node determines the updating control policy of all IoT devices based on the linear approximation in (30), given by:

$$\hat{\pi}_{\lambda, u}^*)A, h[ = \arg \min_{u \in \{0,1\}} \sum_{k=1}^K Q_k)A_k, h_k, u_k; \lambda_k[. \quad (32)$$

Problem (32) can be solved by a brute-force search with complexity of  $O(\mathcal{H})$ . In particular, each IoT device  $k$  observes its AoI state  $A_k$  and channel state  $h_k$  and reports its current per-device Q-factor  $Q_k)A_k, h_k, u_k; \lambda_k[$ ,  $u_k \in \{0,1\}$  to the destination node. Then, the destination node determines the system updating action  $\hat{u}^* = \hat{\pi}_{\lambda, u}^*)A, h[$  according to (32) and broadcasts the updating action  $\hat{u}^* = \hat{u}_k^*)_{k \in \mathcal{K}}$  to the  $K$  IoT devices. Based on the local observation of  $A_k$  and  $h_k$ , as well as the updating action  $\hat{u}_k^*$ , each IoT device  $k$  decides its sampling action  $\hat{s}_k^*$ , which minimizes the right-hand side of (31):

$$\begin{aligned} \hat{s}_k^* = & \arg \min_{s_k \in \{0,1\}} \left\{ L_k)A_k, h_k, s_k, \hat{u}_k^*; \lambda_k[ \right. \\ & + \left. \sum_{h_k^\infty \in \mathcal{H}_k} p_{\mathcal{H}_k}(h_k^\infty) \min_{u_k^\infty \in \{0,1\}} Q_k)A_k^\infty, h_k^\infty, u_k^\infty; \lambda_k[ \right\}. \quad (33) \end{aligned}$$

Note that, to obtain the proposed suboptimal policy  $\hat{\pi}^*$  in (32) and (33), we need to compute  $Q_k)A_k, h_k, u_k; \lambda_k[$  by solving (31) for all IoT devices, which is a total of  $O(\prod_{k=1}^K A_{l,k}^{\max} A_{r,k}^{\max} \mathcal{H}_k)$  values. However, to obtain the optimal policy  $\pi^*$ , computing  $Q_k)A, h, u; \lambda[$  by solving (29) requires a total of  $O(\sum_{k=1}^K A_{l,k}^{\max} A_{r,k}^{\max} \mathcal{H}_k)$  values. Therefore, the complexity of the proposed suboptimal policy decreases from exponential with  $K$  to linear with  $K$ .<sup>5</sup>

2) *Online Stochastic Learning and Convergence Analysis*: We observe that the proposed semi-distributed policy  $\hat{\pi}^*$  requires the knowledge of the per-device Q-factor and the associated Lagrange multiplier, which is challenging to obtain. Thus, we propose an online learning algorithm to estimate  $Q_k)A_k, h_k, u_k; \lambda_k[$  and  $\lambda_k$  at each IoT device  $k$ .

For IoT device  $k$ , based on the locally observed AoI state  $A_k)t[$ , channel state  $h_k)t[$ , the updating action  $\hat{u}_k^*)t[$  from the destination node, and the sampling action  $\hat{s}_k^*)t[$ , the per-device Q-factor and the Lagrange multiplier are respectively updated according to

$$\begin{aligned} Q_k^{t+1})A_k, h_k, u_k; \lambda_k^t[ &= Q_k^t)A_k, h_k, u_k; \lambda_k^t[ + \epsilon_{q,k}^{v_k^t})A_k, h_k, u_k[ \quad (34) \\ &* )F_k)A_k, h_k, u_k, \hat{s}_k^*)t[; \lambda_k^t[ - F_k)A_k^r, h_k^r, u_k^r, \hat{s}_k^*)t^r[; \lambda_k^t[ \\ &Q_k^t)A_k, h_k, \hat{u}_k^*; \lambda_k^t[ \left\{ * \mathbb{1}_{\hat{u}_k^*)A_k)t[, h_k)t[, \hat{u}_k^*)t[} = )A_k, h_k, u_k[[, \right. \\ &\left. \lambda_k^{t+1} = \lambda_k^t + \epsilon_{\lambda,k}^t)C_k)w_k)t[ - C_k^{\max}[\cdot, \right. \quad (35) \end{aligned}$$

where  $\mathbb{1}_{\hat{u}_k^*)A_k)t[, h_k)t[, \hat{u}_k^*)t[}$  is the indicator function,  $v_k^t)A_k, h_k, u_k[ \triangleq \prod_{\tau=1}^t \mathbb{1}_{\hat{u}_k^*)A_k)\tau[, h_k)\tau[, \hat{u}_k^*)\tau[}$  is the number of updates of the state-action pair  $)A_k, h_k, u_k[$  till  $t$ ,  $F_k)A_k, h_k, u_k, s_k; \lambda_k^t[ \triangleq L_k)A_k, h_k, s_k, u_k; \lambda_k^t[ +$

<sup>5</sup>The approximation error analysis of the linear approximation in (30) (which is a feature-based method) remains an open problem for CMDPs. Thus, we only provide numerical comparisons to illustrate its performance in the simulations.

<sup>4</sup>In general, RMAB only works for the problem with only one type of control actions.



---

**Algorithm 2** Semi-Distributed Sampling and Updating Learning Algorithm.

---

- 1: **Initialization:** Set  $t = 1$ . Each IoT device initializes its per-device Q-factor  $Q_k^t$  and Lagrange multiplier  $\lambda_k^t$ .
  - 2: **Updating control at the destination:** At slot  $t$ , each IoT device  $k$  reports  $\{Q_k\}A_k, h_k, u_k; \lambda_k\}t[$  to the destination node. Then, the destination node determines the system updating action according to (32) and broadcast the updating action  $\hat{u}^\bullet\}t[ = \hat{u}_k^\bullet\}t[[_{k \in \mathcal{K}}$  to the  $K$  IoT devices.
  - 3: **Sampling control at each IoT device:** Based on the updating action  $u_k^\bullet\}t[$ , each IoT device  $k$  decides its sampling action  $\hat{s}_k^\bullet\}t[$  according to (33).
  - 4: **Per-device Q-factor and Lagrange multiplier update at each IoT device:** Based on the current observations  $A_k\}t[$  and  $h_k\}t[$ , each IoT device  $k$  updates the per-device Q-factor  $Q_k^{t+1}A_k, h_k, u_k; \lambda_k^t[$  and  $\lambda_k^{t+1}$  according to (34) and (35), respectively.
  - 5: Set  $t \rightarrow t + 1$  and go to Step 2 until the convergence of  $Q_k^t$  and  $\lambda_k^t$ .
- 

$\prod_{h_k \in \mathcal{H}_k, p_{\mathcal{H}_k}} h_k\} \min_{u_k \in \{0,1\}} Q_k^t A_k^\infty, h_k^\infty, u_k^\infty, \lambda_k^t[$  with  $A_k^\infty$  and  $A_k$  satisfying the relations in (1) and (2),  $A_k^r, h_k^r, u_k^r[$  is some fixed reference state-action pair,  $t^r \triangleq \sup\{t \mid \hat{u}_k^\bullet\}t[[ = A_k^r, h_k^r, u_k^r[[$ , and  $[x]^+ = \max\{x, 0\}$ .  $\{\epsilon_{q,k}^t\}$  and  $\{\epsilon_{\lambda,k}^t\}$  are the sequences of step sizes satisfying:

$$\begin{aligned} \epsilon_{q,k}^t &= \epsilon, \quad \epsilon_{q,k}^t > 0, \quad \lim_{t \rightarrow \infty} \epsilon_{q,k}^t = 0, \quad \epsilon_{\lambda,k}^t = \epsilon, \quad \epsilon_{\lambda,k}^t > 0, \\ \lim_{t \rightarrow \infty} \epsilon_{\lambda,k}^t &= 0, \quad \sum_t \epsilon_{q,k}^t [^2 + \sum_t \epsilon_{\lambda,k}^t [^2 < \epsilon, \quad \text{and} \quad \lim_{t \rightarrow \infty} \frac{\epsilon_{\lambda,k}^t}{\epsilon_{q,k}^t} = 0. \end{aligned} \quad (36)$$

Here, (34) is formulated following the asynchronous relative value Q-learning algorithm [35].

From (34) and (35), to implement the proposed online learning algorithm at each IoT device, we only need the local AoI and channel states, as well as the updating control action from the destination. The proposed algorithm is illustrated in Algorithm 2. It can be seen that the proposed algorithm is essentially a grant-based uplink transmission protocol (see examples in [36], [37]), which involves the exchange of messages between the IoT devices and the destination node. However, this will not incur any notable overhead, because in each slot, each IoT device needs to only transmit a few bits to exchange its per-device Q-factor value with the destination. Note that, we need to update both the per-device Q-factors and the Lagrange multipliers simultaneously. Thus, conventional value iteration and policy iteration algorithms [24], and the Q-learning algorithm under which the Lagrange multipliers are determined offline [38] are not applicable to our case.

Now, we show the almost-sure convergence of Algorithm 2. From (36), we can see that the per-device Q-factor and the Lagrange multiplier are updated concurrently, albeit over two different timescales [39]. During the update of the per-device Q-factor (timescale I), we have  $\lambda_k^{t+1} - \lambda_k^t = O(\epsilon_{\lambda,k}^t) = o(\epsilon_{q,k}^t)$ ,

and thus,  $\lambda_k^t$  can be seen as quasi-static [39] when updating  $Q_k^t A_k, h_k, u_k; \lambda_k^t$  in (34). We first have the following lemma on the convergence of the per-device Q-factor learning over timescale I.

**Lemma 6:** For step sizes  $\{\epsilon_{q,k}^t\}$  and  $\{\epsilon_{\lambda,k}^t\}$  satisfying the conditions in (36), the update of the per-device Q-factor at each IoT device  $k$  converges almost surely to the solution of the fixed point equation in (31), under any initial per-device Q-factor  $Q_k^1$  and the Lagrange multiplier vector  $\lambda$ , i.e.,  $\lim_{t \rightarrow \infty} Q_k^t A_k, h_k, u_k; \lambda_k^t = Q_k^\epsilon A_k, h_k, u_k; \lambda_k^\epsilon$ , a.s.,  $\forall A_k, h_k, u_k, k$ .

*Proof:* See Appendix F. ■

During the update of the Lagrange multiplier (timescale II) in (35), the per-device Q-factor can be seen as nearly equilibrated [39]. Then, we have the following convergence result.

**Lemma 7:** The update of the vector of the Lagrange multipliers  $\lambda$  converges almost surely, i.e.,  $\lim_{t \rightarrow \infty} \lambda^t = \lambda^\infty$  a.s., where the policy under  $\lambda^\infty$  satisfies the constraints in (26b).

*Proof:* See Appendix G. ■

Based on Lemma 6 and Lemma 7, we summarize the convergence of the proposed semi-distributed online sampling and updating algorithm in Algorithm 2 in the following Theorem.

**Theorem 3:** For step sizes  $\{\epsilon_{q,k}^t\}$  and  $\{\epsilon_{\lambda,k}^t\}$  satisfying the conditions in (36), the iterations of the per-device Q-factor and the Lagrange multipliers in Algorithm 2 converge w.p. 1, i.e.,  $Q_k^t, \lambda_k^t \Rightarrow Q_k^\epsilon, \lambda_k^\epsilon$  almost surely, for each IoT device  $k$ , where  $Q_k^\epsilon, \lambda_k^\epsilon$  satisfies the fixed-point equation in (31) and the sampling and updating policy under  $Q_k^\epsilon, \lambda_k^\epsilon$  satisfies the average energy cost constraints in (26b).

In a nutshell, we have proposed a low-complexity semi-distributed learning algorithm to find a suboptimal sampling and updating policy so as to minimize the average AoI at the destination for the case of multiple IoT devices. The proposed semi-distributed learning algorithm can be implemented at each device, requiring only the local knowledge and simple signaling from the destination, and, thus, is highly desirable for practical implementations.

## IV. SIMULATION RESULTS AND ANALYSIS

### A. Case of A Single IoT Device

We first illustrate the structural properties of the optimal sampling and updating policy for the single IoT device case. In the simulations, we set  $\mathcal{H} = \{0.0131, 0.0418, 0.0753, 0.1157, 0.1661, 0.2343, 0.3407, 0.6200\}$  and the corresponding probabilities are  $p_{\mathcal{H}} = [1, 1, 2, 3, 3, 2, 1, 1] \cdot 14$  [40]. Similar to [40], we assume that the updating cost is  $C_u h = C_u(h)$ , where  $C_u = 0.2$ . For the sampling cost, we adopt the local-computing model in [41] and assume that  $C_s = 0.2$ . We set the upper limits of the AoI at the device and the AoI at the destination  $\hat{A}_l$  and  $\hat{A}_r$  be 10.

Fig. 3 illustrates the effects of the sampling and updating costs on the structural properties of the optimal policy  $\pi_\lambda^*$  for a given  $\lambda$  as shown in Theorem 1. In particular, Fig. 3(a) shows the relationship between the threshold  $\phi_{j,1,0} A_r, h; \lambda$

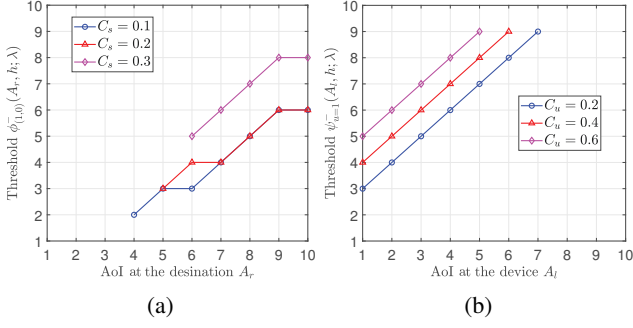


Fig. 3: Impacts of the sampling and updating costs on structures of the optimal policy  $\pi_\lambda^\star$  for a given  $\lambda$  in the single IoT device case.  $\lambda = 0.1$ . (a) Sampling cost.  $h = 0.0418$ . (b) Updating cost.  $h = 0.1157$ .

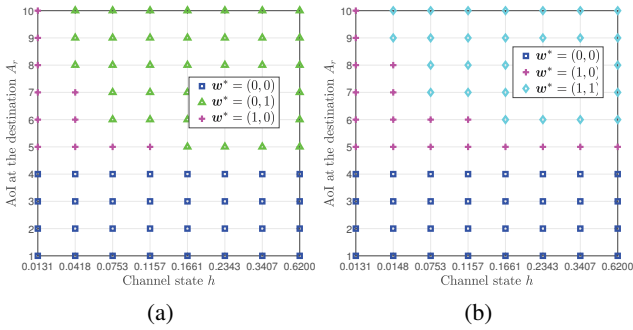


Fig. 4: Structure of the optimal policy  $\pi_\lambda^\star$  for given  $A_l$  and  $\lambda$ .  $\lambda = 0.1$ . (a)  $A_l = 4$ . (b)  $A_l = 5$ .

of choosing action  $)1, 0[$  and  $A_r$  under different sampling costs  $C_s$ , for given  $h$  and  $\lambda$ . From Fig. 3(a), we can see,  $\phi_{1,0}^-(A_r, h; \lambda)$  is non-decreasing with  $C_s$ . This indicates that the IoT device is unlikely to sample the physical process, if the sampling cost is high. Fig. 3(b) shows the relationship between  $\psi_{u=1}^-(A_l, h; \lambda) \triangleq \min\{\psi_{0,1}^-(A_l, h; \lambda), \psi_{1,1}^-(A_l, h; \lambda)\}$  and  $A_l$  under different updating costs  $C_u$ , for given  $h$  and  $\lambda$ . According to Theorem 1, if  $A_r \geq \psi_{u=1}^-(A_l, h; \lambda)$ , then the optimal updating action is  $u = 1$ , as  $\pi_\lambda^\star(A, h) = )0, 1[$  or  $)1, 1[$ . We observe that,  $\phi_{1,0}^-(A_r, h; \lambda)$  is non-decreasing with  $C_u$ . This indicates that the IoT device is not willing to send the status packet to the destination, if the updating cost is high.

Fig. 4 shows the structure of the optimal sampling and updating policy  $\pi_\lambda^\star$  for given  $A_l$  and  $\lambda$ . From Fig. 4(a) and Fig. 4(b), we can see that the scheduling of action  $)0, 1[$  or action  $)1, 1[$  is threshold-based with respect to the channel state  $h$ . In particular, Fig. 4 shows that, if the channel state is poor, it is not efficient for the IoT device to send the status packet to the destination, as a high updating cost will be incurred. Therefore, the optimal policy can fully exploit the random nature of the wireless channel by seizing good transmission opportunities to optimize the system performance. We also notice that the optimal actions  $)0, 1[$  and  $)1, 1[$  do not concurrently appear in the whole state space of  $)A_r, h[$ , under a given  $A_l$ . This is due to the fact that the decisions of choosing  $)0, 1[$  or  $)1, 1[$  are threshold-based with respect to  $A_r$  and  $h$ , as seen in the upper right corners of Fig. 4(a) and Fig. 4(b).

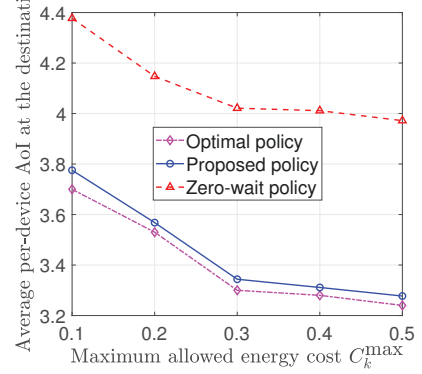


Fig. 5: Performance comparison among the optimal policy, the proposed semi-distributed policy, and the zero-wait baseline policy.  $K = 2$ ,  $A_{l,k}^{\max} = A_{r,k}^{\max} = 20, \forall k = 1, 2$ .

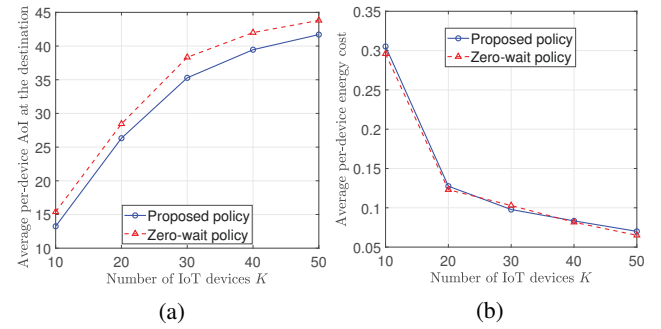


Fig. 6: Performance comparison between the proposed semi-distributed policy and the zero-wait baseline policy.  $C_k^{\max} = 0.3, \forall k$ .  $A_{l,k}^{\max} = A_{r,k}^{\max} = 100, \forall k$ . (a) Average per-device AoI at the destination. (b) Average per-device energy cost.

## B. Case of Multiple IoT Devices

Next, we evaluate the performance of the proposed semi-distributed online sampling and updating policy in Algorithm 2. The system parameters are analogous to those for the single IoT device case. For each device  $k$ , the sampling cost  $C_{s,k}$  is randomly selected from  $]0.2, 0.3[$  and the updating cost is  $C_{u,k}(h_k) = C_{u,k}(h_k)$ , where  $C_{u,k}$  is randomly selected from  $]0.2, 0.3[$ . We assume that the channel statistics of all IoT devices are the same, as given in Section IV-A. For comparison, consider a zero-wait baseline policy, i.e., in each slot, if an IoT device is scheduled to update its status packet, then it will sample the physical process immediately, which takes one slot. For the zero-wait baseline policy, the updating control and the updates of the per-device Q-factors and the Lagrange multipliers are similar to those of the proposed suboptimal policy, i.e., Step 2 and Step 4 in Algorithm 2. This is a commonly used baseline in the literature on AoI minimization, e.g., see [11] and references therein.

In Fig. 5, we compare the average AoI at the destination, resulting from the optimal policy, the proposed semi-randomized policy, and the zero-wait baseline policy, for two IoT devices under different values of  $C_k^{\max}$ . From Fig. 5, we can see that the proposed semi-distributed policy achieves a near-optimal performance and significantly outperforms the zero-wait policy.

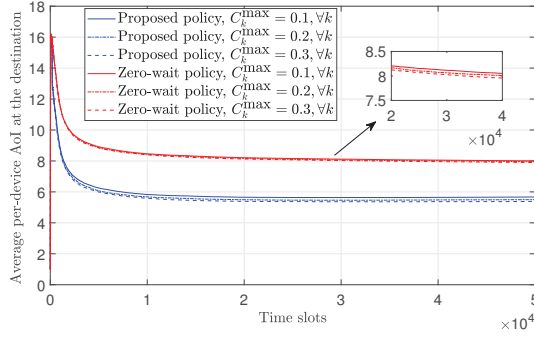


Fig. 7: Illustration of the convergence property. The number of IoT devices is  $K = 5$ .  $A_{l,k}^{\max} = A_{r,k}^{\max} = 100, \forall k$ .

Fig. 6 shows the average, per-device AoI at the destination and the average, per-device energy cost, resulting from the proposed semi-distributed policy and the zero-wait baseline policy. The simulation results are obtained by averaging over 100,000 time slots. In the simulations, for  $K = 10, 20, 30, 40, 50$ , it takes about 17000, 23000, 30000, 41000, 50000 time slots for the convergence of the proposed suboptimal policy, respectively. From Fig. 6, we can see that the proposed semi-distributed policy can achieve up-to 20% reduction of the average AoI at the destination over the zero-wait baseline policy, with similar energy costs. Thus the proposed policy can make better use of the limited energy at the IoT device. Moreover, for both policies, we observe that, as the number of the IoT devices increases, the average per-device AoI at the destination increases, while the average energy cost decreases. This is due to the fact that the transmission opportunities for each IoT device become lower with more IoT devices.

In Fig. 7, we show the evolution of the average per-device AoI at the destination, resulting from the proposed semi-distributed policy and the zero-wait baseline policy, under different  $C_k^{\max}$ . The convergence of the proposed semi-distributed learning algorithm can be clearly observed (after about 15,000 time slots). Moreover, with the increase of  $C_k^{\max}$ , the average per-device AoI at the destination for the two policies decreases. The performance gain of the proposed policy over the baseline policy can be as much as 33% when  $C_k^{\max} = 0.3$ .

## V. CONCLUSION

In this paper, we have studied the optimal sampling and updating processes that enable IoT devices to minimize the average AoI at the destination under an average energy constraint for each IoT device in a real-time IoT monitoring system. We have formulated this problem as an infinite horizon average cost CMDP and transformed it into an unconstrained MDP. For the single IoT device case, we have shown that the optimal sampling and updating policy is of threshold type, which reveals a fundamental tradeoff between the average AoI at the destination and the sampling and updating costs. Based on this optimality property, we have proposed a structure-aware algorithm to obtain the optimal policy for the CMDP. We have also studied the effects of the wireless channel fading

and shown that channels with large mean channel gain and less scattering can achieve better AoI performance. For the case of multiple IoT devices, we have shown that the optimal sampling and updating policy is a function of the Q-factors of the unconstrained MDP. To reduce the complexity in obtaining the optimal Q-factors, we have developed a semi-distributed low-complexity suboptimal policy by approximating the optimal Q-factors by a linear form of the per-device Q-factors. We have proposed an online algorithm for each device to estimate and learn its per-device Q-factor based on the locally observed AoI and channel states. We have shown the almost surely convergence of the proposed learning algorithm to the proposed suboptimal policy. Simulation results have shown that, for the single IoT device case, the optimal thresholds for sampling (updating) are non-decreasing with the sampling (updating) cost and the optimal action is threshold-based with respect to the channel state; and the proposed semi-distributed suboptimal policy for multiple IoT devices yields significant performance gain in terms of the average AoI compared to a zero-wait baseline policy. Future work will address key extensions such as providing an approximation analysis of the considered linear decomposition method and proposing grant-free uplink transmission protocols.

## APPENDIX

### A. Proof of Lemma 3

We prove Lemma 3 using the value iteration algorithm (VIA) and mathematical induction. First, we introduce the VIA [24, Chapter 4.3]. For notational convenience, we omit  $\lambda$  in the notation of  $V)A, h; \lambda$ . For each state  $)A, h[ \in \mathcal{A} * \mathcal{H}$ , let  $V_m)A, h[$  be the value function at iteration  $m$ . Define the state-action cost function at iteration  $m$  as:

$$J_{m+1})A, h, w[ \triangleq L)A, h, w; \lambda[ + \sum_{h' \in \mathcal{H}} p_{\mathcal{H}}(h|h') V_m)A, h', \lambda[ \quad (37)$$

where  $A^\infty$  is given by Lemma 2. Note that  $J_{m+1})A, h, w[$  is related to the right-hand side of the Bellman equation in (11). For each  $)A, h[$ , VIA calculates  $V_{m+1})A, h[$  according to

$$V_{m+1})A, h[ = \min_{w \in \mathcal{W}} J_{m+1})A, h, w[, \quad \forall l. \quad (38)$$

Under any initialization of  $V_0)A, h[$ , the generated sequence  $\{V_m)A, h[\}$  converges to  $V)A, h[$  [24, Proposition 4.3.1], i.e.,

$$\lim_{m \rightarrow \infty} V_m)A, h[ = V)A, h[, \quad \forall )A, h[ \in \mathcal{A} * \mathcal{H}, \quad (39)$$

where  $V)A, h[$  satisfies the Bellman equation in (11). Let  $\pi_m^*)A, h[$  denote the control that attains the minimum of the first term in (38) at iteration  $m$  for all  $)A, h[$ , i.e.,

$$\pi_m^*)A, h[ = \arg \min_{w \in \mathcal{W}} J_{m+1})A, h, w[, \quad \forall )A, h[ \in \mathcal{A} * \mathcal{H}. \quad (40)$$

We refer to  $\pi_m^* = )\pi_{s,m}^*, \pi_{u,m}^*[$  as the optimal policy for iteration  $m$ .

Now, consider two AoI states,  $A^1 = )A_l^1, A_r^1[$  and  $A^2 = )A_l^2, A_r^2[$ . To prove Lemma 3, we only need to show that for any  $A^1, A^2 \in \mathcal{A}$ , such that  $A_l^1 \geq A_l^2$  and  $A_r^1 \geq A_r^2$ ,

$$V_m)A^2, h[ \geq V_m)A^1, h[, \quad (41)$$

holds for all  $m = 0, 1, \infty$

First, we initialize  $V_0)A, h[ = 0$  for all  $A, h$ . Thus, (41) holds for  $m = 0$ . Assume that (41) holds for some  $m > 0$ . We will prove that (41) also holds for  $m + 1$ . By (38), we have

$$\begin{aligned} & V_{m+1})A^1, h[ = J_{m+1})A^1, h, \pi_m^\bullet)A^1, h[ \left\{ \right. \\ & \quad \left. \geq J_{m+1})A^1, h, \pi_m^\bullet)A^2, h[ \left\{ \right. \right. \\ & \quad \left. \stackrel{b[}{=} A_r^1 + \lambda C(\pi_m^\bullet)A^2, h[ + \underset{h \in \mathcal{H}}{p_{\mathcal{H}}(h)} V_m)A_l^1, A_r^1, h^\circ, \right. \end{aligned} \quad (42)$$

where  $a[$  is due to the optimality of  $\pi_m^\bullet)A^1, h[$  for  $)A^1, h[$  in the  $m$ -th iteration,  $b[$  directly follows from (37),  $A_l^1 = \min\{\pi_{s,m}^\bullet)A^2, h[ + )1 - \pi_{s,m}^\bullet)A^2, h[|)A_l^1 + 1[, \hat{A}_l|$  and  $A_r^1 = \min\{\pi_{u,m}^\bullet)A^2, h[|)A_l^1 + 1[ + )1 - \pi_{u,m}^\bullet)A^2, h[|)A_r^1 + 1[, \hat{A}_r|$ . By (37) and (38), we also have

$$\begin{aligned} & V_{m+1})A^2, h[ = J_{m+1})A^2, h, \pi_m^\bullet)A^2, h[ \left\{ \right. \\ & \quad \left. = A_r^2 + \lambda C(\pi_m^\bullet)A^2, h[ + \underset{h \in \mathcal{H}}{p_{\mathcal{H}}(h)} V_m)A_l^2, A_r^2, h^\circ, \right. \end{aligned} \quad (43)$$

where  $A_l^2 = \min\{\pi_{s,m}^\bullet)A^2, h[ + )1 - \pi_{s,m}^\bullet)A^2, h[|)A_l^2 + 1[, \hat{A}_l|$  and  $A_r^2 = \min\{\pi_{u,m}^\bullet)A^2, h[|)A_l^2 + 1[ + )1 - \pi_{u,m}^\bullet)A^2, h[|)A_r^2 + 1[, \hat{A}_r|$ .

It can be seen that  $A_l^2 \geq A_l^1$  and  $A_r^2 \geq A_r^1$  for all possible  $\pi_m^\bullet)A^2, h[ \in \mathcal{W}$ , which implies that  $V_m)A_l^2, A_r^2, h^\circ \geq V_m)A_l^1, A_r^1, h^\circ$  by induction. Thus, we have  $V_{m+1})A^2, h[ \geq V_{m+1})A^1, h[$ , i.e., (41) holds for  $m + 1$ . Therefore, by induction, we can show that (41) holds for any  $m$ . By taking limits on both sides of (41) and by (39), we complete the proof of Lemma 3.

### B. Proof of Lemma 4

First, we derive the general relation between  $\Delta J_{w, w^\infty})A^1, h[$  and  $\Delta J_{w, w^\infty})A^2, h[$  for any  $w, w^\infty \in \mathcal{W}$ ,  $h \in \mathcal{H}$ , and  $A^1, A^2 \in \mathcal{A}$ . Here,  $\lambda$  is also omitted in the notation of  $\Delta J_{w, w^\infty})A^1, h[; \lambda[$  for notational convenience. By (13), we have

$$\begin{aligned} & \Delta J_{w, w^\infty})A^1, h[ - \Delta J_{w, w^\infty})A^2, h[ \\ & = \underset{h \in \mathcal{H}}{L)A^1, h, w[ + \underset{h \in \mathcal{H}}{p_{\mathcal{H}}(h)} V)A^{1, w}, h^\circ - L)A^1, h, w^\circ} \\ & \quad + \underset{h \in \mathcal{H}}{p_{\mathcal{H}}(h)} V)A^{1, w^\infty}, h^\circ \left\{ \right. \\ & \quad \left. \underset{h \in \mathcal{H}}{L)A^2, h, w[ + \underset{h \in \mathcal{H}}{p_{\mathcal{H}}(h)} V)A^{2, w}, h^\circ - L)A^2, h, w^\circ} \right. \\ & \quad + \underset{h \in \mathcal{H}}{p_{\mathcal{H}}(h)} V)A^{2, w^\infty}, h^\circ \left\{ \right. \\ & = \underset{h \in \mathcal{H}}{p_{\mathcal{H}}(h)} V)A^{1, w}, h^\circ - V)A^{1, w^\infty}, h^\circ \\ & \quad \left. V)A^{2, w}, h^\circ + V)A^{2, w^\infty}, h^\circ \right\}, \end{aligned} \quad (44)$$

where  $A_l^{1, w} = \min\{s + )1 - s[|)A_l^1 + 1[, \hat{A}_l|$ ,  $A_r^{1, w} = \min\{u)A_l^1 + 1[ + )1 - u[|)A_r^1 + 1[, \hat{A}_r|$ ,  $A_l^{1, w^\infty} = \min\{s^\infty + )1 - s^\circ[|)A_l^1 + 1[, \hat{A}_l|$ ,  $A_r^{1, w^\infty} = \min\{u^\infty)A_l^1 + 1[ + )1 - u^\circ[|)A_r^1 + 1[, \hat{A}_r|$ ,  $A_l^{2, w} = \min\{s + )1 - s[|)A_l^2 + 1[, \hat{A}_l|$ ,  $A_r^{2, w} = \min\{u)A_l^2 + 1[ + )1 - u[|)A_r^2 + 1[, \hat{A}_r|$ ,  $A_l^{2, w^\infty} = \min\{s^\infty + )1 - s^\circ[|)A_l^2 + 1[, \hat{A}_l|$ , and  $A_r^{2, w^\infty} = \min\{u^\infty)A_l^2 + 1[ + )1 - u^\circ[|)A_r^2 + 1[, \hat{A}_r|$ .

Next, based on (44), we show that  $\Delta J_{w, w^\infty})A^1, h[$  is non-decreasing with  $A_l$  for  $w^\infty = )1, 0[$ . Consider  $w = )0, 0[$ ,  $w^\infty = )1, 0[$ , and  $A^1$  and  $A^2$  satisfying  $A_l^1 \geq A_l^2$  and  $A_r^1 = A_r^2$ . We can see that,  $A_l^{1, w} \geq A_l^{2, w}$ ,  $A_r^{1, w} = A_r^{2, w}$ ,  $A_l^{1, w^\infty} = A_l^{2, w^\infty}$ , and  $A_r^{1, w^\infty} = A_r^{2, w^\infty}$ . Thus, we have  $V)A^{1, w^\infty}, h^\circ = V)A^{2, w^\infty}, h^\circ$  and by Lemma 3, we have  $V)A^{1, w}, h^\circ \geq V)A^{2, w}, h^\circ$ . Therefore, by (44), we have  $\Delta J_{w, w^\infty})A^1, h[ - \Delta J_{w, w^\infty})A^2, h[ \geq 0$ , which completes the proof. Similarly, we can also prove the remaining properties of  $\Delta J_{w, w^\infty})A^1, h[$  in Lemma 4.

### C. Proof of Theorem 1

We first prove Property A) of Theorem 1. Consider action  $w = )0, 0[$ , action  $w^\infty = )1, 0[$ , channel state  $h$ , AoI state  $A$  where  $A_l = \phi_{0,0}^+)A_r, h[$ . ( $\lambda$  is omitted here.) Note that, we only need to consider  $\phi_{0,0}^+)A_r, h[ > \epsilon$ . According to the definition of  $\phi_w^+)A_r, h[$  in (16), we can see that  $\Delta J_{w, w^\infty})A, h[ \geq 0$ , i.e.,  $w$  dominates  $w^\infty$  for state  $A, h[$ . Now, consider another AoI state  $A^\infty$  where  $A_l^\infty \geq A_l$  and  $A_r^\infty = A_r$ . By Lemma 4, we obtain that

$$\Delta J_{w, w^\infty})A^\infty, h[ \geq \Delta J_{w, w^\infty})A, h[ \geq 0, \quad (45)$$

i.e.,  $w = )0, 0[$  also dominates  $w^\infty = )1, 0[$  for state  $A^\infty, h[$ . Now, we consider  $w^\infty = )0, 1[$  or  $)1, 1[$ , channel state  $h$ , AoI state  $A$  where  $A_l = \psi_{0,0}^+)A_l, h[$ , AoI state  $A^\infty$  where  $A_l^\infty = A_l$  and  $A_r^\infty \geq A_r$ . According to the definition of  $\psi_w^+)A_l, h[$  in (18) and Lemma 4, we can prove that (45) still holds, i.e.,  $w = )0, 0[$  also dominates  $w^\infty = )1, 0[$  or  $)1, 1[$  for state  $A^\infty, h[$ . By the definition of  $A_0)h[$ , we can see that if  $A \in \mathcal{A}_0)h[$ , then  $w = )0, 0[$  dominates all other actions, i.e.,  $\pi^\bullet)A, h[ = )0, 0[$ . We complete the proof of Property A).

Next, we prove Property B) of Theorem 1. Consider action  $w = )0, 1[$ , channel state  $h$ , AoI state  $A$  where  $A_r = \psi_{0,1}^+)A_l, h[$ . We only need consider that  $\psi_{0,1}^+)A_l, h[ < +\epsilon$ . By the definition of  $\psi_{0,1}^+)A_l, h[$  in (19), we have  $\Delta J_{w, w^\infty})A, h[ \geq 0$  for all  $w^\infty \neq w$ , i.e.,  $\pi^\bullet)A, h[ = )0, 1[$ . Now consider another AoI state  $A^\infty$  where  $A_l^\infty = A_l$  and  $A_r^\infty \geq A_r$ . By Lemma 4, we can see that  $\Delta J_{w, w^\infty})A^\infty, h[ \geq \Delta J_{w, w^\infty})A, h[ \geq 0$  holds for all  $w^\infty \neq w$ , i.e.,  $\pi^\bullet)A^\infty, h[ = )0, 1[$ . We complete the proof of Property B). Following the proof of Property B), we can also prove Properties C) and D). This completes the proof of Theorem 1.

### D. Proof of Theorem 2

To prove Theorem 2, we first prove that for any channels  $I$  and  $J$  such that  $h^I$  first-order stochastically dominates  $h^J$ ,

$$V^I)A, h; \lambda[ \geq V^J)A, h; \lambda[, \quad (46)$$

holds for all  $)A, h[$ , where  $V^I)A, h; \lambda[$  and  $V^J)A, h; \lambda[$  are the value functions under channels  $I$  and  $J$ , respectively. We prove (46) through mathematical induction and the VIA in Appendix A. Similar to Appendix A, we introduce  $V_m^I)A, h; \lambda[, V_m^J)A, h; \lambda[, J_m^I)A, h, w[, J_m^J)A, h, w[, \pi_m^{I\bullet} = )\pi_{s,m}^{I\bullet}, \pi_{u,m}^{I\bullet}[$ , and  $\pi_m^{J\bullet} = )\pi_{s,m}^{J\bullet}, \pi_{u,m}^{J\bullet}[$  for channels  $I$  and  $J$ . Since  $C_u)h[$  is non-increasing with  $h$ , it can be easily shown that  $V_m^I)A, h; \lambda[$  and  $V_m^J)A, h; \lambda[$  are non-increasing with  $h$ , by

using induction and the VIA. To show (46), by (39), we only need to show that

$$V_m^I \mathbf{A}, h; \lambda \geq V_m^J \mathbf{A}, h; \lambda, \quad (47)$$

holds for  $m = 0, 1, \dots$ . We initialize  $V_0^I \mathbf{A}, h; \lambda = V_0^J \mathbf{A}, h; \lambda = 0$ , for all  $\mathbf{A}, h$ , i.e., (47) holds for  $m = 0$ . Assume that (47) holds for some  $m > 0$ . We will show that (47) also holds for  $m + 1$ . By (37) and (38), we have

$$\begin{aligned} & V_{m+1}^I \mathbf{A}, h; \lambda = J_{m+1}^I \mathbf{A}, h; \pi_m^{\bullet} \mathbf{A}, h; \lambda \Big\{ \\ & \quad \Big\{ \geq J_{m+1}^I \mathbf{A}, h; \pi_m^{\bullet} \mathbf{A}, h; \lambda \Big\} \\ & = L \mathbf{A}, h; \pi_m^{\bullet} \mathbf{A}, h; \lambda + \int_{\mathcal{H}} p_{\mathcal{H}}^I h^{\circ} V_m^I \mathbf{A}^{\infty}; h^{\circ}; \lambda \\ & \quad \Big\{ \geq L \mathbf{A}, h; \pi_m^{\bullet} \mathbf{A}, h; \lambda + \int_{\mathcal{H}} p_{\mathcal{H}}^J h^{\circ} V_m^I \mathbf{A}^{\infty}; h^{\circ}; \lambda \\ & \quad \Big\{ \geq L \mathbf{A}, h; \pi_m^{\bullet} \mathbf{A}, h; \lambda + \int_{\mathcal{H}} p_{\mathcal{H}}^J h^{\circ} V_m^J \mathbf{A}^{\infty}; h^{\circ}; \lambda \\ & = V_{m+1}^J \mathbf{A}, h; \lambda, \end{aligned} \quad (48)$$

where  $\Big\{ \geq$  is due to the optimality of  $\pi_m^{\bullet} \mathbf{A}, h$  for  $\mathbf{A}, h$  under channel  $I$  in the  $m$ -th iteration,  $\Big\{ \geq$  is due to that  $V_m^I$  first-order stochastically dominates  $V_m^J$  and  $V_m^I \mathbf{A}, h; \lambda$  is non-increasing with  $h$ ,  $\Big\{ \geq$  follows from the induction hypothesis  $V_m^I \mathbf{A}, h; \lambda \geq V_m^J \mathbf{A}, h; \lambda$ ,  $A_l^{\infty} = \min \{ \pi_{s,m}^{\bullet} \mathbf{A}, h; \lambda + 1, \pi_{s,m}^{\bullet} \mathbf{A}, h; \lambda \}$ , and  $A_r^{\infty} = \min \{ \pi_{u,m}^{\bullet} \mathbf{A}, h; \lambda + 1, \pi_{u,m}^{\bullet} \mathbf{A}, h; \lambda \}$ . Thus, we prove (47) holds for  $m + 1$ . Then, by induction and (39), we can show that (46) holds. Based on (46), Lemma 2 and Propositions 4.3.1 in [24], we have  $\bar{L}^{\bullet} \lambda = \theta_{\lambda}^{\bullet} \geq \theta_{\lambda}^J = \bar{L}^J \lambda$ . Finally, by Lemma 1, we can see that  $A_r^{\bullet} = \max_{\lambda \geq 0} \bar{L}^J \lambda$ .  $\lambda C^{\max} \geq \max_{\lambda \geq 0} \bar{L}^J \lambda$ .  $\lambda C^{\max} = \bar{A}_r^{\bullet}$ , which completes the proof of Theorem 2.

#### E. Proof of Lemma 5

For a given  $\lambda$ , by Propositions 4.2.1, 4.2.3, and 4.2.5 in [24] (see Propositions 4.2.3 and 4.2.5 in Appendix H), the optimal average Lagrange cost for the unconstrained MDP in (28) is the same for all initial states and the optimal policy can be obtained by solving the following Bellman equation with respect to  $\theta_{\lambda}, \forall \mathbf{A}, h; \lambda$  [

$$\begin{aligned} \theta_{\lambda} + V \mathbf{A}, h; \lambda &= \min_{\mathbf{w} \in \mathcal{W}} \Big\{ L \mathbf{A}, h, \mathbf{w}; \lambda \\ &+ \int_{\mathcal{H}} p_{\mathcal{H}} h^{\circ} V \mathbf{A}^{\infty}; h^{\circ}; \lambda, \forall \mathbf{A}, h \in \mathcal{A} * \mathcal{H}, \end{aligned} \quad (49)$$

where  $V \mathbf{A}, h; \lambda$  is the value function. Since  $\pi_{\lambda}^{\bullet} \mathbf{A}, h = \pi_{\lambda,s}^{\bullet} \mathbf{A}, h, \pi_{\lambda,u}^{\bullet} \mathbf{A}, h$ , we introduce the Q-factor of state  $\mathbf{A}, h$  under updating action  $u$  as:

$$Q \mathbf{A}, h, u; \lambda \triangleq \min_{s \in \mathcal{S}} \Big\{ L \mathbf{A}, h, u; \lambda + \int_{\mathcal{H}} p_{\mathcal{H}} h^{\circ} V \mathbf{A}^{\infty}; h^{\circ}; \lambda \Big\} \quad (50)$$

Thus, we have  $V \mathbf{A}, h; \lambda = \min_{u \in \mathcal{U}} Q \mathbf{A}, h, u; \lambda$  for all  $\mathbf{A}, h$  and  $\theta_{\lambda}, \forall \mathbf{A}, h; \lambda$  [ satisfies the Bellman equation in (29). We complete the proof.

#### F. Proof of Lemma 6

Under a unichain policy defined in Definition 3, the induced random process  $\{ \mathbf{A} \}_t, \{ h \}_t$  is a controlled Markov chain with a single recurrent class and possibly some transient states [24]. According to the explanation for the condition of Proposition 4.3.2 in [24], the condition of Lemma 2 in [34] is satisfied for our problem. Then, by following the proofs of Lemma 2 in [34] and Proposition 4.3.2 in [24], we can prove the Lemma 6. The detailed proof is omitted due to page limitations.

#### G. Proof of Lemma 7

Due to the separation of the two timescales of the updates in (34) and (35), the update of the Q-factors can be regarded as converged to  $Q_k^{\epsilon} \lambda^t$  under  $\lambda^t$  [39]. Then, by the theory of stochastic approximation [34], [39], [42], the iterations of the update of the Lagrange multiplier in (35) can be described by the following Ordinary Differential Equation (ODE):

$$\lambda^t = \mathbb{E}^{\pi_{\lambda^t}^{\bullet}} [C_1] \hat{\mathbf{w}}_1(t) - [C_1^{\max}, \dots, C_K] \hat{\mathbf{w}}_K(t) - [C_K^{\max}], \quad (51)$$

where  $\pi_{\lambda^t}^{\bullet}$  is the converged control policy in Algorithm 2 under  $\lambda^t$  and the expectation is taken with respect to the measure induced by the policy  $\pi_{\lambda^t}^{\bullet}$ . Denote  $\bar{L} \lambda^t = \mathbb{E}^{\pi_{\lambda^t}^{\bullet}} [\prod_{k=1}^K A_{r,k} + \lambda_k C_k] \hat{\mathbf{w}}_k$ . Since the sampling and updating actions are discrete, we have  $\pi_{\lambda^t}^{\bullet} = \pi_{\lambda^t + \delta_{\lambda}}^{\bullet}$ . By chain rule, it can be seen that  $\frac{\partial \bar{L} \lambda^t}{\partial \lambda_k^t} = \mathbb{E}^{\pi_{\lambda^t}^{\bullet}} [C_k] \hat{\mathbf{w}}_k(t) - [C_k^{\max}]$ . Thus, the ODE in (51) can be expressed as  $\lambda^t = \nabla \bar{L} \lambda^t$ . Therefore, the ODE in (51) will converge to  $\arg \max_{\lambda \in \mathbb{R}^+} \bar{L} \lambda$ , which corresponds to  $\nabla \bar{L} \lambda = 0$ . In other words, the policy under  $Q^{\epsilon}, \lambda^{\epsilon}$  satisfies the constraint in (26b). This completes the proof.

#### H. Some preliminaries on MDP

Proposition 4.2.3 in [24]: Let the weak accessibility (WA) condition hold. Then the optimal average cost is the same for all initial states.

Proposition 4.2.5 in [24]: If all stationary policies are unichain, the WA condition holds.

Theorem 8.6.6 in [26]: Suppose all stationary policies are unichain, and the set of states and actions are finite, then policy iteration converges in a finite number of iterations to the optimal policy satisfying the Bellman equation.

#### REFERENCES

- [1] B. Zhou and W. Saad, "Optimal sampling and updating for minimizing age of information in the Internet of Things," in *Proc. of IEEE Global Communications Conference (GLOBECOM)*, Abu Dhabi, UAE, Dec. 2018.
- [2] L. Atzori, A. Iera, and G. Morabito, "The Internet of Things: A survey," *Comput. Networks*, vol. 54, no. 15, pp. 2787–2805, 2010.
- [3] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Unmanned aerial vehicle with underlaid device-to-device communications: Performance and tradeoffs," *IEEE Trans. Wireless Commun.*, vol. 15, no. 6, pp. 3949–3963, June 2016.
- [4] S. Kaul, R. Yates, and M. Gruteser, "Real-time status: How often should one update?" in *Proc. of IEEE International Conference on Computer Communications (INFOCOM)*, Orlando, FL, USA, March 2012, pp. 2731–2735.

- [5] S. Kaul, R. D. Yates, and M. Gruteser, "Status updates through queues," in *Proc. of 46th Annual Conference on Information Sciences and Systems (CISS)*, Princeton, NJ, USA, March 2012, pp. 1–6.
- [6] Y.-P. Hsu, E. Modiano, and L. Duan, "Scheduling algorithms for minimizing age of information in wireless broadcast networks with random arrivals," *arXiv preprint arXiv:1712.07419*, 2017.
- [7] I. Kadota, A. Sinha, E. Uysal-Biyikoglu, R. Singh, and E. Modiano, "Scheduling policies for minimizing age of information in broadcast wireless networks," *arXiv preprint arXiv:1801.01803*, 2018.
- [8] Y.-P. Hsu, "Age of information: Whittle index for scheduling stochastic arrivals," in *Proc. of IEEE International Symposium on Information Theory (ISIT)*, Colorado, USA, June 2018.
- [9] Z. Jiang, B. Krishnamachari, S. Zhou, and Z. Niu, "Can decentralized status update achieve universally near-optimal age-of-information in wireless multiaccess channels?" in *Proc. of IEEE The International Teletraffic Congress (ITC)*, Vienna, Austria, Sep. 2018.
- [10] Z. Jiang, B. Krishnamachari, X. Zheng, S. Zhou, and Z. Niu, "Timely status update in massive IoT systems: Decentralized scheduling for wireless uplinks," *arXiv preprint arXiv:1801.03975*, 2018.
- [11] Y. Sun, E. Uysal-Biyikoglu, R. D. Yates, C. E. Koksal, and N. B. Shroff, "Update or wait: How to keep your data fresh," *IEEE Trans. Inf. Theory*, vol. 63, no. 11, pp. 7492–7508, Nov 2017.
- [12] A. M. Bedewy, Y. Sun, S. Kompella, and N. B. Shroff, "Age-optimal sampling and transmission scheduling in multi-source systems," *arXiv preprint arXiv:1812.09463*, 2018.
- [13] X. Wu, J. Yang, and J. Wu, "Optimal status update for age of information minimization with an energy harvesting source," *IEEE Trans. Green Commun. and New.*, vol. 2, no. 1, pp. 193–204, March 2018.
- [14] B. T. Bacinoglu and E. Uysal-Biyikoglu, "Scheduling status updates to minimize age of information with an energy harvesting sensor," *arXiv preprint arXiv:1701.08354*, 2017.
- [15] S. Feng and J. Yang, "Minimizing age of information for an energy harvesting source with updating failures," in *Proc. of IEEE International Symposium on Information Theory (ISIT)*, Colorado, USA, June 2018, pp. 2431–2435.
- [16] E. T. Ceran, D. Gunduz, and A. Gyorgy, "Average age of information with hybrid ARQ under a resource constraint," in *Proc. of IEEE Wireless Communications and Networking Conference (WCNC)*, Barcelona, Spain, April 2018.
- [17] Y. Sun, Y. Polyanskiy, and E. Uysal-Biyikoglu, "Remote estimation of the wiener process over a channel with random delay," *arXiv preprint arXiv:1701.06734*, 2017.
- [18] Nest Cam IQ indoor security camera, <https://nest.com/cameras/nest-cam-iq-indoor/overview/>.
- [19] M. Chen, U. Challita, W. Saad, C. Yin, and M. Debbah, "Machine learning for wireless networks with artificial intelligence: A tutorial on neural networks," *arXiv preprint arXiv:1710.02913*, 2017.
- [20] S. Teerapittayanon, B. McDanel, and H. Kung, "Distributed deep neural networks over the cloud, the edge and end devices," in *Proc. of IEEE International Conference on Distributed Computing Systems (ICDCS)*, Atlanta, GA, USA, June 2017, pp. 328–339.
- [21] D. V. Djonin and V. Krishnamurthy, "MIMO transmission control in fading channels—a constrained Markov decision process formulation with monotone randomized policies," *IEEE Trans. Signal Process.*, vol. 55, no. 10, pp. 5069–5083, 2007.
- [22] E. Altman, *Constrained Markov Decision Processes*. London, U.K.: Chapman & Hall, 1999.
- [23] F. J. Beutler and K. W. Ross, "Optimal policies for controlled Markov chains with a constraint," *Journal of Mathematical Analysis and Applications*, vol. 112, no. 1, pp. 236 – 252, 1985.
- [24] D. P. Bertsekas, *Dynamic programming and optimal control, 3rd edition, volume II*. Belmont, MA: Athena Scientific, 2011.
- [25] G. Koole, "Monotonicity in Markov reward and decision chains: Theory and applications," *Foundations and Trends in Stochastic Systems*, vol. 1, no. 1, pp. 1–76, 2006.
- [26] M. L. Puterman, *Markov decision processes: discrete stochastic dynamic programming*. New York, NY, USA: Wiley, 2009, vol. 414.
- [27] M. L. Littman, T. L. Dean, and L. P. Kaelbling, "On the complexity of solving markov decision problems," in *Proceedings of the Eleventh conference on Uncertainty in artificial intelligence*, 1995, pp. 394–402.
- [28] M. H. Ngo and V. Krishnamurthy, "Monotonicity of constrained optimal transmission policies in correlated fading channels with ARQ," *IEEE Trans. Signal Process.*, vol. 58, no. 1, pp. 438–451, Jan 2010.
- [29] H. Robbins and S. Monro, "A stochastic approximation method," *Ann. Math. Statist.*, vol. 22, no. 3, pp. 400–407, 09 1951.
- [30] J. E. Smith and K. F. McCardle, "Structural properties of stochastic dynamic programs," *Operations Research*, vol. 50, no. 5, pp. 796–809, 2002.
- [31] A. D. Zayas and P. Merino, "The 3GPP NB-IoT system architecture for the Internet of Things," in *Proc. of IEEE International Conference on Communications Workshops (ICC Workshops)*, Paris, France, May 2017, pp. 277–282.
- [32] K. W. Ross, "Randomized and past-dependent policies for Markov decision processes with multiple constraints," *Operations Research*, vol. 37, no. 3, pp. 474–477, 1989.
- [33] J. Gittins, K. Glazebrook, and R. Weber, *Multi-armed bandit allocation indices*. John Wiley & Sons, 2011.
- [34] Y. Cui and V. K. N. Lau, "Distributive stochastic learning for delay-optimal OFDMA power and subband allocation," *IEEE Trans. Signal Process.*, vol. 58, no. 9, pp. 4848–4858, Sept 2010.
- [35] J. Abounadi, D. Bertsekas, and V. S. Borkar, "Learning algorithms for markov decision processes with average cost," *SIAM J. Control Optim.*, vol. 40, no. 3, pp. 681–698, 2001.
- [36] Y. Wang, X. Lin, A. Adhikary, A. Grovlen, Y. Sui, Y. Blankenship, J. Bergman, and H. S. Razaghi, "A primer on 3GPP narrowband Internet of Things," *IEEE Commun. Mag.*, vol. 55, no. 3, pp. 117–123, March 2017.
- [37] M. Hasan, E. Hossain, and D. Niyato, "Random access for machine-to-machine communication in LTE-advanced networks: Issues and approaches," *IEEE Commun Mag.*, vol. 51, no. 6, pp. 86–93, 2013.
- [38] D. V. Djonin and V. Krishnamurthy, "Q-learning algorithms for constrained markov decision processes with randomized monotone policies: Application to mimo transmission control," *IEEE Trans. Signal Processing*, vol. 55, no. 5-2, pp. 2170–2181, 2007.
- [39] V. Borkar, *Stochastic Approximation: A Dynamical Systems Viewpoint*. Hindustan Book Agency, 2008.
- [40] K. T. Phan, T. Le-Ngoc, M. van der Schaar, and F. Fu, "Optimal scheduling over time-varying channels with traffic admission control: Structural results and online learning algorithms," *IEEE Trans. Wireless Commun.*, vol. 12, no. 9, pp. 4434–4444, September 2013.
- [41] C. You, K. Huang, H. Chae, and B. H. Kim, "Energy-efficient resource allocation for mobile-edge computation offloading," *IEEE Trans. Wireless Commun.*, vol. 16, no. 3, pp. 1397–1411, March 2017.
- [42] V. S. Borkar, "Stochastic approximation with two time scales," *Systems & Control Letters*, vol. 29, no. 5, pp. 291 – 294, 1997.