

Multi-Source AoI-Constrained Resource Minimization under HARQ: Heterogeneous Sampling Processes

Saeid Sadeghi Vilni*, Mohammad Moltafet*, Markus Leinonen*, and
Marian Codreanu[†]

Abstract

We consider a multi-source hybrid automatic repeat request (HARQ) based system, where a transmitter sends status update packets of *random arrival* (i.e., uncontrollable sampling) and *generate-at-will* (i.e., controllable sampling) sources to a destination through an error-prone channel. We develop transmission scheduling policies to minimize the average number of transmissions subject to an average age of information (AoI) constraint. First, we consider known environment (i.e., known system statistics) and develop a near-optimal deterministic transmission policy and a low-complexity dynamic transmission (LC-DT) policy. The former policy is derived by casting the main problem into a constrained Markov decision process (CMDP) problem, which is then solved using the Lagrangian relaxation, relative value iteration algorithm, and bisection. The LC-DT policy is developed via the drift-plus-penalty (DPP) method by transforming the main problem into a sequence of per-slot problems. Finally, we consider unknown environment and devise a learning-based transmission policy by relaxing the CMDP problem into an MDP problem using the DPP method and then adopting the deep Q-learning algorithm. Numerical results show that the proposed policies achieve near-optimal performance and illustrate the benefits of HARQ in status updating.

Index Terms: AoI, multi-source status update, CMDP, Lagrangian, Lyapunov, machine learning.

I. INTRODUCTION

There is a growing demand for Internet-of-Things (IoT) and cyber-physical systems such as autonomous vehicles, wireless industrial automation, and health monitoring that rely heavily on real-time (fresh) status updates. In these systems, a source (containing a sensor) monitors a physical phenomenon such as temperature, pressure, or motion and sends status updates to a

*Centre for Wireless Communications–Radio Technologies, University of Oulu, 90014 Oulu, Finland (e-mail: firstname.lastname@oulu.fi). [†]Department of Science and Technology, Linköping University, Sweden (e-mail: marian.codreanu@liu.se).

destination (e.g., a remote controller) for decision-making [1], [2]. The Age of Information (AoI) [1]–[3] is a metric used to evaluate the freshness of information in the status update systems. AoI is defined as the difference between the current time and the generation time of the last received packet at a destination [1]–[3]. Each status update packet contains a timestamp representing the time when the sample was generated and the measured value of the monitored process. At time instant t , denoting the timestamp of the last received status update packet by U_t , the AoI, δ_t , is defined as $\delta_t = t - U_t$ [1]–[5].

The reliability of data transmissions under an unreliable communication channel can be enhanced via retransmission protocols [6]. Automatic repeat request (ARQ) protocols are standard error control methods, where after each transmission, the transmitter receives a feedback about the reception status of the packet as acknowledgement/negative-acknowledgement (ACK/NACK) [6]. The transmitter keeps retransmitting each packet until it receives an ACK or reaches the maximum allowed number of retransmissions. The ARQ protocols use only the last received version of a packet for decoding, whereas the hybrid ARQ (HARQ) protocols use all received versions, thus increasing the probability of successfully decoding the packet [6], [7].

In this paper, we consider a multi-source HARQ-based status update system, where the sources are connected to a transmitter that sends status update packets to a receiver over an unreliable wireless channel (see Fig. 1). We assume a slotted communication, in which the transmitter can send at most one packet per slot. The sources, which monitor some time-varying random processes, are classified into two categories based on their sampling processes: 1) *random arrival* sources (i.e., uncontrollable sampling) which generate status update packets according to a Bernoulli process, and 2) *generate-at-will* sources (i.e., controllable sampling) which can be commanded to generate a status update packets at any slot. Our considered system may represent a multi-source node such as a multi-sensor IoT device that is equipped with a single transmitter to communicate all the different sensed data to a remote location through a wireless channel. In such scenario, the transmitter has control over the sampling of some sources, e.g., on-demand requesting of samples from a temperature or moisture sensor. On the contrary, the sampling of some sources depends on other grounds, such as energy to generate a sample (e.g., the source needs to harvest energy) or time to generate a sample (e.g., the source needs to scan an area which takes a random amount of time), leading to generating packets at random times. We further consider that each random arrival source has a buffer to retain the last (randomly) generated packet. Furthermore, in order to benefit from the HARQ, the transmitter has a memory to store the last transmitted but not successfully decoded packet of each source as this packet

has a higher chance of being decoded than a new packet.

Apart from freshness requirements, the radio resources (e.g., power and channel utilization) also play an essential role in the operation of status update systems [8]. Hence, we investigate the problem of minimizing the average number of transmissions subject to the average AoI constraint. The solution of the problem determines the transmission status at each time slot, i.e., transmit a fresh packet from a source, retransmit the previously transmitted but not successfully decoded packet from a source, or stay idle.

A scenario where the controller knows the probability of possible outcomes after making a decision in a system is called a known environment [9]. In our considered system, the known environment corresponds to the case where the transmitter knows the packet arrival rate of the random arrival sources and the probability of successful decoding after each transmission attempt. Since, in some cases, the known environment is not accessible, we investigate the problem in both the known environment and the unknown environment. We propose three solutions to the problem, namely, a (stationary) deterministic transmission policy and a low-complexity dynamic transmission (LC-DT) policy for the known environment and a learning-based transmission policy for the unknown environment.

To obtain the deterministic transmission policy in the known environment, we cast the main problem as a constrained Markov decision process (CMDP) problem. Then, we transform the CMDP problem into an MDP problem via the Lagrangian relaxation. In general, an optimal policy for the CMDP problem is a randomized mixture of two deterministic policies, where one deterministic policy is feasible (satisfies the constraint) and the another policy is infeasible [10]; see also recent applications [11]–[13]. However, since obtaining such randomized policy is often computationally intractable, we propose a near-optimal practical deterministic transmission policy (feasible deterministic policy) and a lower-bound policy (infeasible deterministic policy) for benchmarking purposes using relative value iteration algorithm (RVIA) and the bisection algorithm. Since the number of states to explore in RVIA increases exponentially in the number of sources and RVIA is run at each bisection iteration, obtaining the deterministic transmission policy is inefficient computationally. Therefore, we propose the LC-DT policy by using the drift-plus-penalty (DPP) method [14]. According to the DPP method, the average AoI constraint is transformed into a queue stability constraint, and subsequently, the time average main problem is transformed into an optimization problem that is to be solved at each time slot. To obtain the learning-based transmission policy, we use the DPP method to transform the CMDP problem into an MDP problem, in which we minimize the time average DPP function. Then, we develop

the policy by solving the MDP problem with a deep Q-learning (DQL) algorithm [15].

In the numerical results, we analyze the effectiveness of the proposed policies. We compare the effectiveness and complexity of the policies and study the effect of employing HARQ.

The main contributions of the paper are summarized as follows:

- We consider a multi-source HARQ-based status update system that consists of random arrival and generate-at-will sources. We minimize the average number of transmissions under the average AoI constraint in the known and unknown environment.
- For the known environment, we develop a deterministic transmission policy using the Lagrangian relaxation, RVIA, and the bisection. Moreover, we propose a low-complexity dynamic transmission policy using the DPP method.
- For the unknown environment, we develop a learning-based transmission policy by using the DPP method to cast the main problem as an MDP problem, which is then solved by applying DQL.
- The numerical results demonstrate the near-optimal performance of the proposed transmission policies compared to the lower-bound policy, and a significant improvement respect to a baseline policy. The results corroborate that HARQ improves the performance of the system and illustrate that the learning-based transmission policy performs close to the policies developed for the known environment.

A. Related Work

AoI characterization has extensively been studied from the perspective of queueing theory; see, e.g., [16]–[20] and the references therein. One of the earliest studies to analyze AoI under an HARQ protocol is [18], where the authors derived the closed-form expression of the average AoI for an HARQ-based M/G/1/1 queueing system. Considering the queueing system and the derived AoI result from [18], the work [19] studied the age-optimal redundancy allocation problem under a constraint on the decoding error probability for both chase combining and incremental redundancy HARQ protocols. An M/M/1 queueing system with network-code-HARQ protocol is considered in [20], where the closed-form expression of AoI is derived.

Besides the analysis, the AoI has been studied in the retransmission-based status update systems from the perspective of sampling and transmission policies [12], [13], [21]–[26]. In [21], the authors considered a multi-source and generate-at-will-based status update system and minimized the average AoI by proposing a source selection policy under three pre-defined transmission policies. In [22], the authors derived the closed-form expression of the average AoI

in an HARQ-based status update system in which two energy harvesting sources send the same information for providing diversity at the destination. The work [23] considered a multi-source status update system in which the transmitter harvests energy and uses a greedy retransmission policy. They minimized the average AoI by determining a set of transmission times and choosing a source to send status update packet. An HARQ-based non-orthogonal multiple access system with two users are considered in [24], where the average AoI is minimized by determining the transmit power and transmission status, i.e., transmitting a new packet or retransmitting the previously transmitted but not successfully decoded packet, at each slot. In [25], the authors investigated the average AoI minimization problem in a status update system with a pre-defined retransmission policy to find the times for updating the destination. The work [26] studied the average AoI minimization problem in an HARQ-based status update system. They calculated the probability of decoding failure through an erasure channel and developed a threshold-based transmission policy that decides between the transmission of a new packet and the retransmission of the previously transmitted one.

The most related works to our paper are [12], [13]. The work [13] considered a similar HARQ-based status update system to ours, yet with the following differences. The authors in [13] considered a single generate-at-will source, while we consider both random arrival and generate-at-will sources as a multi-source system. Considering the random arrival sources makes the system more complicated, as the transmitter does not know the availability of the fresh packets at the subsequent slots. We study the problem of minimizing the average number of transmissions subject to the average AoI constraint, while they studied the average AoI minimization problem subject to the average number of transmissions constraint. Similarly as in [13], we use the CMDP approach along with the Lagrangian relaxation to solve the problem in the known environment; however, we also propose the low-complexity Lyapunov-based dynamic transmission policy. Furthermore, for the unknown environment, they proposed a learning-based transmission policy by the Lagrangian relaxation which involves running the learning procedure for several Lagrangian multipliers. Differently, our learning-based transmission policy utilizes the Lyapunov optimization theory, and thus, the learning procedure needs to be run only once. In [12], which is an extension of [13], the authors considered an HARQ-based status update system that contains one generate-at-will source and several users (destinations), in which at most one user is served at each slot. They constructed a CMDP problem for minimizing the weighted average AoI subject to the average number of transmissions constraint. They solved the CMDP problem with the Lagrangian relaxation for the known environment; for the unknown

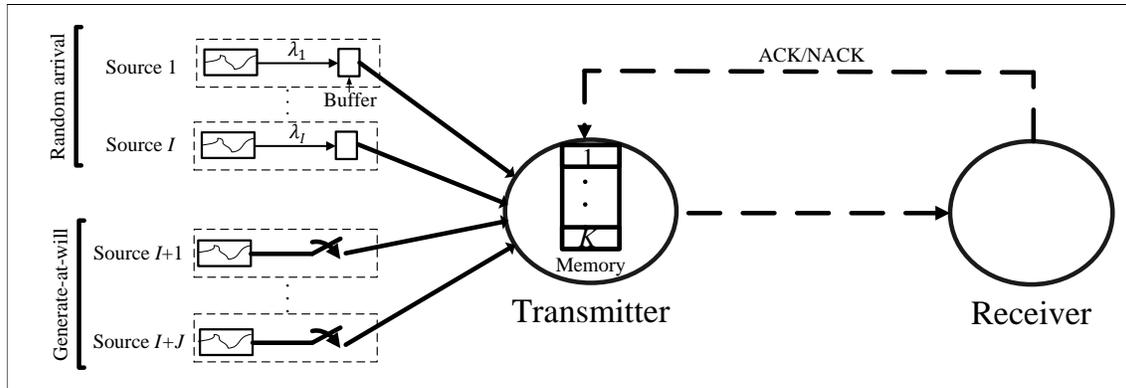


Figure 1: The considered HARQ-based multi-source status update system with two groups of sources: 1) random arrival sources with buffers which receive fresh packets with probability λ_k , 2) generate-at-will sources. The previously transmitted but not successfully decoded packets of each source are stored in the transmitter's memory. After each transmission attempt, the receiver sends a feedback signal as ACK (successful decoding) or NACK (unsuccessful decoding) to the transmitter.

environment, they proposed different learning-based transmission policies by the Lagrangian relaxation.

B. Organization

The rest of this paper is organized as follows. The system model and problem formulation are presented in Section II. The CMDP formulation and its solution are presented in Section III. In Section IV, we present the LC-DT and the learning-based transmission policies. Numerical results are presented in Section V. Finally, concluding remarks are made in Section VI.

II. SYSTEM MODEL AND PROBLEM FORMULATION

A. System Model

We consider a multi-source status update system that consists of K sources, one transmitter, and one receiver, as depicted in Fig. 1. The receiver is interested in timely information about different random processes monitored by the K sources. The transmitter sends status update packets¹ to the receiver through an error-prone wireless channel with the aid of an HARQ protocol. The system operates in discrete time with unit time slots $t \in \{1, 2, \dots\}$.

The K sources are divided into two classes based on their sampling processes: 1) a set \mathcal{I} of I *random arrival* sources whose sampling processes are uncontrollable and 2) a set \mathcal{J} of J *generate-at-will* sources, where the transmitter can sample the process at any time. Each source $k \in \mathcal{I}$ generates status update packets randomly and independently at the beginning of slots

¹Each status update packet contains a timestamp representing the time when the sample was generated and the measured value of the monitored process.

according to a Bernoulli random process with parameter λ_k . We denote the set of all sources by $\mathcal{K} = \mathcal{I} \cup \mathcal{J} = \{1, \dots, K\}$, where $K = I + J$.

Each random arrival source has a *buffer* of size one to store the last arrived packet. As long as a new packet does not arrive, the buffer keeps the last arrived packet. The transmitter has a *memory* of size K packets to store the previously transmitted but not successfully decoded packets of each source. Note that, after a number of unsuccessful transmission attempts of a packet from a source, the transmitter may decide to transmit a packet from the other sources. In this case, the transmitter retains the previously transmitted but not successfully decoded packet of each source in the transmitter's memory for possible future retransmissions, since this packet is more likely to be decoded than a new packet from that source due to the HARQ protocol. We term a packet in the transmitter's memory an *under-process packet*. Thus, the maximum number of packets stored in the system is $I + K$ packets, i.e., I packets at the buffers of random arrival sources and K packets at the transmitter's memory.

We assume that the transmitter² can transmit at most one packet per slot. At each slot, the transmitter decides whether to send a packet or stay idle. The possible transmission options for a random arrival source $k \in \mathcal{I}$ are either transmitting the packet from its buffer or retransmitting the under-process packet from the transmitter's memory. The possible transmission options for a generate-at-will source $k \in \mathcal{J}$ are either generating and transmitting a new sample or retransmitting the under-process packet from the transmitter's memory. We refer to the packets in the buffers of the random arrival sources and to the newly generated packets of the generate-at-will sources as *fresh packets*. If the transmitter decides to transmit a fresh packet from a given source, this packet replaces the source's under-process packet in the transmitter's memory.

1) *Transmission Model:* At each slot t , the transmitter takes one of the following actions: 1) transmit a fresh packet from a source, 2) retransmit an under-process packet of a source, or 3) stay idle. Let $u_{t,k} \in \{0, 1\}$ denote the decision variable about transmitting a fresh packet from source k at slot t , where $u_{t,k} = 1$ indicates that the transmitter sends the fresh packet, and $u_{t,k} = 0$ otherwise. Let $r_{t,k} \in \{0, 1\}$ denote the decision variable about retransmitting the under-process packet of source k at slot t , where $r_{t,k} = 1$ indicates that the transmitter sends the under-process packet, and $r_{t,k} = 0$ otherwise. Since the transmitter can transmit at most one packet per slot, we have $\sum_{k \in \mathcal{K}} u_{t,k} + r_{t,k} \leq 1$.

²A mathematically equivalent system is the one where each source is equipped with an own transmitter while at most one source is allowed to transmit at each slot.

HARQ protocol: In the considered HARQ protocol, every packet transmission attempt is followed by an instantaneous error-free ACK/NACK feedback signal from the receiver. Let $d_t \in \{0, 1\}$ denote the packet reception status at slot t , where $d_t = 1$ indicates that the transmitted packet was decoded successfully (ACK), and $d_t = 0$ indicates that either the transmitted packet was not decoded successfully (NACK) or the transmitter remained idle. In the HARQ protocol, the receiver uses all previously received versions of a packet to decode it. Therefore, the probability of successfully decoding a packet is an increasing function of the number of attempted transmissions of the packet. Let $x_{t,k}$ denote the number of attempted transmissions of a packet of source k up to slot t . The evolution of $x_{t,k}$ is given as

$$x_{t+1,k} = \begin{cases} 1 & u_{t,k} = 1 \\ x_{t,k} & u_{t,k} + r_{t,k} = 0 \\ x_{t,k} + 1 & r_{t,k} = 1. \end{cases} \quad (1)$$

To account for the fact that most practical HARQ protocols allow only a finite number of re-transmissions, we limit the number of transmission attempts of a packet to x^{\max} , i.e., $x_{t,k} \leq x^{\max}$. The function representing the probability of successful decoding after $x_{t,k}$ transmissions is denoted by $f(x_{t,k})$. In practice, $f(\cdot)$ is a complicated function of several parameters such as the channel conditions, the channel coding methods, and the combining technique utilized in the HARQ protocol [27], [28].

2) *Age of Information:* The AoI is defined as the time elapsed since the generation of the most recently received status update packet at a destination. Let $\delta_{t,k}$ denote the AoI of source k at the receiver at slot t ; we refer to this simply as the AoI of source k hereinafter. We use the common assumption (see, e.g., [11]–[13], [29], [30]) that all AoI values in the system are upper bounded by δ^{\max} . Besides making the analysis tractable, this supports the fact that an AoI value exceeding a high enough upper bound carries the same timeliness information as the upper bound for the destination's decision making (e.g., control actions for a drone). To characterize the AoI of each source, we need to define the age of a fresh packet at a source and the age of an under-process packet in the transmitter's memory. These are defined in the following.

Age of the fresh packets: Let $\delta_{t,k}^f$ denote the age of the fresh packet of source k at slot t . For a random arrival source, if a packet arrives at the buffer at the beginning of slot t , the age of the fresh packet becomes zero, otherwise it is incremented by one. Let $b_{t,k} \in \{0, 1\}$ denote the packet arrival status of source $k \in \mathcal{I}$ at slot t , where $b_{t,k} = 1$ indicates a packet arrives at the buffer, and $b_{t,k} = 0$ otherwise. Note that $\Pr(b_{t,k} = 1) = \lambda_k$. For the generate-at-will sources, the transmitter can generate a fresh packet at any time so that the age of the fresh packet is always

zero. Thus, the evolution of $\delta_{t,k}^f$ with the initial value $\delta_{0,k}^f = 0$ is given as

$$\delta_{t,k}^f = \begin{cases} 0 & b_{t,k} = 1, k \in \mathcal{I} \\ \min\{\delta_{t-1,k}^f + 1, \delta^{\max}\} & b_{t,k} = 0, k \in \mathcal{I} \\ 0 & k \in \mathcal{J}, \end{cases} \quad (2)$$

Age of the under-process packets: Let $\delta_{t,k}^p$ denote the age of the under-process packet of source k at slot t . If the transmitter sends a fresh packet of source k at slot t , the age of the under-process packet of the source at the next slot drops to $\min\{\delta_{t,k}^f + 1, \delta^{\max}\}$. In other cases (i.e., retransmission or staying idle), the age of the under-process packet is incremented by one. The evolution of $\delta_{t,k}^p$ with the initial value $\delta_{0,k}^p = 0$ is given by

$$\delta_{t+1,k}^p = \begin{cases} \min\{\delta_{t,k}^f + 1, \delta^{\max}\} & u_{t,k} = 1 \\ \min\{\delta_{t,k}^p + 1, \delta^{\max}\} & \text{otherwise.} \end{cases} \quad (3)$$

AoI at the receiver: Having defined $\delta_{t,k}^f$ and $\delta_{t,k}^p$, we now characterize the evolution of the AoI at the receiver. If the transmitter sends a fresh packet of source k at slot t (i.e., $u_{t,k} = 1$) and the packet is decoded successfully at the receiver (i.e., $d_t = 1$), the AoI of the source at the next slot drops to $\min\{\delta_{t,k}^f + 1, \delta^{\max}\}$, otherwise (i.e., $d_t = 0$), the AoI increases by one. If the transmitter retransmits the under-process packet of source k (i.e., $r_{t,k} = 1$) and it is decoded successfully at the receiver, the AoI of the source at the next slot drops to $\min\{\delta_{t,k}^p + 1, \delta^{\max}\}$, otherwise (i.e., $d_t = 0$), the AoI increases by one. If, at slot t , the transmitter does not transmit any packet of source k (i.e., $u_{t,k} + r_{t,k} = 0$), the AoI of the source at the next slot increases by one. The evolution of $\delta_{t,k}$ with the initial value $\delta_{0,k} = 0$ is given as

$$\delta_{t+1,k} = \begin{cases} \min\{\delta_{t,k}^f + 1, \delta^{\max}\} & u_{t,k}d_t = 1 \\ \min\{\delta_{t,k}^p + 1, \delta^{\max}\} & r_{t,k}d_t = 1 \\ \min\{\delta_{t,k} + 1, \delta^{\max}\} & u_{t,k}(1 - d_t) = 1 \\ \min\{\delta_{t,k} + 1, \delta^{\max}\} & r_{t,k}(1 - d_t) = 1 \\ \min\{\delta_{t,k} + 1, \delta^{\max}\} & u_{t,k} + r_{t,k} = 0. \end{cases} \quad (4)$$

Note that the conditions in (4) are mutually exclusive and collectively exhaustive.

B. Problem Formulation

Our main goal is to minimize the average number of transmissions subject to the average AoI constraint by finding a transmission policy that determines the transmission decision variables at each slot t , $\{u_{t,k}, r_{t,k}\}_{k \in \mathcal{K}}$. The transmission decision for slot t is based on the age of the fresh packets, $\delta_{t,k}^f$, the age of the under-process packets, $\delta_{t,k}^p$, the AoI of each source, $\delta_{t,k}$, and the number of previous transmission attempts of each under-process packet, $x_{t,k}$.

Let $\tau_t \in \{0, 1\}$ denote the transmission status at slot t , where $\tau_t = 1$ indicates that the transmitter sends a packet, and $\tau_t = 0$ otherwise. Thus, we have

$$\tau_t = \begin{cases} 1 & \sum_{k \in \mathcal{K}} u_{t,k} + r_{t,k} = 1 \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

Let $\bar{\tau}$ denote the expected long-term time average number of transmissions, defined as

$$\bar{\tau} = \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{\tau_t\}, \quad (6)$$

where $\mathbb{E}\{\cdot\}$ is the expectation with respect to the randomness of the system (i.e., packet arrival processes of the random arrival sources and randomness in the communication channel) and the decision variables $\{u_{t,k}, r_{t,k}\}_{k \in \mathcal{K}}$. Finally, let $\bar{\delta}$ denote the expected long-term time average of AoI, given as

$$\bar{\delta} = \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{\hat{\delta}_t\}, \quad (7)$$

where $\hat{\delta}_t$ is average AoI over all sources at slot t , given as

$$\hat{\delta}_t = \frac{1}{K} \sum_{k=1}^K \delta_{t,k}. \quad (8)$$

Using (6) and (7), the main problem of this paper is formulated as the following stochastic optimization problem:

$$\text{minimize } \bar{\tau} \quad (9a)$$

$$\text{subject to } \bar{\delta} \leq \Delta^{\max} \quad (9b)$$

$$x_{t+1,k} \leq x^{\max}, \quad k \in \mathcal{K}, \quad t \in \mathbb{N} \quad (9c)$$

$$\sum_{k \in \mathcal{K}} u_{t,k} + r_{t,k} \leq 1, \quad t \in \mathbb{N} \quad (9d)$$

$$u_{t,k}, r_{t,k} \in \{0, 1\}, \quad k \in \mathcal{K}, \quad t \in \mathbb{N}, \quad (9e)$$

with variables $\{u_{t,k}, r_{t,k}\}_{k \in \mathcal{K}}$ for all $t \in \mathbb{N}$, where Δ^{\max} is the maximum allowed average AoI. The constraints of problem (9) are as follows. Inequality (9b) represents the average AoI constraint. Inequality (9c) ensures that the number of transmission attempts for each packet cannot exceed x^{\max} . Inequality (9d) ensures that the transmitter can transmit at most one packet per slot. Expression (9e) indicates the binary nature of the decision variables.

III. DETERMINISTIC TRANSMISSION POLICY

In this section, we propose a (near-optimal) solution to main problem (9) for the known environment, i.e., the packet arrival probability of each random arrival source, λ_k , and the probability of successful decoding function, $f(\cdot)$, are known. We cast problem (9) as a constrained

Markov decision process (CMDP) problem. Then, we use the Lagrangian relaxation to find a near-optimal deterministic transmission policy.

A. CMDP Formulation

The CMDP is defined by a tuple of five elements $(\mathcal{S}, \mathcal{A}_s, \mathcal{P}, c, d)$: state space, action space, state transition probabilities, and two cost functions, which are defined in the following.

State: Let $s_{t,k} = \{\delta_{t,k}^f, \delta_{t,k}^p, \delta_{t,k}, x_{t,k}\}$ denote the state of source k at slot t . The system state at slot t is defined as $s_t = \{s_{t,k}\}_{k \in \mathcal{K}} \in \mathcal{S}$, where \mathcal{S} is the state space. The initial state is denoted with $s_0 = \{s_{0,k}\}_{k \in \mathcal{K}}$, where $s_{0,k} = \{0, 0, 0, 0\}$ for all $k \in \mathcal{K}$.

Action: Let $a_t = \{a_{t,k}\}_{k \in \mathcal{K}} \in \mathcal{A}_{s_t}$ denote the action of the transmitter at slot t , where $a_{t,k} = \{u_{t,k}, r_{t,k}\}$ represents the action for source k , and \mathcal{A}_{s_t} is a space of feasible actions in state s_t , defined as $\mathcal{A}_{s_t} = \{u_{t,k}, r_{t,k} \in \{0, 1\} \mid k \in \mathcal{K}, \sum_{k \in \mathcal{K}} u_{t,k} + r_{t,k} \leq 1, r_{t,k}(x_{t,k} + 1) \leq x^{\max}\}$.

Cost functions: The CMDP has two cost functions: 1) transmission cost, defined as $c(a_t) = \tau_t$, i.e., $c(a_t) = 1$ if the transmitter makes a transmission attempt at slot t , otherwise $c(a_t) = 0$, and 2) AoI cost, defined as $d(s_t) = \hat{\delta}_t$, i.e., the average AoI over sources at slot t .

State transition probabilities: Let $\mathcal{P}(s' \mid s, a) = \Pr(s' = s_{t+1} \mid s = s_t, a = a_t)$ denote the state transition probabilities, defined as the probability of moving from current state $s = s_t$ to a next state $s' = s_{t+1}$ under action $a = a_t$. Given an action, the one-slot evolution of the AoI values (at the source, memory, and destination) and of the number of transmissions of the under-process packets is independent among the sources. Therefore, the state transition probability factorizes as $\mathcal{P}(s' \mid s, a) = \prod_{k \in \mathcal{K}} \Pr(s_{t+1,k} \mid s_{t,k}, a_{t,k})$. Let us denote $\tilde{\delta}_{t,k}^f \triangleq \min\{\delta_{t,k}^f + 1, \delta^{\max}\}$, $\tilde{\delta}_{t,k}^p \triangleq \min\{\delta_{t,k}^p + 1, \delta^{\max}\}$, $\tilde{\delta}_{t,k} \triangleq \min\{\delta_{t,k} + 1, \delta^{\max}\}$, $\bar{f}(\cdot) \triangleq 1 - f(\cdot)$, and $\bar{\lambda}_k \triangleq 1 - \lambda_k$. Given the state $s_{t,k} = \{\delta_{t,k}^f, \delta_{t,k}^p, \delta_{t,k}, x_{t,k}\}$, the state transition probabilities for a random arrival source $k \in \mathcal{I}$ under different actions can be expressed as

$$\Pr(\{0, \tilde{\delta}_{t,k}^f, \tilde{\delta}_{t,k}^f, 1\} \mid s_{t,k}, a_{t,k} = \{1, 0\}) = f(1)\lambda_k \quad (10a)$$

$$\Pr(\{\tilde{\delta}_{t,k}^f, \tilde{\delta}_{t,k}^f, \tilde{\delta}_{t,k}^f, 1\} \mid s_{t,k}, a_{t,k} = \{1, 0\}) = f(1)\bar{\lambda}_k \quad (10b)$$

$$\Pr(\{0, \tilde{\delta}_{t,k}^f, \tilde{\delta}_{t,k}^f, 1\} \mid s_{t,k}, a_{t,k} = \{1, 0\}) = \bar{f}(1)\lambda_k \quad (10c)$$

$$\Pr(\{\tilde{\delta}_{t,k}^f, \tilde{\delta}_{t,k}^f, \tilde{\delta}_{t,k}^f, 1\} \mid s_{t,k}, a_{t,k} = \{1, 0\}) = \bar{f}(1)\bar{\lambda}_k \quad (10d)$$

$$\Pr(\{0, \tilde{\delta}_{t,k}^p, \tilde{\delta}_{t,k}^p, x_{t,k} + 1\} \mid s_{t,k}, a_{t,k} = \{0, 1\}) = f(x_{t,k} + 1)\lambda_k \quad (10e)$$

$$\Pr(\{\tilde{\delta}_{t,k}^f, \tilde{\delta}_{t,k}^p, \tilde{\delta}_{t,k}^p, x_{t,k} + 1\} \mid s_{t,k}, a_{t,k} = \{0, 1\}) = f(x_{t,k} + 1)\bar{\lambda}_k \quad (10f)$$

$$\Pr(\{0, \tilde{\delta}_{t,k}^p, \tilde{\delta}_{t,k}^p, x_{t,k} + 1\} \mid s_{t,k}, a_{t,k} = \{0, 1\}) = \bar{f}(x_{t,k} + 1)\lambda_k \quad (10g)$$

$$\Pr(\{\tilde{\delta}_{t,k}^f, \tilde{\delta}_{t,k}^p, \tilde{\delta}_{t,k}^p, x_{t,k} + 1\} \mid s_{t,k}, a_{t,k} = \{0, 1\}) = \bar{f}(x_{t,k} + 1)\bar{\lambda}_k \quad (10h)$$

$$\Pr(\{0, \tilde{\delta}_{t,k}^p, \tilde{\delta}_{t,k}^p, x_{t,k}\} \mid s_{t,k}, a_{t,k} = \{0, 0\}) = \lambda_k \quad (10i)$$

$$\Pr(\{\tilde{\delta}_{t,k}^f, \tilde{\delta}_{t,k}^p, \tilde{\delta}_{t,k}^p, x_{t,k}\} \mid s_{t,k}, a_{t,k} = \{0, 0\}) = \bar{\lambda}_k, \quad (10j)$$

whereas the other cases are zero. The state transition probabilities for the generate-at-will source $k \in \mathcal{J}$ are obtained by substituting $\lambda_k = 1$ in (10).

Let π denote a policy that determines the action taken at each state. A stationary randomized policy is mapping from each state to a distribution over actions, $\pi(a \mid s) : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, $\sum_{a \in \mathcal{A}_s} \pi(a \mid s) = 1$. A (stationary) deterministic policy chooses an action at a given state with probability one, which is a special case of the stationary randomized policy. With a slight abuse of notation, we denote the action taken in state s by a deterministic policy π with $\pi(s)$. Let $\bar{\tau}^\pi = \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{c(a_t) \mid s_0\}$ denote the average number of transmissions (see (6)), obtained under policy π starting from the initial state s_0 . Let $\bar{\delta}^\pi = \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{d(s_t) \mid s_0\}$ denote average AoI (see (7)), obtained under policy π starting from the initial state s_0 . Having constructed the CMDP, problem (9) is equivalently cast as the CMDP problem

$$\begin{aligned} & \underset{\pi}{\text{minimize}} && \bar{\tau}^\pi \\ & \text{subject to} && \bar{\delta}^\pi \leq \Delta^{\max}. \end{aligned} \quad (11)$$

An optimal policy that solves CMDP problem (11) is denoted with π^* , and the optimal value of the problem is denoted with $\bar{\tau}^*$.

Similarly to [11], [31], [32], to solve problem (11), we need to make extra assumptions about the CMDP structure. Specifically, we assume that given the initial state (s_0), all policies will induce a Markov chain with only one recurrent class and a (possibly empty) set of transient states. This assumption makes problem (11) well-posed so that we can use the tools associated

with the unichain MDPs, as described in the next section.

B. Solution of the CMDP Problem

To solve CMDP problem (11), we apply the Lagrangian relaxation method to transform the CMDP problem to an (unconstrained) MDP problem, parametrized by a Lagrange dual variable [33, Sec. 3.3]. In comparison to the CMDP problem, the MDP problem has only one cost function that is defined as $L(s, a, \beta) = c(a_t) + \beta d(s_t)$, whereas the other elements, i.e., the state space, action space, and state transition probabilities, are the same. Let $\bar{L}(\pi, \beta) = \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{c(a_t) + \beta(d(s_t) - \Delta^{\max})\}$ denote the Lagrangian corresponding to CMDP problem (11), where β is the Lagrangian multiplier. Following the standard Lagrangian relaxation procedure, we restrict to the set of deterministic policies and construct the following MDP problem associated with the CMDP problem (11)

$$\underset{\pi \in \Pi_D}{\text{minimize}} \quad \bar{L}(\pi, \beta), \quad (12)$$

where Π_D is the set of all deterministic policies. Let π_β^* denote an optimal policy that solves problem (12) for a given β , which is called a β -optimal policy.

The following remark expresses the relation between the optimal values of CMDP problem (11) and the MDP problem (12).

Remark 1. *The cost function in the objective of CMDP problem (11) is bounded below, i.e., $c(a_t) \geq 0$ for all $t \in \mathbb{N}$. Moreover, the state space, S , is finite. Therefore, the two conditions in [33, Corollary 12.2] are satisfied in our CMDP formulation, and we have*

$$\bar{\tau}^* = \sup_{\beta \geq 0} \min_{\pi \in \Pi_D} \bar{L}(\pi, \beta). \quad (13)$$

According to Remark 1, the optimal value of CMDP problem (11), $\bar{\tau}^*$, is obtained via the solution of the right hand side of (13), which means finding the optimal Lagrangian multiplier β^* and its corresponding β^* -optimal policy, $\pi_{\beta^*}^*$. If policy $\pi_{\beta^*}^*$ satisfies the constraint of CMDP problem (11) with equality, i.e., $\bar{\delta}^{\pi_{\beta^*}^*} = \Delta^{\max}$, then $\pi_{\beta^*}^*$ is an optimal policy for the CMDP problem, i.e., $\pi^* = \pi_{\beta^*}^*$. However, due to the discrete nature of the action space, in general, there is no guarantee that $\pi_{\beta^*}^*$ satisfies the constraint with equality. To elaborate this further, the following remark presents the structure of an optimal policy π^* .

Remark 2. *An optimal policy for CMDP problem (11), π^* , is a randomized mixture of two deterministic $\tilde{\beta}$ -optimal policies, from which one policy satisfies the constraint and the other*

one violates it. The two policies are mixed with a randomization factor such that the obtained optimal policy satisfies $\bar{\delta}^{\pi^*} = \Delta^{\max}$ [10], [11], [13].

According to Remark 2, two deterministic $\tilde{\beta}$ -optimal policies and the optimal randomization factor to mix between these policies should be found to obtain an optimal policy, π^* . However, finding these becomes readily computationally intractable even for the moderate number of states, especially because oftentimes, the optimal randomization factor can be found only numerically [34, Section 3.2]. Therefore, in order to solve CMDP problem (11), we propose a practical deterministic policy, which is numerically shown to provide near-optimal performance in Section V. More specifically, we develop an iterative algorithm based on bisection and the relative value iteration algorithm (RVIA), as summarized in Algorithm 1. In brief, at each iteration, we find a β -optimal policy for a given β via the RVIA and subsequently update β according to the bisection rule. The iterative procedure continues until the best β -optimal policy among the feasible β -optimal policies is found. In the next two subsections, we delve into details of this procedure.

1) *Algorithm to Find a β -optimal Policy:* To obtain an optimal policy π_β^* for a given β , we solve the MDP problem (12) via RVIA. By [35, Theorem 8.4.3], there exists a relative value function $h(s)$, $s \in \mathcal{S}$, that satisfies

$$\bar{L}^*(\beta) + h(s) = \min_{a \in \mathcal{A}_s} \left[L(s, a, \beta) + \sum_{s' \in \mathcal{S}} \Pr(s' | s, a) h(s') \right], \text{ for all } s \in \mathcal{S}, \quad (14)$$

where $\bar{L}^*(\beta)$ is the optimal value of the MDP problem (12) for a given β , defined as $\bar{L}^*(\beta) = \min_{\pi \in \Pi_D} \bar{L}(\pi, \beta)$. Subsequently, the β -optimal policy, π_β^* , is obtained as [35, Theorem 8.4.4]

$$\pi_\beta^*(s) = \arg \min_{a \in \mathcal{A}_s} \left[L(s, a, \beta) + \sum_{s' \in \mathcal{S}} \Pr(s' | s, a) h(s') \right], \text{ for all } s \in \mathcal{S}. \quad (15)$$

To obtain the β -optimal policy, we use the RVIA, in which the relative value function for all states $s \in \mathcal{S}$ at each iteration $i \in \{0, 1, \dots\}$ is updated as $h^i(s) = v^i(s) - v^i(s^{\text{ref}})$. Where $s^{\text{ref}} \in \mathcal{S}$ is an arbitrary reference state which remains unchanged throughout the iterations. The term $v^i(s)$, called value function, is obtained at each iteration as

$$v^i(s) = \min_{a \in \mathcal{A}_s} \left[L(s, a, \beta) + \sum_{s' \in \mathcal{S}} \Pr(s' | s, a) h^{i-1}(s') \right]. \quad (16)$$

For any state $s \in \mathcal{S}$ and initialization $v^0(s)$, the sequences $\{h^i(s)\}_{i=1,2,\dots}$ and $\{v^i(s)\}_{i=1,2,\dots}$ converge, i.e., $\lim_{i \rightarrow \infty} h^i(s) = h(s)$ and $\lim_{i \rightarrow \infty} v^i(s) = v(s)$. The RVI algorithm to find a β -optimal policy is presented in Steps 3-12 of Algorithm 1. After the convergence of RVIA, i.e., convergence of the relative value function, $h(\cdot)$, and the value function, $v(\cdot)$ (see Steps 3-9 in Algorithm 1), we obtain the β -optimal policy, π_β^* , according to (15) (see Steps 10-12 in

Algorithm 1). It is worth noting that the optimal value of the MDP problem (12) for a given β is given by $\bar{L}^*(\beta) = v(s^{\text{ref}})$.

2) *Algorithm to Find the Optimal Lagrangian Multiplier:* According to [10, Lemma 3.1], for a given β -optimal policy (π_β^*), the objective function of the CMDP problem, $\bar{\tau}^{\pi_\beta^*}$, and the objective function of the MDP problem, $\bar{L}^*(\beta)$, are increasing in β , while the constraint of the CMDP problem, $\bar{\delta}^{\pi_\beta^*}$, is decreasing in β . Therefore, we are interested in the smallest Lagrangian multiplier that satisfies the constraint in CMDP problem (11), defined as

$$\tilde{\beta} \triangleq \inf \{ \beta \geq 0 \mid \bar{\delta}^{\pi_\beta^*} \leq \Delta^{\max} \}. \quad (17)$$

To search for $\tilde{\beta}$, we use the bisection algorithm which takes advantage of the monotonicity of $\bar{\delta}^{\pi_\beta^*}$ with respect to β , as presented in Algorithm 1 (see Steps 1-18). We initialize the bisection algorithm with β_u and β_l in such a way that $\bar{\delta}^{\pi_{\beta_u}^*} \leq \Delta^{\max}$ and $\bar{\delta}^{\pi_{\beta_l}^*} \geq \Delta^{\max}$, which also implies $\beta_u \geq \beta_l$. The algorithm termination criterion is $\beta_u - \beta_l < \kappa$, where κ is a sufficiently small constant. After termination of the bisection algorithm, we set $\tilde{\beta} = \beta_u$ and obtain the best feasible β -optimal policy as $\pi_{\tilde{\beta}}^* = \pi_{\beta_u}^*$. Moreover, the algorithm returns the infeasible policy associated with β_l , which represents a lower-bound to an optimal solution of (11).

IV. LYAPUNOV-BASED TRANSMISSION SCHEDULING POLICIES

In this section, we use the Lyapunov optimization theory to derive a solution for problem (9). According to the deterministic transmission policy (Algorithm 1), RVIA needs to explore all states and actions. When the number of sources increases, the number of states grows exponentially. Besides this, the RVIA is run for each bisection iteration. These make obtaining the policy computationally inefficient. Therefore, we develop a low-complexity dynamic transmission (LC-DT) policy using the DPP method for the known environment in Section IV-A. The numerical results in Section V show that LC-DT policy performs close to the optimal solution.

Furthermore, we develop a Lyapunov-based transmission policy for the unknown environment in Section IV-B. We transform CMDP problem (11) into an MDP problem using the DPP method. Then, by using deep Q-learning (DQL), we solve the MDP problem and provide the learning-based transmission policy.

A. Low-complexity Dynamic Transmission Policy: Known Environment

We use the Lyapunov drift-plus-penalty (DPP) method [14], where average AoI constraint (9b) is enforced by transforming it into a queue stability constraint. In particular, the constraint

Algorithm 1: The deterministic transmission policy to solve CMDP problem (11)

Input: 1) System parameters: Δ^{\max} , $f(\cdot)$, and λ_k for all $k \in \mathcal{I}$, 2) RVI parameters: s^{ref} , ϵ , and 3) Bisection parameters: β_u , β_l , κ

/*Bisection algorithm

```

1 while  $\beta_u - \beta_l \geq \kappa$  do
2    $\bar{\beta} = \frac{\beta_u + \beta_l}{2}$ 
   /*RVIA for the given  $\bar{\beta}$ 
   Initialize: 1)  $i = 1$ , 2) set  $h^0(s) = 1, h^1(s) = 0, v^0(s) = 0$  for all  $s \in \mathcal{S}$ 
3   while  $\max_{s \in \mathcal{S}} |h^i(s) - h^{i-1}(s)| \geq \epsilon$  do
4      $i = i + 1$ 
5     for  $s \in \mathcal{S}$  do
6        $v^i(s) = \min_{a \in \mathcal{A}_s} [L(s, a, \bar{\beta}) + \sum_{s' \in \mathcal{S}} \Pr(s' | s, a) h^{i-1}(s')]$ 
7        $h^i(s) = v^i(s) - v^i(s^{\text{ref}})$ 
8     end
9   end
   /*An optimal policy for given  $\bar{\beta}$ 
10  for  $s \in \mathcal{S}$  do
11     $\pi_{\bar{\beta}}^*(s) = \arg \min_{a \in \mathcal{A}_s} [L(s, a, \bar{\beta}) + \sum_{s' \in \mathcal{S}} \Pr(s' | s, a) h^i(s')]$ 
12  end
13  if  $\bar{\delta}^{\pi_{\bar{\beta}}^*} \leq \Delta^{\max}$  then
14     $\beta_u = \bar{\beta}$ 
15  else
16     $\beta_l = \bar{\beta}$ 
17  end
18 end

```

Output: Lagrangian multiplier: $\tilde{\beta} = \beta_u$, feasible policy: $\pi_{\tilde{\beta}}^* = \pi_{\beta_u}^*$, (infeasible)
lower-bound policy: $\pi_{\beta_l}^*$

is mapped into a virtual queue so that the stability of the virtual queue implies the feasibility of the constraint.

Let Q_t denote the virtual queue associated with average AoI constraint (9b) at slot t . The virtual queue evolves as follows:

$$Q_{t+1} = \max\{Q_t - \Delta^{\max} + \hat{\delta}_{t+1}, 0\}. \quad (18)$$

To ensure that average AoI constraint (9b) is satisfied, we use the notion of strong stability, where the virtual queue is (strongly) stable if $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}\{Q_t\} < \infty$ [14, Chapter 2]. According to the DPP method, the strong stability of queue (18) implies that average AoI constraint (9b) is satisfied [14, Chapter 4].

To define the queue stability condition, we introduce the Lyapunov function and its drift. A quadratic Lyapunov function is defined as $L(Q_t) = \frac{1}{2}Q_t^2$ [14, Chapter 3]. The Lyapunov function

measures the network congestion and thus, by minimizing the expected change of the Lyapunov function from one slot to the next slot, the virtual queue can be stabilized [14, Chapter 4].

Let $o_t = \{\{\delta_{t,k}^f, \delta_{t,k}^p, \delta_{t,k}, x_{t,k}\}_{k \in \mathcal{K}}, Q_t\}$ denote the network state at slot t . The conditional Lyapunov drift, $\alpha(o_t)$, is defined as the expected change in the Lyapunov function over one slot, given the network state at slot t , i.e., $\alpha(o_t) = \mathbb{E}\{L(Q_{t+1}) - L(Q_t) \mid o_t\}$ [14, Chapter 4].

By following the DPP minimization approach [14, Chapter 3], a solution for (9) can be derived by solving the following problem at each slot t

$$\text{minimize } V\mathbb{E}\{\tau_t \mid o_t\} + \alpha(o_t) \quad (19a)$$

$$\text{subject to } x_{t+1,k} \leq x^{\max}, \quad k \in \mathcal{K} \quad (19b)$$

$$\sum_{k \in \mathcal{K}} u_{t,k} + r_{t,k} \leq 1 \quad (19c)$$

$$u_{t,k}, r_{t,k} \in \{0, 1\}, \quad k \in \mathcal{K}, \quad (19d)$$

with variables $\{u_{t,k}, r_{t,k}\}_{k \in \mathcal{K}}$. The objective function of problem (19) represents the DPP function, in which the positive parameter V is used to adjust the tradeoff between minimizing original objective function (9a) and the size of the (virtual) queue backlog. A larger value of V puts more emphasis on the original objective function, i.e., minimizing the average number of transmissions.

According to the standard procedure used in the DPP method, an upper bound for the drift part is derived, whereas the penalty part (i.e., the original objective function) remains unchanged [14, Chapter 4]. It is worth stressing that using such upper bound of the conditional Lyapunov drift in the optimization procedure does not affect the virtual queue's stabilizing logic. We derive the upper bound using the following inequality, where for any $\tilde{\gamma} \geq 0$, $\bar{\gamma} \geq 0$, and $\hat{\gamma} \geq 0$, we have [29]

$$(\max\{\tilde{\gamma} - \bar{\gamma} + \hat{\gamma}, 0\})^2 \leq \tilde{\gamma}^2 + \bar{\gamma}^2 + \hat{\gamma}^2 + 2\tilde{\gamma}(\hat{\gamma} - \bar{\gamma}). \quad (20)$$

By applying (20) to (18), an upper bound for Q_{t+1}^2 is given as

$$Q_{t+1}^2 \leq Q_t^2 + (\Delta^{\max})^2 + \hat{\delta}_{t+1}^2 + 2Q_t(\hat{\delta}_{t+1} - \Delta^{\max}). \quad (21)$$

By applying (21) to (15), the upper bound of objective function (19a) is given as

$$\begin{aligned} V\mathbb{E}\{\tau_t \mid o_t\} + \alpha(o_t) &\leq V \sum_{k \in \mathcal{K}} \mathbb{E}\{u_{t,k} \mid o_t\} + \mathbb{E}\{r_{t,k} \mid o_t\} + \frac{1}{2}\mathbb{E}\{(\Delta^{\max})^2 + \hat{\delta}_{t+1}^2 \\ &+ 2Q_t(\hat{\delta}_{t+1} - \Delta^{\max}) \mid o_t\} = V \sum_{k \in \mathcal{K}} \mathbb{E}\{u_{t,k} \mid o_t\} + \mathbb{E}\{r_{t,k} \mid o_t\} \\ &+ \frac{1}{2}((\Delta^{\max})^2 + \mathbb{E}\{\hat{\delta}_{t+1}^2 \mid o_t\} + 2Q_t(\mathbb{E}\{\hat{\delta}_{t+1} \mid o_t\} - \Delta^{\max})). \end{aligned} \quad (22)$$

To complete the derivation of (22), we calculate $\mathbb{E}\{\hat{\delta}_{t+1} \mid o_t\}$ and $\mathbb{E}\{\hat{\delta}_{t+1}^2 \mid o_t\}$, which are given by the following lemmas.

Lemma 1. *The conditional expectation $\mathbb{E}\{\hat{\delta}_{t+1} \mid o_t\}$ is given as*

$$\begin{aligned} \mathbb{E}\{\hat{\delta}_{t+1} \mid o_t\} &= \frac{1}{K} \sum_{k \in \mathcal{K}} \mathbb{E}\{u_{t,k} \mid o_t\} f(1) \tilde{\delta}_{t,k}^f + \mathbb{E}\{r_{t,k} \mid o_t\} f(x_{t,k} + 1) \tilde{\delta}_{t,k}^p \\ &+ [1 - f(1) \mathbb{E}\{u_{t,k} \mid o_t\} - f(x_{t,k} + 1) \mathbb{E}\{r_{t,k} \mid o_t\}] \tilde{\delta}_{t,k}. \end{aligned} \quad (23)$$

Proof. The proof is presented in Appendix A. □

Lemma 2. *The conditional expectation $\mathbb{E}\{\hat{\delta}_{t+1}^2 \mid o_t\}$ is given as*

$$\begin{aligned} \mathbb{E}\{\hat{\delta}_{t+1}^2 \mid o_t\} &= \frac{1}{K^2} \left[\sum_{k \in \mathcal{K}} \mathbb{E}\{u_{t,k} \mid o_t\} f(1) (\tilde{\delta}_{t,k}^f)^2 + \mathbb{E}\{r_{t,k} \mid o_t\} f(x_{t,k} + 1) (\tilde{\delta}_{t,k}^p)^2 \right. \\ &+ [1 - f(1) \mathbb{E}\{u_{t,k} \mid o_t\} - f(x_{t,k} + 1) \mathbb{E}\{r_{t,k} \mid o_t\}] (\tilde{\delta}_{t,k})^2 \\ &+ \sum_{k \in \mathcal{K}} \sum_{k' \in \mathcal{K}, k' \neq k} \mathbb{E}\{u_{t,k} \mid o_t\} f(1) \tilde{\delta}_{t,k}^f \tilde{\delta}_{t,k'} + \mathbb{E}\{r_{t,k} \mid o_t\} f(x_{t,k} + 1) \tilde{\delta}_{t,k}^p \tilde{\delta}_{t,k'} \\ &+ \mathbb{E}\{u_{t,k'} \mid o_t\} f(1) \tilde{\delta}_{t,k'}^f \tilde{\delta}_{t,k} + \mathbb{E}\{r_{t,k'} \mid o_t\} f(x_{t,k'} + 1) \tilde{\delta}_{t,k'}^p \tilde{\delta}_{t,k} - \mathbb{E}\{u_{t,k} \mid o_t\} f(1) \tilde{\delta}_{t,k} \tilde{\delta}_{t,k'} \\ &- \mathbb{E}\{r_{t,k} \mid o_t\} f(x_{t,k} + 1) \tilde{\delta}_{t,k} \tilde{\delta}_{t,k'} - \mathbb{E}\{u_{t,k'} \mid o_t\} f(1) \tilde{\delta}_{t,k'} \tilde{\delta}_{t,k} \\ &\left. - \mathbb{E}\{r_{t,k'} \mid o_t\} f(x_{t,k'} + 1) \tilde{\delta}_{t,k'} \tilde{\delta}_{t,k} + \tilde{\delta}_{t,k} \tilde{\delta}_{t,k} \right]. \end{aligned} \quad (24)$$

Proof. The proof is presented in Appendix B. □

Having derived the upper bound of the DPP in (22), our goal is to minimize (22) subject to constraints (19b)–(19d) with variables $\{u_{t,k}, r_{t,k}\}_{k \in \mathcal{K}}$. According to the standard procedure, we drop the expectations in (22) [14, Page 59]. We denote the upper bound of the DPP after dropping the expectations at slot t by W_t , which, as shown in Appendix C, can be expressed as

$$\begin{aligned} W_t &= V \sum_{k \in \mathcal{K}} u_{t,k} + r_{t,k} + \frac{1}{2K^2} \left[\sum_{k \in \mathcal{K}} u_{t,k} f(1) (\tilde{\delta}_{t,k}^f)^2 + r_{t,k} f(x_{t,k} + 1) (\tilde{\delta}_{t,k}^p)^2 \right. \\ &+ [1 - u_{t,k} f(1) - r_{t,k} f(x_{t,k} + 1)] (\tilde{\delta}_{t,k})^2 + \sum_{k \in \mathcal{K}} \sum_{k' \in \mathcal{K}, k' \neq k} u_{t,k} f(1) \tilde{\delta}_{t,k}^f \tilde{\delta}_{t,k'} \\ &+ r_{t,k} f(x_{t,k} + 1) \tilde{\delta}_{t,k}^p \tilde{\delta}_{t,k'} - u_{t,k} f(1) \tilde{\delta}_{t,k} \tilde{\delta}_{t,k'} - r_{t,k} f(x_{t,k} + 1) \tilde{\delta}_{t,k} \tilde{\delta}_{t,k'} + u_{t,k'} f(1) \tilde{\delta}_{t,k'}^f \tilde{\delta}_{t,k} \\ &+ r_{t,k'} f(x_{t,k'} + 1) \tilde{\delta}_{t,k'}^p \tilde{\delta}_{t,k} - u_{t,k'} f(1) \tilde{\delta}_{t,k'} \tilde{\delta}_{t,k} - r_{t,k'} f(x_{t,k'} + 1) \tilde{\delta}_{t,k'} \tilde{\delta}_{t,k} \\ &+ \tilde{\delta}_{t,k} \tilde{\delta}_{t,k} + 2KQ_t \sum_{k \in \mathcal{K}} u_{t,k} f(1) \tilde{\delta}_{t,k}^f + r_{t,k} f(x_{t,k} + 1) \tilde{\delta}_{t,k}^p \\ &\left. + [1 - u_{t,k} f(1) - r_{t,k} f(x_{t,k} + 1)] \tilde{\delta}_{t,k} \right] + \frac{1}{2} [(\Delta^{\max})^2 - 2Q_t \Delta^{\max}]. \end{aligned}$$

Accordingly, the solution of (19) can be derived by solving the following problem

$$\text{minimize } W_t \quad (25a)$$

$$\text{subject to } (19b) - (19d), \quad (25b)$$

with variables $\{u_{t,k}, r_{t,k}\}_{k \in \mathcal{K}}$.

The proposed LC-DT policy is summarized in Algorithm 2. In Step 1, the transmitter updates the age of random arrival sources' fresh packets at slot t , $\delta_{t,k}^f$, using (2). In Step 2, given the current network state, it solves problem (25) to find the optimal transmission decision at slot t . In Step 3, the network state (except $\delta_{t,k}^f$ for all sources) is updated based on the current network state and the obtained transmission decision.

Algorithm 2: Proposed low-complexity dynamic transmission (LC-DT) policy

Initialize: Set V , and initialize o_1 .

- 1 **for** $t = 1, 2, 3, \dots$ **do**
- 2 Step 1: Update $\delta_{t,k}^f$ using (2)
- 3 Step 2: Find decision variables $\{u_{t,k}, r_{t,k}\}_{k \in \mathcal{K}}$ by solving problem (25)
- 4 Step 3: Update $x_{t+1,k}$ using (1), $\delta_{t+1,k}^p$ using (3), $\delta_{t+1,k}$ using (4), and Q_{t+1} using (18)
- 5 **end**

Regarding finding the solution to problem (25), we observe that the number of feasible transmission decisions at each slot is a moderate number. Namely, the transmission options over the K sources are sending a fresh packet (i.e., K options), sending an under-process packet (i.e., K options), or staying idle (i.e., 1 option). On the other hand, if the under-process packet of a source was sent x^{\max} times, the transmitter cannot send that packet. Thus, the number of feasible actions at each slot, i.e., the number of feasible combinations of variables $\{u_{t,k}, r_{t,k}\}_{k \in \mathcal{K}}$ in problem (25), is at most $2K + 1$. Because $2K + 1$ is a small value for a reasonable system (increasing only linearly in K), we use the exhaustive search algorithm to solve problem (25).

B. Learning-based Transmission Policy: Unknown Environment

In this subsection, we assume that the environment is unknown, i.e., the transmitter does not know the system statistics, namely: 1) the packet arrival probability of each random arrival source, λ_k , and 2) the probability of successful decoding function, $f(\cdot)$. In this scenario, we transform CMDP problem (11) into an unconstrained MDP problem via the Lyapunov DPP method. Then, we utilize the DQL algorithm [15] to solve the MDP problem. Even though the DQL algorithm cannot ensure the optimality of the solution, the algorithm can be applied 1) without knowing the system statistics, 2) in a system with a large state and action space, and 3) without upper bounding AoI.

Inspired by [36] and using the results of Section IV-A, we transform CMDP problem (11) into an MDP problem, in which our goal is to minimize the time average DPP function

$$\underset{\pi}{\text{minimize}} \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \{V\tau_t + L(Q_{t+1}) - L(Q_t) \mid o_t\}. \quad (26)$$

The system state of the MDP at slot t is $o_t = \{\{\delta_{t,k}^f, \delta_{t,k}^p, \delta_{t,k}, x_{t,k}\}_{k \in \mathcal{K}}, Q_t\}$ (defined in Section IV-A), and the action at time slot t is $a_t = \{u_{t,k}, r_{t,k}\}_{k \in \mathcal{K}} \in \mathcal{A}_{s_t}$ (defined in Section III-A). The cost function \tilde{c}_t , which is defined as the DPP function, is given as $\tilde{c}_t = V\tau_t + L(Q_{t+1}) - L(Q_t)$.

Note that the state transition probabilities of the MDP problem (26) are unknown, as the environment is unknown.

To solve the MDP problem (26), we use the DQL algorithm in [15, Algorithm 1]. According to the DQL algorithm, an action is selected at each state to maximize a cumulative discounted immediate reward. As we aim to minimize cost function \tilde{c}_t , the immediate reward is defined as $r_t = -\tilde{c}_t$. The implementation of DQL, along with the parameters, is presented in Section V.

V. NUMERICAL RESULTS

In this section, we evaluate the performance of the proposed transmission scheduling policies, namely: 1) the deterministic transmission policy presented in Algorithm 1 in Section III, 2) the low-complexity dynamic transmission (LC-DT) policy presented in Algorithm 2 in Section IV-A, and 3) the learning-based transmission policy presented in Section IV-B.

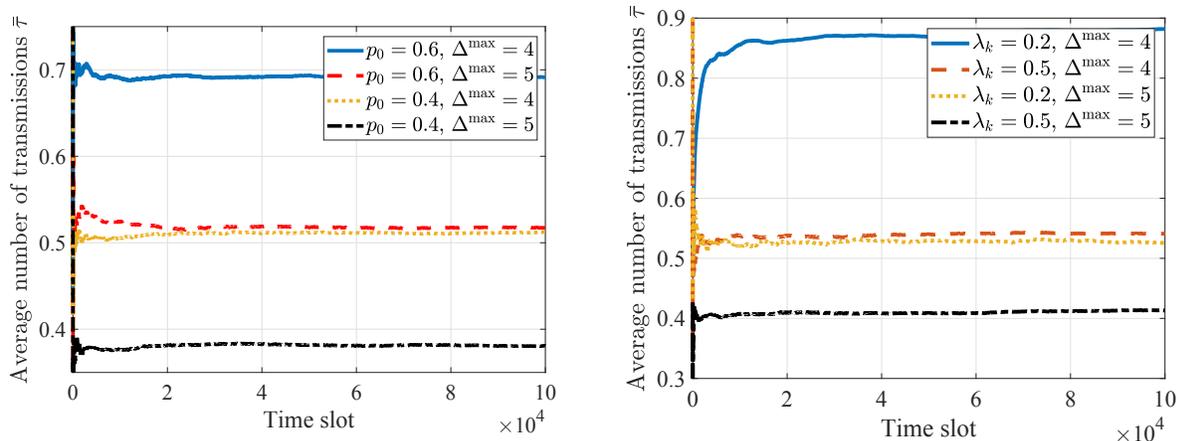
For the probability of successful decoding, we use the function in [13], i.e., $f(x_{t,k}) = 1 - p_0 \eta^{x_{t,k}-1}$, where $p_0 \in [0, 1]$ is the error probability of the first transmission of a packet and $\eta \in [0, 1]$ determines the effectiveness of the HARQ protocol. Unless otherwise specified, we consider one random arrival source and one generate-at-will source, i.e., $K = 2$, and we set $\delta^{\max} = 18$ and $x^{\max} = 5$. The rest of the parameters are specified in each figure.

A. Deterministic Transmission Policy

In Fig. 2, we evaluate the deterministic transmission policy and the impact of system parameters on the performance. For Algorithm 1, we set the bounds on the Lagrangian multiplier as $\beta_u = 1$, $\beta_l = 0$, the bisection stopping criterion as $\kappa = 0.005$, and the RVIA stopping criterion as $\epsilon = 0.01$.

Fig. 2(a) illustrates the evolution of average number of transmissions, $\bar{\tau}$, with respect to time slots for different values of maximum allowable average AoI, Δ^{\max} , and the error probability of the first transmission, p_0 . It can be seen from Fig. 2(a) that $\bar{\tau}$ increases when p_0 increases. This behavior is due to the fact that when p_0 increases, the probability of successful decoding decreases, and thus, more transmission attempts are needed to meet the average AoI constraint. For example, when $\Delta^{\max} = 4$, by increasing p_0 from 0.4 to 0.6, $\bar{\tau}$ increases by about 50 %. In addition, $\bar{\tau}$ decreases when Δ^{\max} increases, because the transmitter needs fewer transmission attempts to satisfy the AoI constraint.

Fig. 2(b) shows the evolution of $\bar{\tau}$ with respect to time slots for the different packet arrival rate, λ_k , and Δ^{\max} . From Fig. 2(b), it can be seen that when λ_k decreases, the average number



(a) The evolution of $\bar{\tau}$ versus time slots for different p_0 with $\lambda_k = 0.7$ for all $k \in \mathcal{I}$. (b) The evolution of $\bar{\tau}$ versus time slots for different λ_k for $k \in \mathcal{I}$ with $p_0 = 0.4$.

Figure 2: The performance of the deterministic transmission policy for different Δ^{\max} with $\eta = 0.4$.

of transmissions increases dramatically. For example, when $\Delta^{\max} = 4$, by decreasing λ_k from 0.5 to 0.2, the value of $\bar{\tau}$ increases by about 75 %. This is because when λ_k decreases, the availability of the fresh packets at the random arrival source decreases, and consequently, the AoI of this source increases. In this case, to satisfy the AoI constraint, the transmitter must send the generate-at-will source's packets more frequently to compensate for the negative effect of the random arrival sources on the average AoI.

B. LC-DT Policy

In this subsection, we evaluate the performance of the LC-DT policy. Fig. 3(a) depicts the evolution of the average number of transmissions, $\bar{\tau}$, with respect to time slot for different values of the DPP trade-off parameter, V . Fig. 3(b) depicts the average AoI with respect to time slot for different values of V .

According to the figure, when V increases, the average number of transmissions decreases and the average AoI increases. This is because by increasing V , the algorithm puts more emphasis on minimizing the average number of transmissions. In addition, Fig. 3(b) shows that for any value of V , the algorithm satisfies the constraint. Furthermore, for $V \geq 20$, the obtained average AoI is about the maximum allowable average AoI; also, increasing V beyond this value leads to negligible improvement in the objective function (see Fig. 3(a)). Therefore, it is beneficial to set V larger than 30 when utilizing the LC-DT policy.

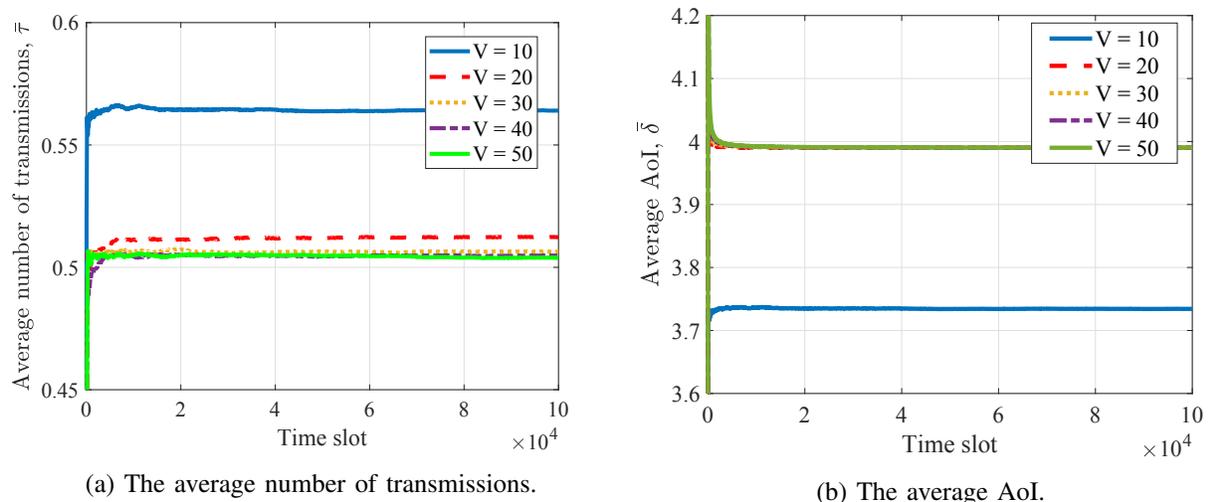


Figure 3: The performance of the LC-DT policy versus time slots for different values of V , where $\eta = 0.4$, $p_0 = 0.4$, $\lambda_k = 0.7$ for all $k \in \mathcal{I}$, and $\Delta^{\max} = 4$.

C. Learning-based Transmission Policy

In this section, we examine the learning-based transmission policy by evaluating the average number of transmissions, $\bar{\tau}$. To implement DQL algorithm, we use a fully-connected deep neural network with two hidden layers. Each hidden layer has 256 neurons with *ReLU* activation function. The optimizer is *Adam*, the mini-batch size is 32, and the replay memory size is 100000. We set the learning rate as 0.001, the discount factor as 0.99, the number of steps per episode as 1000, and the target network update rate as $\frac{1}{500}$. The convergence of the policy for different values of V are shown in Fig. 4(a). It can be seen that by taking about 800 episodes, the DQL-agent (transmitter) is learned. After this, we stop the learning process and use the learned transmitter to operate in the system. Fig. 4(b) shows the performance of the algorithm by evaluating $\bar{\tau}$ with respect to time slots for different values of V . Similar to Fig. 3, it can be seen that the performance is not considerably improved for V larger than 30.

D. Comparison of the Proposed Policies

Fig. 5 shows the average number of transmissions, $\bar{\tau}$, as a function of Δ^{\max} under different policies. In this figure, we plot the (infeasible) lower-bound policy, obtained in Algorithm 1, as a benchmark. Furthermore, we consider a (feasible) baseline policy, where the transmitter sends a packet whenever the average AoI reaches Δ^{\max} . In every transmission attempt, the source with larger AoI is selected; if there are multiple sources with the largest AoI, one of them is selected randomly. The policy employs an HARQ protocol where the transmitter persistently re-transmits the packet at consecutive slots until it is transmitted successfully or reaches the

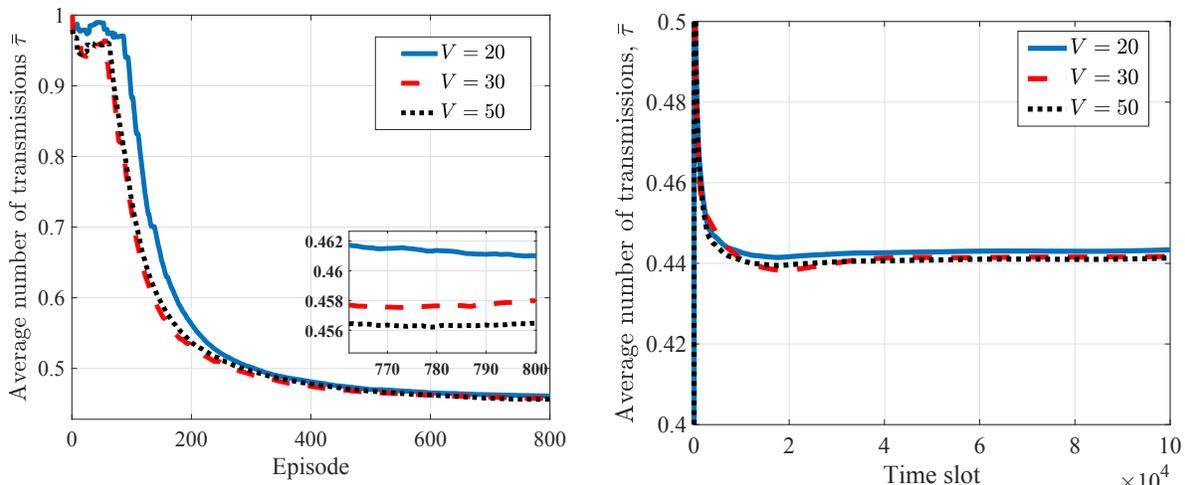
(a) The evolution of $\bar{\tau}$ during the learning process.(b) The evolution of $\bar{\tau}$ for the learned transmitter.

Figure 4: The performance of the learning-based transmission policy for different values of V , where $\eta = 0.4$, $p_0 = 0.3$, $\lambda_k = 0.7$ for all $k \in \mathcal{I}$, and $\Delta^{\max} = 4$.

maximum allowed number of transmissions x^{\max} . According to Fig. 5, as expected, the lower-bound policy outperforms the other proposed policies as it does not satisfy the constraint. Both the LC-DT policy and the deterministic transmission policy have a small gap with the lower-bound policy, which implies their near-optimal performance. The learning-based transmission policy, which is obtained without knowing the system statistics, performs relatively close to the policies for the known environment, especially for $\Delta^{\max} \geq 4$. In general, compared to the baseline policy, the proposed policies improve the system performance considerably, e.g., the LC-DT policy provides about 40 % improvement.

In Table I, we compare the time complexity of the proposed transmission policies by evaluating the consumed time 1) in offline phase, i.e., initial processing time to find a policy and 2) in online phase, i.e., running time to find the optimal action at each slot. Regarding the offline phase, both the deterministic transmission policy and the learning-based transmission policy need to explore a large number of states, thus, they have a large initial processing time. Since the deterministic transmission policy needs to search for the optimal Lagrangian multiplier (bisection) and explore all the states (RVIA), its initial processing time is the longest. On the contrary, the LC-DT policy does not involve any initial processing (besides simply setting up problem (25)). In the online phase, the running time of the deterministic transmission policy is the smallest as the action selection is done through the lookup table. The running time of the LC-DT policy is larger, because it needs to solve optimization problem (25). The learning-based transmission policy has the largest running time, where a forward pass through the neural network is executed to select

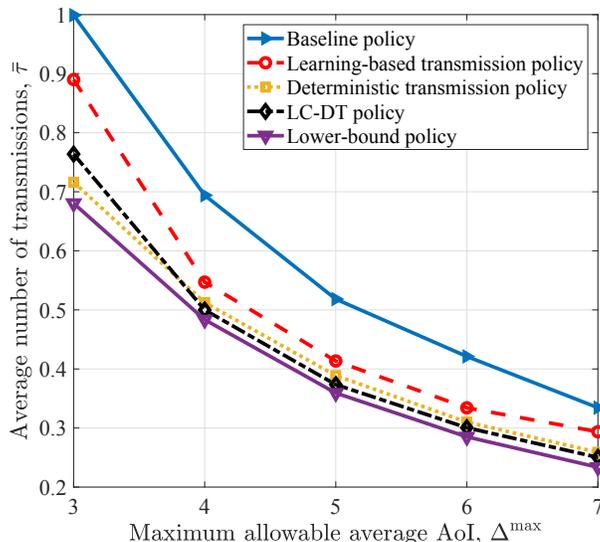


Figure 5: The average number of transmissions, $\bar{\tau}$, for the proposed transmission policies versus Δ^{\max} where $p_0 = 0.4$, $\eta = 0.4$, and $\lambda_k = 0.7$ for all $k \in \mathcal{I}$.

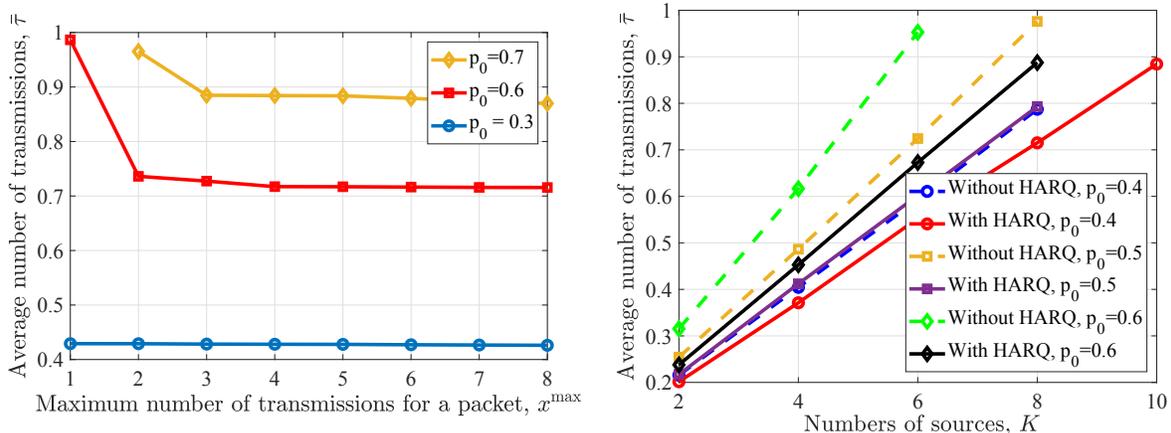
Table I: Time complexity of the proposed transmission scheduling policies

Policy	Initial processing time (s)	Running time (ms)
Deterministic transmission policy	105013.745	0.031
LC-DT policy	0	0.342
Learning-based transmission policy	4029.38	0.934

an action.

E. Impact of HARQ

In Fig. 6, we study the impact of HARQ on the system's performance by evaluating the average number of transmissions, $\bar{\tau}$. Here, without loss of generality, we utilize the LC-DT policy. Fig. 6(a) shows the effect of the maximum allowed number of transmissions for a packet, x^{\max} , where $x^{\max} = 1$ indicates that the system operates without HARQ. As it can be seen, for high probability of successful decoding (small p_0), i.e., for good channel conditions, HARQ is not that beneficial. This is due to the fact that with high probability of successful decoding, most of the packets are successfully decoded in the first transmission attempt. However, in bad channel conditions, HARQ plays an important role. For example, when $p_0 = 0.6$, the performance improves substantially by increasing x^{\max} from 1 to 2, i.e., by merely activating the HARQ with only one allowed retransmission. It is worth noting that when $p_0 = 0.7$, the transmitter cannot satisfy the average AoI constraint without HARQ ($x^{\max} = 1$). This is because when p_0 is large,



(a) The average number of transmissions, $\bar{\tau}$, versus x^{\max} , where $\eta = 0.3$ and $\Delta^{\max} = 4$.

(b) The average number of transmissions, $\bar{\tau}$, versus K , where $\eta = 0.4$, $x^{\max} = 5$, and $\Delta^{\max} = 10$.

Figure 6: The impact of employing HARQ on the performance of the system for different p_0 .

the first transmit attempts tend to fail, in which case the retransmissions would be the crucial enabler for successful receptions of the packets.

In Fig. 6(b), we study the effect of employing HARQ by evaluating $\bar{\tau}$ with respect to the number of sources K , where $K = I + J$ and $I = J$. According to the figure, $\bar{\tau}$ increases by increasing K . This is due to the fact that because the transmitter has to send each source's packets regularly to satisfy the AoI constraint, the increased number of sources inevitably leads to more transmissions in the system. Moreover, similar to Fig. 6(a), it can be seen in Fig. 6(b) that employing HARQ decreases $\bar{\tau}$, as HARQ benefits from the increased probability of successful decoding via retransmissions. For example, with $p_0 = 0.6$ and $K = 6$, $\bar{\tau}$ with HARQ is 30 % less than in the case without HARQ. Interestingly, the system cannot support $K = 10$ sources for $p_0 \geq 0.4$, unless HARQ is employed.

VI. CONCLUSION

We studied an HARQ-based multi-source status update system with random arrival and generate-at-will sources, communicating through an error-prone channel. We solved the problem of minimizing the average number of transmissions subject to the average AoI constraint in the known and unknown environments. We developed a deterministic transmission policy using the RVIA and the bisection for the known environment. For the sake of reduced computational complexity, we developed the LC-DT policy using the DPP method for the known environment. For the unknown environment, we utilized the DPP method and the DQL algorithm to develop a learning-based transmission policy. The numerical results showed the near-optimal performance

of the deterministic transmission policy and the LC-DT policy. Also, the learning-based policy attained performance relatively close to the policies developed for the known environment. Overall, the results showed about 40 % performance gain for the proposed policies over a baseline scheduling policy and demonstrated the great potential of HARQ to improve information freshness in multi-source status update systems.

APPENDIX A

PROOF OF LEMMA 1

To derive $\mathbb{E}\{\hat{\delta}_{t+1} \mid o_t\}$, we use the definition in (8), $\hat{\delta}_t = \frac{1}{K} \sum_{k=1}^K \delta_{t,k}$, and re-express $\delta_{t+1,k}$ via (4) as

$$\begin{aligned} \delta_{t+1,k} &= u_{t,k}d_t\tilde{\delta}_{t,k}^f + r_{t,k}d_t\tilde{\delta}_{t,k}^p + [u_{t,k}(1-d_t) + r_{t,k}(1-d_t) + (1-u_{t,k}-r_{t,k})]\tilde{\delta}_{t,k} \\ &= u_{t,k}d_t\tilde{\delta}_{t,k}^f + r_{t,k}d_t\tilde{\delta}_{t,k}^p + [u_{t,k} - u_{t,k}d_t + r_{t,k} - r_{t,k}d_t + 1 - u_{t,k} - r_{t,k}]\tilde{\delta}_{t,k} \\ &= u_{t,k}d_t\tilde{\delta}_{t,k}^f + r_{t,k}d_t\tilde{\delta}_{t,k}^p + [1 - d_t(u_{t,k} + r_{t,k})]\tilde{\delta}_{t,k}. \end{aligned} \quad (27)$$

Taking conditional expectation in (27), $\mathbb{E}\{\delta_{t+1,k} \mid o_t\}$ is expressed as

$$\begin{aligned} \mathbb{E}\{\delta_{t+1,k} \mid o_t\} &= \mathbb{E}\{u_{t,k}d_t\tilde{\delta}_{t,k}^f \mid o_t\} + \mathbb{E}\{r_{t,k}d_t\tilde{\delta}_{t,k}^p \mid o_t\} + \mathbb{E}\{[1 - d_t(u_{t,k} + r_{t,k})]\tilde{\delta}_{t,k} \mid o_t\} \\ &\stackrel{(a)}{=} \mathbb{E}\{u_{t,k}d_t \mid o_t\}\tilde{\delta}_{t,k}^f + \mathbb{E}\{r_{t,k}d_t \mid o_t\}\tilde{\delta}_{t,k}^p + [1 - \mathbb{E}\{d_t(u_{t,k} + r_{t,k}) \mid o_t\}]\tilde{\delta}_{t,k}, \end{aligned} \quad (28)$$

where equality (a) follows from the fact that $\delta_{t,k}^f$, $\delta_{t,k}^p$, and $\delta_{t,k}$ are given by the network state.

We need to calculate $\mathbb{E}\{u_{t,k}d_t \mid o_t\}$ and $\mathbb{E}\{r_{t,k}d_t \mid o_t\}$ in (28) which are given by the following two lemmas.

Lemma 3. *For any source k , the conditional expectation $\mathbb{E}\{u_{t,k}d_t \mid o_t\}$ is given as*

$$\mathbb{E}\{u_{t,k}d_t \mid o_t\} = \mathbb{E}\{u_{t,k} \mid o_t\}f(1). \quad (29)$$

Proof. Based on the law of iterated expectations, we have

$$\begin{aligned} \mathbb{E}\{u_{t,k}d_t \mid o_t\} &= \mathbb{E}\{\mathbb{E}\{u_{t,k}d_t \mid o_t, u_{t,k}\}\} = \mathbb{E}\{1d_t \mid o_t, u_{t,k} = 1\}\Pr(u_{t,k} = 1 \mid o_t) \\ &\quad + \mathbb{E}\{0d_t \mid o_t, u_{t,k} = 0\}\Pr(u_{t,k} = 0 \mid o_t) = \left(1f(1) + 0(1-f(1))\right)\Pr(u_{t,k} = 1 \mid o_t) \\ &= f(1)\Pr(u_{t,k} = 1 \mid o_t) \stackrel{(a)}{=} \mathbb{E}\{u_{t,k} \mid o_t\}f(1), \end{aligned}$$

where the equality (a) comes from the following equality

$$\mathbb{E}\{u_{t,k} \mid o_t\} = 0\Pr(u_{t,k} = 0 \mid o_t) + 1\Pr(u_{t,k} = 1 \mid o_t) = \Pr(u_{t,k} = 1 \mid o_t). \quad (30)$$

□

Lemma 4. *For any source k , the conditional expectation $\mathbb{E}\{r_{t,k}d_t \mid o_t\}$ is given as*

$$\mathbb{E}\{r_{t,k}d_t \mid o_t\} = \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k} + 1). \quad (31)$$

Proof. Following the same steps as in the proof of Lemma 3, we have

$$\begin{aligned}
\mathbb{E}\{r_{t,k}d_t \mid o_t\} &= \mathbb{E}\{\mathbb{E}\{r_{t,k}d_t \mid o_t, r_{t,k}\}\} = \mathbb{E}\{1d_t \mid o_t, r_{t,k} = 1\}\Pr(r_{t,k} = 1 \mid o_t) \\
&+ \mathbb{E}\{0d_t \mid o_t, r_{t,k} = 0\}\Pr(r_{t,k} = 0 \mid o_t) \\
&= \left(1f(x_{t,k} + 1) + 0(1 - f(x_{t,k} + 1))\right)\Pr(r_{t,k} = 1 \mid o_t) \\
&= f(x_{t,k} + 1)\Pr(r_{t,k} = 1 \mid o_t) \stackrel{(a)}{=} \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k} + 1),
\end{aligned}$$

where the equality (a) comes from the following equality

$$\mathbb{E}\{r_{t,k} \mid o_t\} = 0\Pr(r_{t,k} = 0 \mid o_t) + 1\Pr(r_{t,k} = 1) = \Pr(r_{t,k} = 1 \mid o_t). \quad (32)$$

□

Using Lemmas 3 and 4, the expression in (28) becomes

$$\begin{aligned}
\mathbb{E}\{\hat{\delta}_{t+1}^2 \mid o_t\} &= \frac{1}{K} \sum_{k \in \mathcal{K}} \mathbb{E}\{u_{t,k} \mid o_t\}f(1)\tilde{\delta}_{t,k}^f + \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k} + 1)\tilde{\delta}_{t,k}^p \\
&+ [1 - f(1)\mathbb{E}\{u_{t,k} \mid o_t\} - f(x_{t,k} + 1)\mathbb{E}\{r_{t,k} \mid o_t\}]\tilde{\delta}_{t,k}.
\end{aligned} \quad (33)$$

APPENDIX B

PROOF OF LEMMA 2

To derive $\mathbb{E}\{\hat{\delta}_{t+1}^2 \mid o_t\}$, we use the definition in (8), $\hat{\delta}_t = \frac{1}{K} \sum_{k=1}^K \delta_{t,k}$, and re-express it as

$$\begin{aligned}
\mathbb{E}\{\hat{\delta}_{t+1}^2 \mid o_t\} &= \frac{1}{K^2} \mathbb{E}\{(\delta_{t+1,1} + \dots + \delta_{t+1,K})^2 \mid o_t\} \\
&= \frac{1}{K^2} \left[\sum_{k \in \mathcal{K}} \mathbb{E}\{\delta_{t+1,k}^2 \mid o_t\} + \sum_{k \in \mathcal{K}} \sum_{k' \in \mathcal{K}, k' \neq k} \mathbb{E}\{\delta_{t+1,k}\delta_{t+1,k'} \mid o_t\} \right].
\end{aligned} \quad (34)$$

We need to calculate the terms $\mathbb{E}\{\delta_{t+1,k}^2 \mid o_t\}$ and $\mathbb{E}\{\delta_{t+1,k}\delta_{t+1,k'} \mid o_t\}$. As the conditions in (4) are mutually exclusive and collectively exhaustive and the square of a binary variable equals the variable itself, $\delta_{t+1,k}^2$ can be calculated from (4) and (27) as

$$\delta_{t+1,k}^2 = u_{t,k}d_t(\tilde{\delta}_{t,k}^f)^2 + r_{t,k}d_t(\tilde{\delta}_{t,k}^p)^2 + [1 - d_t(u_{t,k} + r_{t,k})](\tilde{\delta}_{t,k})^2. \quad (35)$$

Taking conditional expectation in (35), $\mathbb{E}\{\delta_{t+1,k}^2 \mid o_t\}$ is expressed as

$$\begin{aligned}
\mathbb{E}\{\delta_{t+1,k}^2 \mid o_t\} &= \mathbb{E}\{u_{t,k}d_t(\tilde{\delta}_{t,k}^f)^2 \mid o_t\} + \mathbb{E}\{r_{t,k}d_t(\tilde{\delta}_{t,k}^p)^2 \mid o_t\} \\
&+ \mathbb{E}\{[1 - d_t(u_{t,k} + r_{t,k})](\tilde{\delta}_{t,k})^2 \mid o_t\} \stackrel{(a)}{=} \mathbb{E}\{u_{t,k}d_t \mid o_t\}(\tilde{\delta}_{t,k}^f)^2 + \mathbb{E}\{r_{t,k}d_t \mid o_t\}(\tilde{\delta}_{t,k}^p)^2 \\
&+ [1 - \mathbb{E}\{d_t(u_{t,k} + r_{t,k}) \mid o_t\}](\tilde{\delta}_{t,k})^2,
\end{aligned} \quad (36)$$

where equality (a) follows from the fact that $\delta_{t,k}^f$, $\delta_{t,k}^p$, and $\delta_{t,k}$ are given by the network state.

Using Lemmas 3 and 4, the expression in (36) is calculated as

$$\begin{aligned}
\mathbb{E}\{\delta_{t+1,k}^2 \mid o_t\} &= \mathbb{E}\{u_{t,k} \mid o_t\}f(1)(\tilde{\delta}_{t,k}^f)^2 + \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k} + 1)(\tilde{\delta}_{t,k}^p)^2 \\
&+ [1 - f(1)\mathbb{E}\{u_{t,k} \mid o_t\} - f(x_{t,k} + 1)\mathbb{E}\{r_{t,k} \mid o_t\}](\tilde{\delta}_{t,k})^2.
\end{aligned} \quad (37)$$

Now, we calculate $\delta_{t+1,k}\delta_{t+1,k'}$ in expression (34). Based on (27), we can express it as

$$\begin{aligned} \delta_{t+1,k}\delta_{t+1,k'} &= u_{t,k}d_t(1-d_t(u_{t,k'}+r_{t,k'}))\tilde{\delta}_{t,k}^f\tilde{\delta}_{t,k'} + r_{t,k}d_t(1-d_t(u_{t,k'}+r_{t,k'}))\tilde{\delta}_{t,k}^p\tilde{\delta}_{t,k'} \\ &\quad + (1-d_t(u_{t,k}+r_{t,k}))(d_tu_{t,k'})\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'}^f + (1-d_t(u_{t,k}+r_{t,k}))(d_tr_{t,k'})\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'}^p \\ &\quad + (1-d_t(u_{t,k}+r_{t,k}))(1-d_t(u_{t,k'}+r_{t,k'}))\tilde{\delta}_{t,k'}\tilde{\delta}_{t,k}. \end{aligned} \quad (38)$$

Because the transmitter can transmit one packet per slot, we have $u_{t,k}u_{t,k'} = 0$ and $u_{t,k}r_{t,k'} = 0$.

Thus, the expression in (38) is rewritten as

$$\begin{aligned} \delta_{t+1,k}\delta_{t+1,k'} &= u_{t,k}d_t\tilde{\delta}_{t,k}^f\tilde{\delta}_{t,k'} + r_{t,k}d_t\tilde{\delta}_{t,k}^p\tilde{\delta}_{t,k'} + u_{t,k'}d_t\tilde{\delta}_{t,k}^f\tilde{\delta}_{t,k} + r_{t,k'}d_t\tilde{\delta}_{t,k}^p\tilde{\delta}_{t,k} \\ &\quad - u_{t,k}d_t\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'} - r_{t,k}d_t\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'} - u_{t,k'}d_t\tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} - r_{t,k'}d_t\tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} + \tilde{\delta}_{t,k'}\tilde{\delta}_{t,k}. \end{aligned} \quad (39)$$

Using Lemmas 3 and 4, the conditional expectation of the expression in (39) is calculated as

$$\begin{aligned} \mathbb{E}\{\delta_{t+1,k}\delta_{t+1,k'} \mid o_t\} &= \mathbb{E}\{u_{t,k} \mid o_t\}f(1)\tilde{\delta}_{t,k}^f\tilde{\delta}_{t,k'} + \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k}+1)\tilde{\delta}_{t,k}^p\tilde{\delta}_{t,k'} \\ &\quad + \mathbb{E}\{u_{t,k'} \mid o_t\}f(1)\tilde{\delta}_{t,k'}^f\tilde{\delta}_{t,k} + \mathbb{E}\{r_{t,k'} \mid o_t\}f(x_{t,k'}+1)\tilde{\delta}_{t,k'}^p\tilde{\delta}_{t,k} \\ &\quad - \mathbb{E}\{u_{t,k} \mid o_t\}f(1)\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'} - \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k}+1)\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'} \\ &\quad - \mathbb{E}\{u_{t,k'} \mid o_t\}f(1)\tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} - \mathbb{E}\{r_{t,k'} \mid o_t\}f(x_{t,k'}+1)\tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} + \tilde{\delta}_{t,k'}\tilde{\delta}_{t,k}. \end{aligned} \quad (40)$$

Substituting (37) and (40) into (34), we derive

$$\begin{aligned} \mathbb{E}\{\hat{\delta}_{t+1}^2 \mid o_t\} &= \frac{1}{K^2} \left[\sum_{k \in \mathcal{K}} \mathbb{E}\{u_{t,k} \mid o_t\}f(1)(\tilde{\delta}_{t,k}^f)^2 + \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k}+1)(\tilde{\delta}_{t,k}^p)^2 \right. \\ &\quad + [1 - f(1)\mathbb{E}\{u_{t,k} \mid o_t\} - f(x_{t,k}+1)\mathbb{E}\{r_{t,k} \mid o_t\}](\tilde{\delta}_{t,k})^2 \\ &\quad + \sum_{k \in \mathcal{K}} \sum_{k' \in \mathcal{K}, k' \neq k} \mathbb{E}\{u_{t,k} \mid o_t\}f(1)\tilde{\delta}_{t,k}^f\tilde{\delta}_{t,k'} + \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k}+1)\tilde{\delta}_{t,k}^p\tilde{\delta}_{t,k'} \\ &\quad + \mathbb{E}\{u_{t,k'} \mid o_t\}f(1)\tilde{\delta}_{t,k'}^f\tilde{\delta}_{t,k} + \mathbb{E}\{r_{t,k'} \mid o_t\}f(x_{t,k'}+1)\tilde{\delta}_{t,k'}^p\tilde{\delta}_{t,k} \\ &\quad - \mathbb{E}\{u_{t,k} \mid o_t\}f(1)\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'} - \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k}+1)\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'} \\ &\quad \left. - \mathbb{E}\{u_{t,k'} \mid o_t\}f(1)\tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} - \mathbb{E}\{r_{t,k'} \mid o_t\}f(x_{t,k'}+1)\tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} + \tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} \right]. \end{aligned}$$

APPENDIX C

THE EXPRESSION FOR THE OBJECTIVE FUNCTION OF PROBLEM (25)

To derive W_t , we rewrite (22) using Lemmas 1 and 2 as

$$\begin{aligned} V\mathbb{E}\{\tau_t \mid o_t\} + \alpha(o_t) &\leq V \sum_{k \in \mathcal{K}} \mathbb{E}\{u_{t,k} \mid o_t\} + \mathbb{E}\{r_{t,k} \mid o_t\} + \frac{1}{2}((\Delta^{\max})^2 \\ &\quad + \frac{1}{K^2} \left[\sum_{k \in \mathcal{K}} \mathbb{E}\{u_{t,k} \mid o_t\}f(1)(\tilde{\delta}_{t,k}^f)^2 + \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k}+1)(\tilde{\delta}_{t,k}^p)^2 \right. \\ &\quad + [1 - f(1)\mathbb{E}\{u_{t,k} \mid o_t\} - f(x_{t,k}+1)\mathbb{E}\{r_{t,k} \mid o_t\}](\tilde{\delta}_{t,k})^2 \\ &\quad + \sum_{k \in \mathcal{K}} \sum_{k' \in \mathcal{K}, k' \neq k} \mathbb{E}\{u_{t,k} \mid o_t\}f(1)\tilde{\delta}_{t,k}^f\tilde{\delta}_{t,k'} + \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k}+1)\tilde{\delta}_{t,k}^p\tilde{\delta}_{t,k'} \\ &\quad + \mathbb{E}\{u_{t,k'} \mid o_t\}f(1)\tilde{\delta}_{t,k'}^f\tilde{\delta}_{t,k} + \mathbb{E}\{r_{t,k'} \mid o_t\}f(x_{t,k'}+1)\tilde{\delta}_{t,k'}^p\tilde{\delta}_{t,k} - \mathbb{E}\{u_{t,k} \mid o_t\}f(1)\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'} \\ &\quad - \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k}+1)\tilde{\delta}_{t,k}\tilde{\delta}_{t,k'} - \mathbb{E}\{u_{t,k'} \mid o_t\}f(1)\tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} \\ &\quad \left. - \mathbb{E}\{r_{t,k'} \mid o_t\}f(x_{t,k'}+1)\tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} + \tilde{\delta}_{t,k'}\tilde{\delta}_{t,k} \right] + 2Q_t \left(\frac{1}{K} \sum_{k \in \mathcal{K}} \mathbb{E}\{u_{t,k} \mid o_t\}f(1)\tilde{\delta}_{t,k}^f \right. \\ &\quad + \mathbb{E}\{r_{t,k} \mid o_t\}f(x_{t,k}+1)\tilde{\delta}_{t,k}^p \\ &\quad \left. + [1 - f(1)\mathbb{E}\{u_{t,k} \mid o_t\} - f(x_{t,k}+1)\mathbb{E}\{r_{t,k} \mid o_t\}]\tilde{\delta}_{t,k} - \Delta^{\max} \right). \end{aligned} \quad (41)$$

Dropping the expectation in (41), W_t is derived as

$$\begin{aligned}
W_t = & V \sum_{k \in \mathcal{K}} u_{t,k} + r_{t,k} + \frac{1}{2K^2} \left[\sum_{k \in \mathcal{K}} u_{t,k} f(1) (\tilde{\delta}_{t,k}^f)^2 + r_{t,k} f(x_{t,k} + 1) (\tilde{\delta}_{t,k}^p)^2 \right. \\
& + [1 - u_{t,k} f(1) - r_{t,k} f(x_{t,k} + 1)] (\tilde{\delta}_{t,k})^2 + \sum_{k \in \mathcal{K}} \sum_{k' \in \mathcal{K}, k' \neq k} u_{t,k} f(1) \tilde{\delta}_{t,k}^f \tilde{\delta}_{t,k'} \\
& + r_{t,k} f(x_{t,k} + 1) \tilde{\delta}_{t,k}^p \tilde{\delta}_{t,k'} - u_{t,k} f(1) \tilde{\delta}_{t,k} \tilde{\delta}_{t,k'} - r_{t,k} f(x_{t,k} + 1) \tilde{\delta}_{t,k} \tilde{\delta}_{t,k'} \\
& + u_{t,k'} f(1) \tilde{\delta}_{t,k'}^f \tilde{\delta}_{t,k} + r_{t,k'} f(x_{t,k'} + 1) \tilde{\delta}_{t,k'}^p \tilde{\delta}_{t,k} - u_{t,k'} f(1) \tilde{\delta}_{t,k'} \tilde{\delta}_{t,k} \\
& - r_{t,k'} f(x_{t,k'} + 1) \tilde{\delta}_{t,k'} \tilde{\delta}_{t,k} + \tilde{\delta}_{t,k'} \tilde{\delta}_{t,k} + 2KQ_t \left(\sum_{k \in \mathcal{K}} u_{t,k} f(1) \tilde{\delta}_{t,k}^f + r_{t,k} f(x_{t,k} + 1) \tilde{\delta}_{t,k}^p \right. \\
& \left. + [1 - u_{t,k} f(1) - r_{t,k} f(x_{t,k} + 1)] \tilde{\delta}_{t,k} \right) + \frac{1}{2} ((\Delta^{\max})^2 - 2Q_t \Delta^{\max}).
\end{aligned}$$

REFERENCES

- [1] S. Kaul, R. Yates, and M. Gruteser, "Real-time status: How often should one update?" in *Proc. IEEE Int. Conf. on Computer Commun. (INFOCOM)*, Orlando, FL, USA, Mar. 25–30, 2012, pp. 2731–2735.
- [2] M. A. Abd-Elmagid, N. Pappas, and H. S. Dhillon, "On the role of age of information in the internet of things," *IEEE Commun. Mag.*, vol. 57, no. 12, pp. 72–77, Dec. 2019.
- [3] S. K. Kaul, R. D. Yates, and M. Gruteser, "Status updates through queues," in *Proc. Conf. Inform. Sciences Syst. (CISS)*, Princeton, NJ, USA, Mar.21–23, 2012, pp. 1–6.
- [4] R. D. Yates, "The age of information in networks: Moments, distributions, and sampling," *IEEE Trans. Inform. Theory*, vol. 66, no. 9, pp. 5712–5728, Sep. 2020.
- [5] A. Kosta, N. Pappas, and V. Angelakis, "Age of information: A new concept, metric, and tool," *Found. Trends Netw.*, vol. 12, no. 3, pp. 162–259, Nov. 2017.
- [6] S. Lin, D. J. Costello, and M. J. Miller, "Automatic-repeat-request error-control schemes," *IEEE Commun. Mag.*, vol. 22, no. 12, pp. 5–17, Dec. 1984.
- [7] "IEEE standard for air interface for broadband wireless access systems," *IEEE Std 802.16-2017 (Revision of IEEE Std 802.16-2012)*, pp. 1–2726, Mar. 2018.
- [8] V. Raghunathan, C. Schurgers, S. Park, and M. Srivastava, "Energy-aware wireless microsensor networks," *IEEE Signal Processing Mag.*, vol. 19, no. 2, pp. 40–50, Mar. 2002.
- [9] S. Russell and P. Norvig, *Artificial intelligence: a modern approach*, 2002.
- [10] 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, Nov. 1985.
- [11] B. Zhou and W. Saad, "Joint status sampling and updating for minimizing age of information in the internet of things," *IEEE Trans. Commun.*, vol. 67, no. 11, pp. 7468–7482, Nov. 2019.
- [12] E. T. Ceran, D. Gündüz, and A. György, "A reinforcement learning approach to age of information in multi-user networks with HARQ," *IEEE J. Select. Areas Commun.*, vol. 39, no. 5, pp. 1412–1426, May 2021.
- [13] E. T. Ceran, D. Gündüz, and A. György, "Average age of information with hybrid ARQ under a resource constraint," *IEEE Trans. Wireless Commun.*, vol. 18, no. 3, pp. 1900–1913, Mar. 2019.
- [14] M. J. Neely, *Stochastic network optimization with application to communication and queueing systems*. Morgan & Claypool Publishers, 2010.
- [15] V. Mnih *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
- [16] R. D. Yates and S. K. Kaul, "The age of information: Real-time status updating by multiple sources," *IEEE Trans. Inform. Theory*, vol. 65, no. 3, pp. 1807–1827, Mar. 2018.
- [17] M. Moltafet, M. Leinonen, and M. Codreanu, "Average AoI in multi-source systems with source-aware packet management," *IEEE Trans. Commun.*, vol. 69, no. 2, pp. 1121–1133, Feb. 2021.

- [18] E. Najm, R. Yates, and E. Soljanin, "Status updates through M/G/1/1 queues with HARQ," in *Proc. IEEE Int. Symp. Inform. Theory*, Aachen, Germany, Jun. 25–30, 2017, pp. 131–135.
- [19] D. Li, S. Wu, Y. Wang, J. Jiao, and Q. Zhang, "Age-optimal HARQ design for freshness-critical satellite-IoT systems," *IEEE Internet of Things J.*, vol. 7, no. 3, pp. 2066–2076, Mar. 2020.
- [20] S. Liu, J. Jiao, Z. Ni, S. Wu, and Q. Zhang, "Age-optimal NC-HARQ protocol for multi-hop satellite-based internet of things," in *Proc. IEEE Wireless Commun. and Networking Conf.*, Nanjing, China, Mar. 29–1, 2021, pp. 1–6.
- [21] S. Farazi, A. G. Klein, and D. R. Brown, "Average age of information in update systems with active sources and packet delivery errors," *IEEE Wireless Commun. Lett.*, vol. 9, no. 8, pp. 1164–1168, Aug. 2020.
- [22] Y. Shi, L. Jing, X. Jia, P. Ji, and N. Wan, "Improvement on age of information for information update systems with HARQ chase combining and sensor harvesting-transmitting diversities," *IEEE Access*, vol. 9, pp. 78 035–78 049, May 2021.
- [23] A. Arafa, J. Yang, S. Ulukus, and H. V. Poor, "Timely status updating over erasure channels using an energy harvesting sensor: Single and multiple sources," *IEEE Trans. Green Commun. Net.*, vol. 6, no. 1, pp. 6–19, Mar. 2022.
- [24] Z. Deng, S. Wu, C. Guo, J. Jiao, N. Zhang, and Q. Zhang, "Age-optimal transmission policy for intelligent HARQ-CC aided NOMA systems," in *Proc. IEEE Int. Conf. Commun.*, Montreal, QC, Canada, Jun. 14–23, 2021, pp. 1–6.
- [25] S. Feng and J. Yang, "Age of information minimization for an energy harvesting source with updating erasures: Without and with feedback," *IEEE Trans. Commun.*, vol. 69, no. 8, pp. 5091–5105, May 2021.
- [26] Y. Wang, S. Wu, J. Jiao, W. Wu, Y. Wang, and Q. Zhang, "Age-optimal transmission policy with HARQ for freshness-critical vehicular status updates in space-air-ground integrated networks," *IEEE Internet of Things J.*, vol. 9, no. 8, pp. 5719–5729, Apr. 2022.
- [27] X. Lagrange, "Throughput of HARQ protocols on a block fading channel," *IEEE Commun. Lett.*, vol. 14, no. 3, pp. 257–259, Mar. 2010.
- [28] P. Frenger, S. Parkvall, and E. Dahlman, "Performance comparison of HARQ with chase combining and incremental redundancy for HSDPA," in *Proc. IEEE Veh. Technol. Conf.*, Atlantic City, NJ, USA, Oct. 2001, pp. 1829–1833.
- [29] M. Moltafet, M. Leinonen, M. Codreanu, and N. Pappas, "Power minimization for age of information constrained dynamic control in wireless sensor networks," *IEEE Trans. Commun.*, vol. 70, no. 1, pp. 419–432, Jan. 2022.
- [30] A. Maatouk, S. Kriouile, M. Assad, and A. Ephremides, "On the optimality of the Whittle's index policy for minimizing the age of information," *IEEE Trans. Wireless Commun.*, vol. 20, no. 2, pp. 1263–1277, Oct. 2021.
- [31] G. Yao, A. Bedewy, and N. B. Shroff, "Age-optimal low-power status update over time-correlated fading channel," *IEEE Trans. Mobile Comput.*, pp. 1–1, Mar. 2022.
- [32] 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 Processing*, vol. 55, no. 10, pp. 5069–5083, Oct. 2007.
- [33] E. Altman, *Constrained Markov decision processes*. CRC Press, 1999.
- [34] D.-j. Ma, A. M. Makowski, and A. Shwartz, "Estimation and optimal control for constrained Markov chains," in *1986 25th IEEE Conf. on Decision and Contr.*, Athens, Greece, Dec. 1986, pp. 994–999.
- [35] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., 1994.
- [36] W. Wu, P. Yang, W. Zhang, C. Zhou, and X. Shen, "Accuracy-guaranteed collaborative DNN inference in industrial IoT via deep reinforcement learning," *IEEE Trans. Ind. Informat.*, vol. 17, no. 7, pp. 4988–4998, Aug. 2020.