

# Optimizing Data Freshness, Throughput, and Delay in Multi-Server Information-Update Systems

Ahmed M. Bedewy<sup>†</sup>, Yin Sun<sup>†</sup>, and Ness B. Shroff<sup>† ‡</sup>

<sup>†</sup>Dept. of ECE, <sup>‡</sup>Dept. of CSE, The Ohio State University, Columbus, OH.

**Abstract**—In this work, we investigate the design of information-update systems, where incoming update packets are forwarded to a remote destination through multiple servers (each server can be viewed as a wireless channel). One important performance metric of these systems is the *data freshness at the destination*, also called the *age-of-information* or simply *age*, which is defined as the time elapsed since the freshest packet at the destination was generated. Recent studies on information-update systems have shown that the age-of-information can be reduced by intelligently dropping stale packets. However, packet dropping may not be appropriate in many applications, such as news and social updates, where users are interested in not just the latest updates, but also past news. Therefore, all packets may need to be successfully delivered. In this paper, we study how to optimize age-of-information without throughput loss. We consider a general scenario where incoming update packets do not necessarily arrive in the order of their generation times. We prove that a preemptive Last Generated First Served (LGFS) policy simultaneously optimizes the age, throughput, and delay performance in infinite buffer queueing systems. We also show age-optimality for the LGFS policy for any finite queue size. These results hold for arbitrary, including non-stationary, arrival processes.

## I. INTRODUCTION

The ubiquity of mobile devices and applications, has increased the demand for real-time information updates, such as news, weather reports, email notifications, stock quotes, social updates, mobile ads, etc. Also, in network-based monitoring and control systems, timely status updates are crucial. These include, but are not limited to, sensor networks used in temperature or other physical phenomenon, and autonomous vehicle systems.

A common objective in these applications is to keep the destination updated with the latest information. To identify the timeliness of the updates, a metric called *data freshness*, also called *age of information*, or simply *age*, was defined in [1]–[4]. At time  $t$ , if  $U(t)$  is the time when the freshest update at the destination was generated, age  $\Delta(t)$  is  $\Delta(t) = t - U(t)$ . Hence, age is the time elapsed since the freshest packet was generated.

There have been several recent works on characterizing the time-average age of different information-update policies under Poisson arrival process, and finding policies with a small time-average age [4]–[10]. In [4]–[6], the update generation rate was optimized to improve data freshness in First-Come First-Served (FCFS) information-update systems. To improve the age, these studies also reduced the update generation rate, which in turn sacrificed the system throughput. In [7], [8], it was found that the age can be improved by discarding old packets waiting in the queue if a new sample arrives. This can greatly reduce the impact of queueing delay on data freshness. However, many applications may not want to discard packets,

This work has been partially supported by grants from the National Foundation CNS-1421576, Army Research Office W911NF-14-1-0368, and Office of Naval Research N00014-15-1-2166.

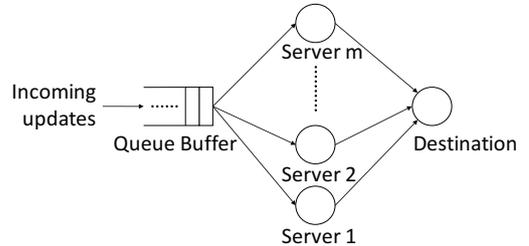


Figure 1: System model.

e.g., where the users are interested in not just the latest updates, but also past news, in which case all packets must be successfully delivered. In [9], [10], the time-average age was characterized for Last-Come First-Served (LCFS) information-update systems with and without preemption; and FCFS with two servers under a Poisson arrival process. Applications of information updates in channel information feedback and sensor networks were considered in [11]–[13].

Another important problem is how to maximize data freshness in information-update systems. This involves jointly controlling both the generation and transmission of packet updates [12]–[14]. An information update policy was developed in [14], which was proven to minimize the time-average age and time-average age penalty among all causally feasible policies. In this setting, a counter-intuitive phenomenon was revealed: While a zero-wait or work-conserving policy, that generates and submits a fresh update once the server becomes idle, achieves the maximum throughput and the minimum average delay, surprisingly, this zero-wait policy does not always minimize the age. This implies that there is no policy that can simultaneously minimize age and maximize throughput, if the generation and transmission of update packets are jointly controlled.

In this paper, we consider an information-update system which enqueues incoming update packets and forwards them to a remote destination through multiple servers, as shown in Fig. 1. In this setting, the updates are generated exogenously to the system, which is different from [14]. We aim to answer the following questions: How to establish age-optimality in a general policy space and under arbitrary arrival process? Is it possible to simultaneously optimize multiple performance metrics, such as age, throughput, and delay? To that end, the following are the key contributions of this paper:

- We consider a general scenario where the update packets do not necessarily arrive in the order of their generation times, which has not been considered before. We prove that, if the packet service times are *i.i.d.* exponentially distributed, then for an arbitrary arrival process and any queue size, a preemptive Last-Generated First-Served (LGFS) policy achieves an age process that is stochastically smaller than any causally feasible policies

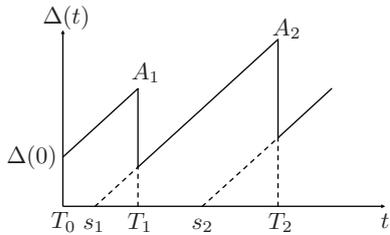


Figure 2: Evolution of the age-of-information  $\Delta(t)$ .

(Theorem 1). This implies that the preemptive LGFS policy minimizes many data freshness metrics, including time-average age, average peak age, and time-average age penalty (Corollary 6). The intuition is that the freshest update packets are served as early as possible in the preemptive LGFS policy. In particular, the distribution of the age process of the preemptive LGFS policy is invariant over all queue sizes.

- In addition, we show that if the buffer has an infinite size, then the preemptive LGFS policy is also throughput-optimal and delay-optimal among all causally feasible policies (Theorem 7). To the best of our knowledge, this is the first study which simultaneously optimizes data freshness, throughput, and delay in information-update systems.

We note that when the incoming update packets are arriving in the same order of their generation times, the proposed LGFS policy is identical to the LCFS policy studied in [9]. In particular, the time-average age of preemptive and non-preemptive LCFS policies are analyzed in [9] for single-server queueing systems with Poisson arrival process and a queue size of one packet. This paper complements and generalizes the results in [9] by (i) allowing the incoming updates to not arrive in the order of their generation times, (ii) considering more general multi-server queueing systems with arbitrary update arrivals and arbitrary queue size, and (iii) providing an age-optimality proof.

## II. SYSTEM MODEL

We consider an information-update system depicted in Fig. 1, where an infinite sequence of update packets are sent to a destination through  $m$  identical servers. Each server could be a wireless channel, a TCP connection, etc. The update packets are generated at time  $s_1, s_2, \dots$ , where  $s_1 = 0 \leq s_2 \leq \dots$  are arbitrarily given. These packets do not need to arrive to the system in the order that they were generated. Let  $a_i$  be the arrival time of the packet generated at time  $s_i$ , such that  $s_i \leq a_i$ . As the packets may arrive out of order, it may happen that  $a_i > a_{i+1}$ . Hence, each packet is identified by the pair  $(s_i, a_i)$ , which will be used to make scheduling decision. After arrival, each update packet is stored in a queue, waiting to be assigned to one of the servers. Let  $B$  denote the buffer size of the queue which can be infinite, finite, or even zero. If  $B$  is finite, the queue buffer may overflow and some packets are dropped, which would incur a throughput loss. The packet service times are exponentially distributed with rate  $\mu$ , which are *i.i.d.* across time and servers. Let  $T_i$  denote the  $i$ -th packet delivery time instant.

The system starts to operate at time  $t = 0$ . Define  $U(t)$  as the time when the freshest update at the destination at time  $t$

was generated, where we set  $U(t) = 0$  at  $t = 0$ . The *age-of-information*, or simply the *age*, is defined as

$$\Delta(t) = t - U(t). \quad (1)$$

From this definition, we can deduce that the age increases linearly with  $t$  but is reset to a smaller value with the delivery of each update that contains fresher information, as shown in Fig. 2. Define  $A_k$  as the  $k$ -th peak value of  $\Delta(t)$  since time  $t = 0$ .

Following [14], we define an age penalty function  $g(\Delta(t))$  to represent the level of “dissatisfaction” for data staleness or the “need” for new information update, where  $g(\cdot)$  is a general non-decreasing function, that is determined based on specific applications.

### A. Scheduling Policy

A scheduling policy, denoted by  $\pi$ , assigns update packets to the servers over time. We first define several classes of policies:

**Definition 1. Service Preemption:** In a *preemptive* policy, a server can switch to send any packet at any time; the preempted packets will be stored back into the queue if there is enough buffer space and sent at a later time when the servers are available again. In contrast, in a *non-preemptive* policy, a server must complete delivering the current packet before starting to send another packet.

**Definition 2. Work Conservation:** A policy is said to be *work-conserving*, if no server is idle when there are packets waiting in the queue.

We use  $\Pi$  to denote the set of all causal policies. Let  $\Pi_{wc} \subseteq \Pi$  be the set of work-conserving feasible policies, and  $\Pi_{nwc} \subseteq \Pi$  be the set of non-work-conserving feasible policies.

### B. Age Optimality

We will need the following definitions: Let  $\vec{x} = (x_1, x_2, \dots, x_n)$  and  $\vec{y} = (y_1, y_2, \dots, y_n)$  be two vectors in  $\mathbb{R}^n$ , then we denote  $\vec{x} \leq \vec{y}$  if  $x_i \leq y_i$  for  $i = 1, 2, \dots, n$ .

**Definition 3. Stochastic Ordering:** [15] Let  $X$  and  $Y$  be two random variables. Then,  $X$  is said to be stochastically smaller than  $Y$  (denoted as  $X \leq_{st} Y$ ), if

$$P\{X > x\} \leq P\{Y > x\}, \quad \forall x \in \mathbb{R}.$$

**Definition 4. Multivariate Stochastic Ordering:** [15] A set  $U \subseteq \mathbb{R}^n$  is called upper if  $\vec{y} \in U$  whenever  $\vec{y} \geq \vec{x}$  and  $\vec{x} \in U$ . Let  $\vec{X}$  and  $\vec{Y}$  be two random vectors. Then,  $\vec{X}$  is said to be stochastically smaller than  $\vec{Y}$  (denoted as  $\vec{X} \leq_{st} \vec{Y}$ ), if

$$P\{\vec{X} \in U\} \leq P\{\vec{Y} \in U\}, \quad \text{for all upper sets } U \subseteq \mathbb{R}^n.$$

**Definition 5. Stochastic Ordering of Stochastic Processes:** [15] Let  $\{X(t), t \in [0, \infty)\}$  and  $\{Y(t), t \in [0, \infty)\}$  be two stochastic processes. Then,  $\{X(t), t \in [0, \infty)\}$  is said to be stochastically smaller than  $\{Y(t), t \in [0, \infty)\}$  (denoted by  $\{X(t), t \in [0, \infty)\} \leq_{st} \{Y(t), t \in [0, \infty)\}$ ), if, for all choices of an integer  $n$  and  $t_1 < t_2 < \dots < t_n$  in  $[0, \infty)$ , it holds that

$$(X(t_1), X(t_2), \dots, X(t_n)) \leq_{st} (Y(t_1), Y(t_2), \dots, Y(t_n)), \quad (2)$$

where the multivariate stochastic ordering in (2) was defined in Definition 4.

**Definition 6. Age Optimality:** A policy  $P \in \Pi$  is said to be *age-optimal*, if for all  $\pi \in \Pi$

$$\{\Delta_P(t), t \in [0, \infty)\} \leq_{\text{st}} \{\Delta_\pi(t), t \in [0, \infty)\}. \quad (3)$$

As we will see in Corollary 6, the results obtained using this definition of age optimality is quite strong.

### III. OPTIMALITY ANALYSIS

---

**Algorithm 1:** Preemptive Last Generated First Served policy.

---

```

1  $\alpha := 0;$ 
2 while the system is ON do
3   if a new packet with time stamp  $s$  arrives then
4     if  $s \leq \alpha$  then // The packet is outdated.
5       Store the packet in the queue;
6     else // The packet carries fresh information.
7       if all servers are busy then
8         The new packet is assigned to a server by
9         preempting the packet with time stamp  $\alpha$ ;
10        The preempted packet with time stamp  $\alpha$ 
11        is stored back to the queue;
12        Set  $\alpha$  as the smallest time stamp of the
13        packets under service;
14      else // At least one of the servers is idle.
15        Assign the new packet to one idle server;
16      end
17    end
18  if a packet is delivered then
19    if the queue is not empty then
20      Pick any packet in the queue and assign it to
21      the idle server;
22    end
23  end

```

---

In this section, we study a LGFS policy, in which the packets under service are the freshest (i.e., are generated the latest) among all packets in the queue. The implementation details of a preemptive LGFS policy is depicted in Algorithm 1, where  $\alpha$  is the smallest time stamp of the packets under service. We next show that the preemptive LGFS policy is age-optimal.

**Theorem 1.** If the packet service times are exponentially distributed and *i.i.d.* across time and servers, then for all packet generation and arrival times, and all queue sizes  $B$ , the preemptive LGFS policy is age-optimal.

We need to define the system state of any policy  $\pi$ :

**Definition 7.** At any time  $t$ , the system state of policy  $\pi$  is specified by  $V_\pi(t) = (U_\pi(t), \alpha_{1,\pi}(t), \dots, \alpha_{m,\pi}(t))$ , where  $U_\pi(t)$  is the largest time stamp of the packets that have already been delivered to the destination. Define  $\alpha_{i,\pi}(t)$  as the  $i$ -th smallest time stamp of the packets being processed by the servers. Without loss of generality, if  $k$  servers are sending stale packets (i.e.,  $\alpha_{1,\pi}(t) \leq \dots \leq \alpha_{k,\pi}(t) \leq U_\pi(t)$ ) or  $k$  servers are idle, then we set  $\alpha_{1,\pi}(t) = \dots = \alpha_{k,\pi}(t) = U_\pi(t)$ . Hence,

$$U_\pi(t) \leq \alpha_{1,\pi}(t) \leq \dots \leq \alpha_{m,\pi}(t). \quad (4)$$

Let  $\{V_\pi(t), t \in [0, \infty)\}$  be the state process of policy  $\pi$ , which is assumed to be right-continuous.

For notational simplicity, let policy  $P$  represent the preemptive LGFS policy. Then, by the construction of policy  $P$ ,  $\alpha_{1,P}(t), \alpha_{2,P}(t), \dots, \alpha_{m,P}(t)$  are the time stamps of  $m$  freshest packets among all packets arrived during  $[0, t]$ .

The key step in the proof of Theorem 1 is the following lemma, where we compare policy  $P$  with any work-conserving policy  $\pi \in \Pi_{wc}$ .

**Lemma 2.** For all packet generation and arrival times, and all queue sizes  $B$ , suppose that  $V_P(0^-) = V_\pi(0^-)$  for all work-conserving policies  $\pi \in \Pi_{wc}$ , then

$$\{V_P(t), t \in [0, \infty)\} \geq_{\text{st}} \{V_\pi(t), t \in [0, \infty)\}. \quad (5)$$

We use coupling and forward induction to prove Lemma 2. For any work-conserving policy  $\pi$ , suppose that stochastic processes  $\hat{V}_P(t)$  and  $\hat{V}_\pi(t)$  have the same stochastic laws as  $V_P(t)$  and  $V_\pi(t)$ . The packet deliveries of the state processes  $\hat{V}_P(t)$  and  $\hat{V}_\pi(t)$  are coupled in the following manner: For each  $i = 1, \dots, m$ , if the packet with time stamp  $\hat{\alpha}_{i,P}(t)$  is delivered at time  $t$  as  $\hat{V}_P(t)$  evolves, then the packet with time stamp  $\hat{\alpha}_{i,\pi}(t)$  is delivered at time  $t$  as  $\hat{V}_\pi(t)$  evolves. Such a coupling is valid since the service time is exponentially distributed and thus memoryless. On the other hand, policy  $P$  and policy  $\pi$  have identical arrival processes. According to Theorem 6.B.30 in [15], if we can show

$$P \left[ \hat{V}_P(t) \geq \hat{V}_\pi(t), t \in [0, \infty) \right] = 1, \quad (6)$$

then (5) is proven. Next, we use the following two lemmas to prove (6):

**Lemma 3.** Suppose that under policy  $P$ ,  $\{\hat{U}'_P, \hat{\alpha}'_{1,P}, \dots, \hat{\alpha}'_{m,P}\}$  is obtained by delivering a packet with time stamp  $\hat{\alpha}_{l,P}$  to the destination in the system whose state is  $\{\hat{U}_P, \hat{\alpha}_{1,P}, \dots, \hat{\alpha}_{m,P}\}$ . Further, suppose that under policy  $\pi$ ,  $\{\hat{U}'_\pi, \hat{\alpha}'_{1,\pi}, \dots, \hat{\alpha}'_{m,\pi}\}$  is obtained by delivering a packet with time stamp  $\hat{\alpha}_{l,\pi}$  to the destination in the system whose state is  $\{\hat{U}_\pi, \hat{\alpha}_{1,\pi}, \dots, \hat{\alpha}_{m,\pi}\}$ . If

$$\hat{U}_P \geq \hat{U}_\pi, \hat{\alpha}_{i,P} \geq \hat{\alpha}_{i,\pi}, \quad \forall i = 1, \dots, m, \quad (7)$$

then,

$$\hat{U}'_P \geq \hat{U}'_\pi, \hat{\alpha}'_{i,P} \geq \hat{\alpha}'_{i,\pi}, \quad \forall i = 1, \dots, m. \quad (8)$$

*Proof.* See Appendix A.  $\square$

**Lemma 4.** Suppose that under policy  $P$ ,  $\{\hat{U}'_P, \hat{\alpha}'_{1,P}, \dots, \hat{\alpha}'_{m,P}\}$  is obtained by adding a packet with time stamp  $s$  to the system whose state is  $\{\hat{U}_P, \hat{\alpha}_{1,P}, \dots, \hat{\alpha}_{m,P}\}$ . Further, suppose that under policy  $\pi$ ,  $\{\hat{U}'_\pi, \hat{\alpha}'_{1,\pi}, \dots, \hat{\alpha}'_{m,\pi}\}$  is obtained by adding a packet with time stamp  $s$  to the system whose state is  $\{\hat{U}_\pi, \hat{\alpha}_{1,\pi}, \dots, \hat{\alpha}_{m,\pi}\}$ . If

$$\hat{U}_P \geq \hat{U}_\pi, \hat{\alpha}_{i,P} \geq \hat{\alpha}_{i,\pi}, \quad \forall i = 1, \dots, m, \quad (9)$$

then

$$\hat{U}'_P \geq \hat{U}'_\pi, \hat{\alpha}'_{i,P} \geq \hat{\alpha}'_{i,\pi}, \quad \forall i = 1, \dots, m. \quad (10)$$

*Proof.* See our technical report [16].  $\square$

*Proof of Lemma 2.* For any sample path, we have that  $\hat{U}_P(0^-) = \hat{U}_\pi(0^-)$  and  $\hat{\alpha}_{i,P}(0^-) = \hat{\alpha}_{i,\pi}(0^-)$  for  $i = 1, \dots, m$ . This, together with Lemma 3 and 4, implies that

$$\hat{U}_P(t) \geq \hat{U}_\pi(t), \hat{\alpha}_{i,P}(t) \geq \hat{\alpha}_{i,\pi}(t),$$

holds for all  $t \in [0, \infty)$  and  $i = 1, \dots, m$ . Hence, (6) follows.

Because  $\{\hat{V}_P(t), t \in [0, \infty)\}$  and  $\{V_P(t), t \in [0, \infty)\}$  are of the same distribution,  $\{\hat{V}_\pi(t), t \in [0, \infty)\}$  and  $\{V_\pi(t), t \in [0, \infty)\}$  are of the same distribution, by Theorem 6.B.30 in [15] we get (5). This completes the proof.  $\square$

*Proof of Theorem 1.* As a result of Lemma 2, we have

$$\{U_P(t), t \in [0, \infty)\} \geq_{st} \{U_\pi(t), t \in [0, \infty)\}, \forall \pi \in \Pi_{wc},$$

which implies

$$\{\Delta_P(t), t \in [0, \infty)\} \leq_{st} \{\Delta_\pi(t), t \in [0, \infty)\}, \forall \pi \in \Pi_{wc}.$$

Finally, since the service times are *i.i.d.* across time and servers, service idling only increases the waiting time of the packet in the system. Therefore, the age under non-work-conserving policies will be greater. As a result, we have

$$\{\Delta_P(t), t \in [0, \infty)\} \leq_{st} \{\Delta_\pi(t), t \in [0, \infty)\}, \forall \pi \in \Pi.$$

This completes the proof.  $\square$

As a result of Theorem 1, we can deduce the following corollaries:

**Corollary 5.** If the packet service times are exponentially distributed and *i.i.d.* across time and servers, then for all packet generation and arrival times, the age performance of the preemptive LGFS policy remains the same for any queue size  $B \geq 0$ .

We can notice from the proof of Lemma 3 and Lemma 4 that, the system states evolution is invariant for any queue size  $B \geq 0$ . Hence, the age performance of the preemptive LGFS policy is invariant for any queue size  $B \geq 0$ .

**Corollary 6.** If the packet service times are *i.i.d.* exponentially distributed across time and servers, then for all packet generation and arrival times, all queue sizes  $B$ ,  $\pi \in \Pi$ , and non-decreasing function  $g$ ,

$$\begin{aligned} \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \Delta_{\text{prmp-LGFS}}(t) dt &\leq_{st} \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T \Delta_\pi(t) dt, \\ \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K A_{k, \text{prmp-LGFS}} &\leq_{st} \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{k=1}^K A_{k, \pi}, \\ \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T g(\Delta_{\text{prmp-LGFS}}(t)) dt &\leq_{st} \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T g(\Delta_\pi(t)) dt, \end{aligned}$$

where the limits are assumed to exist.<sup>1</sup>

Hence, the preemptive LGFS policy minimizes time-average age, average peak age, and time-average age penalty for any non-decreasing penalty function  $g$ . Finally, the delay and throughput optimality of the preemptive LGFS policy is stated as follows:

**Theorem 7.** If the packet service times are *i.i.d.* exponentially distributed across time and servers, then for all packet generation and arrival times, and infinite buffer size  $B = \infty$ , the preemptive LGFS policy is throughput-optimal and mean-delay-optimal among all policies in  $\Pi$ .

*Proof.* See our technical report [16].  $\square$

In particular, any work-conserving policy is throughput optimal and mean-delay-optimal

<sup>1</sup>Please refer to Fig. 2 and Section II for the peak age  $A_{k, \pi}$ .

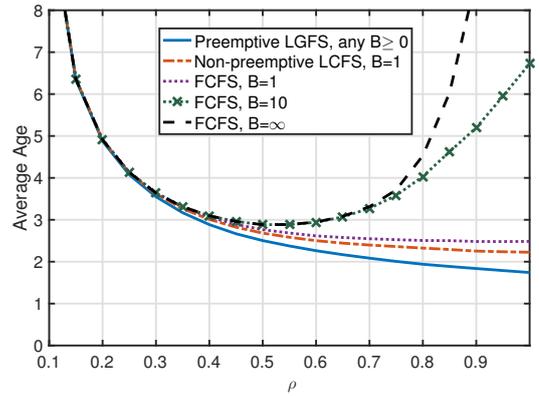


Figure 3: Average age versus traffic intensity  $\rho$  for an update system with  $m = 1$  server and queue size  $B$ .

## IV. NUMERICAL RESULTS

We present some numerical results to illustrate the age performance of different policies. The packet service times are exponentially distributed with mean  $1/\mu = 1$ , which is *i.i.d.* across time and servers. The inter-generation times are *i.i.d.* Erlang-2 distribution with mean  $1/\lambda$ . The number of servers is  $m$ . Hence, the traffic intensity is  $\rho = \lambda/m\mu$ . The queue size is  $B$ , which is a non-negative integer.

Figure 3 illustrates the time-average age versus  $\rho$  for an information-update system with  $m = 1$  server. The time difference between packet generation and arrival ( $a_i - s_i$ ) is zero, i.e., the update packets arrive in the same order of their generation times. One can observe that the preemptive LGFS policy achieves a better (smaller) age than the FCFS policy analyzed in [4], and the non-preemptive LCFS policy with queue size  $B = 1$  [9] which was also named “M/M/1/2\*” in [7]. Note that in these prior studies, the time-average age was characterized only for the special case of Poisson arrival process. Moreover, with ordered arrived packets at the server, the LGFS policy and LCFS policy have the same age performance.

Figure 4 plots the time-average age versus  $\rho$  for an information-update system with  $m = 5$  servers. The time difference between packet generation and arrival, i.e.,  $a_i - s_i$ , is modeled to be either 1 or 100, with equal probability. We found that the age performance of each policy is better than that in Fig. 3, because of the diversity provided by five servers. In addition, the preemptive LGFS policy achieves the best age performance among all plotted policies. It is important to emphasize that the age performance of the preemptive LGFS policy remains the same for any queue size  $B \geq 0$ . However, the age performance of the non-preemptive LGFS policy and FCFS policy varies with the queue size  $B$  when there are multiple servers. We also observe that the average age in case of FCFS policy with  $B = \infty$  blows up when the traffic intensity is high. This is due to the increased congestion in the network which leads to a delivery of stale packets. Moreover, in case of FCFS policy with  $B = 10$ , the average age is high but bounded at high traffic intensity, since the fresh packet has a better opportunity to be delivered in a relatively short period compared with FCFS policy with  $B = \infty$ . These numerical results validate Theorem 1.

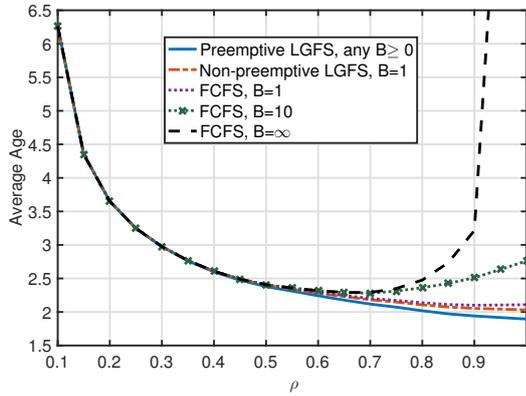


Figure 4: Average age versus traffic intensity  $\rho$  for an update system with  $m = 5$  servers and queue size  $B$ .

## V. CONCLUSION

In this paper, we considered an information-update system, in which update packets are forwarded to a destination through multiple network servers. It was showed that, if the packet service times are *i.i.d.* exponentially distributed, then for any given arrival process and queue size, the preemptive LGFS policy simultaneously optimizes the data freshness, throughput, and delay performance among all causally feasible policies. We will extend these results to more general system settings with general service time distributions.

### APPENDIX A PROOF OF LEMMA 3

Suppose that a packet is delivered in policy  $P$  which has a time stamp  $\hat{\alpha}_{l,P}$ . By coupling, a packet is delivered in policy  $\pi$  at the same time, which has a time stamp  $\hat{\alpha}_{l,\pi}$ . Hence, we have

$$\hat{U}'_P = \hat{\alpha}_{l,P} \geq \hat{\alpha}_{l,\pi} = \hat{U}'_\pi, \quad (11)$$

which holds for any queue size  $B \geq 0$ . Invoking the definition of the system state, we can obtain

$$\begin{aligned} \hat{\alpha}'_{i,P} &= \hat{U}'_P, \quad i = 1, \dots, l-1, \\ \hat{\alpha}'_{i,\pi} &= \hat{U}'_\pi, \quad i = 1, \dots, l-1. \end{aligned} \quad (12)$$

In policy  $P$ , if the queue is empty, then one server will be idle after the packet is delivered (i.e., served). Otherwise, if there exist packets in the queue, one of these packets will be assigned to the available server. Because policy  $P$  follows a preemptive LGFS principle, the time stamp of this assigned packet must be smaller than  $\hat{\alpha}_{l,P}$ , and hence the packet is stale. In both cases, we will have

$$\begin{aligned} \hat{\alpha}'_{i,P} &= \hat{U}'_P, \\ \hat{\alpha}'_{i,P} &= \hat{\alpha}_{i,P}, \quad i = l+1, \dots, m. \end{aligned} \quad (13)$$

We consider two cases for policy  $\pi$ .

Case 1: After the packet delivery, the queue is empty or the packet which will be assigned to the server is outdated. In this case, similar to (13), we have

$$\begin{aligned} \hat{\alpha}'_{l,\pi} &= \hat{U}'_\pi, \\ \hat{\alpha}'_{i,\pi} &= \hat{\alpha}_{i,\pi}, \quad i = l+1, \dots, m. \end{aligned} \quad (14)$$

Combining this with (7), (11), (12) and (13), we can obtain

$$\begin{aligned} \hat{\alpha}'_{i,P} &= \hat{U}'_P \geq \hat{U}'_\pi = \hat{\alpha}'_{i,\pi}, \quad i = 1, \dots, l, \\ \hat{\alpha}'_{i,P} &= \hat{\alpha}_{i,P} \geq \hat{\alpha}_{i,\pi} = \hat{\alpha}'_{i,\pi}, \quad i = l+1, \dots, m. \end{aligned} \quad (15)$$

Case 2: After the packet delivery, the packet which will be assigned to the server is not outdated. Assume this packet has a time stamp  $h$ . Without loss of generality, let us assume that  $h = \hat{\alpha}'_{k,\pi}$ . This packet has three possible scenarios under policy  $P$ : (i) This packet is preempted by a fresher packet and stored back in the queue or a fresher packet is delivered under policy  $P$ . Thus, this packet is outdated regarding to policy  $P$ . Hence, we have  $\hat{\alpha}'_{k,P} \geq \hat{\alpha}'_{1,P} \geq h = \hat{\alpha}'_{k,\pi}$ . (ii) This packet is already delivered in policy  $P$ , then  $\hat{\alpha}'_{k,P} \geq \hat{U}'_P \geq h = \hat{\alpha}'_{k,\pi}$ . (iii) This packet is being processed under policy  $P$ . Without loss of generality, let us assume that  $h = \hat{\alpha}'_{n,P}$ . We must have  $n \leq k$ . Otherwise, policy  $\pi$  would have  $n - k$  fresh packets that do not exist under policy  $P$ , which contradicts the fact that the arrival process is fixed for all policies and that policy  $P$  is preemptive LGFS. Hence, we have  $\hat{\alpha}'_{k,P} \geq \hat{\alpha}'_{n,P} = h = \hat{\alpha}'_{k,\pi}$ . As a result, for all these cases we have

$$\hat{\alpha}'_{i,P} \geq \hat{\alpha}'_{i,\pi}, \quad i = 1, \dots, m. \quad (16)$$

Hence, (8) holds for any queue size  $B \geq 0$ , which complete the proof.

## REFERENCES

- [1] B. Adelberg, H. Garcia-Molina, and B. Kao, "Applying update streams in a soft real-time database system," in *ACM SIGMOD Record*, 1995, vol. 24, pp. 245–256.
- [2] J. Cho and H. Garcia-Molina, "Synchronizing a database to improve freshness," in *ACM SIGMOD Record*, 2000, vol. 29, pp. 117–128.
- [3] L. Golab, T. Johnson, and V. Shkapenyuk, "Scheduling updates in a real-time stream warehouse," in *Proc. IEEE 25th International Conference on Data Engineering*, March 2009, pp. 1207–1210.
- [4] S. Kaul, R. Yates, and M. Gruteser, "Real-time status: How often should one update?," in *Proc. IEEE INFOCOM*, 2012, pp. 2731–2735.
- [5] R. Yates and S. Kaul, "Real-time status updating: Multiple sources," in *Proc. IEEE Int. Symp. Inform. Theory*, July 2012.
- [6] L. Huang and E. Modiano, "Optimizing age-of-information in a multi-class queueing system," in *Proc. IEEE Int. Symp. Inform. Theory*, 2015.
- [7] M. Costa, M. Codreanu, and A. Ephremides, "Age of information with packet management," in *Proc. IEEE Int. Symp. Inform. Theory*, June 2014, pp. 1583–1587.
- [8] N. Pappas, J. Gunnarsson, L. Kratz, M. Kountouris, and V. Angelakis, "Age of information of multiple sources with queue management," in *Proc. IEEE ICC*, June 2015, pp. 5935–5940.
- [9] S. Kaul, R. Yates, and M. Gruteser, "Status updates through queues," in *Conf. on Info. Sciences and Systems*, Mar. 2012.
- [10] C. Kam, S. Kompella, and A. Ephremides, "Effect of message transmission diversity on status age," in *Proc. IEEE Int. Symp. Inform. Theory*, June 2014, pp. 2411–2415.
- [11] M. Costa, S. Valentin, and A. Ephremides, "On the age of channel information for a finite-state markov model," in *Proc. IEEE ICC*, June 2015, pp. 4101–4106.
- [12] T. Bacinoglu, E. T. Ceran, and E. Uysal-Biyikoglu, "Age of information under energy replenishment constraints," in *Proc. Info. Theory and Appl. Workshop*, Feb. 2015.
- [13] R. Yates, "Lazy is timely: Status updates by an energy harvesting source," in *Proc. IEEE Int. Symp. Inform. Theory*, 2015.
- [14] Y. Sun, E. Uysal-Biyikoglu, R. Yates, C. E. Koksal, and N. B. Shroff, "Update or wait: How to keep your data fresh," in *Proc. IEEE INFOCOM*, April 2016.
- [15] M. Shaked and J. G. Shanthikumar, *Stochastic orders*, Springer Science & Business Media, 2007.
- [16] A. M. Bedewy, Y. Sun, and N. B. Shroff, "Optimizing data freshness, throughput, and delay in multi-server information-update systems," *arXiv preprint arXiv:1603.06185*, 2016.