



# An opportunistic and non-anticipating size-aware scheduling proposal for mean holding cost minimization in time-varying channels



Ianire Taboada<sup>a,\*</sup>, Fidel Liberal<sup>a</sup>, Peter Jacko<sup>b</sup>

<sup>a</sup> University of the Basque Country, ETSI Bilbao, Alameda Urquijo s/n, 48013 Bilbao, Spain

<sup>b</sup> Lancaster University Management School, Bailrigg, Lancaster, Lancashire LA1 4YX, UK

## ARTICLE INFO

### Article history:

Available online 11 July 2014

### Keywords:

Opportunistic scheduling  
Non-anticipating size-aware scheduling  
Mean holding cost minimization  
Whittle index  
Markov Decision Process

## ABSTRACT

In this paper we study how to design a scheduling strategy aimed at minimizing the average holding cost for flows with general size distribution when the feasible transmission rate of each user varies randomly over time. We employ a Whittle-index-based approach in order to achieve an opportunistic and non-anticipating size-aware scheduling index rule proposal. When the flow size distribution belongs to the Decreasing Hazard Rate class, we propose the so-called Attained Service Potential Improvement index rule, which consists in giving priority to the flows with the highest ratio between the current attained-service-dependent completion probability and the expected potential improvement of this completion probability. We further analyze the performance of the proposed scheduler, concluding that it outperforms well-known opportunistic disciplines.

© 2014 Elsevier B.V. All rights reserved.

## 1. Introduction

Undoubtedly, due to the massive use of mobile Internet applications, one of the fundamental challenges that networks providers nowadays face is the management for sharing radio resources among users' traffic flows. Thus, motivated by the necessity of obtaining an implementable scheduler in channels with randomly time-varying capacity such as wireless links, in this paper we aim at characterizing in closed-form a novel channel-aware or opportunistic scheduler for the problem of minimizing the expected holding cost in a scenario where flows arrive and depart upon service completion.

Although in time-varying transmission conditions taking advantage of the channel opportunistic gains seems good, short-term disciplines that serve the user with the best instantaneous rate, such as Max Rate, perform very poorly in this setting (see for example [1]). Moreover, due to the complexity of the present problem, flow-level opportunistic scheduling in time-varying systems has been analyzed by approximate techniques [2–4] to design simple schedulers, and in the asymptotic regimes to study optimality and maximal stability [5,6]. Nevertheless, all these works deal with unrealistic assumptions in reference to traffic flow sizes. On the one hand, exponential flow size distributions are considered for traffic modeling, which, even though they simplify the resolution of those problems, are far from reality. On the other hand, it is assumed that flow sizes are known by the scheduler, while in current network systems they are not.

In this paper we take a step forward towards removing the assumption of exponential sizes. Furthermore, we incorporate non-anticipating size-awareness by taking into account the bits that have been transferred of a flow: the attained service.

\* Corresponding author. Tel.: +34 94 6017352; fax: +34 94 6014259.

E-mail addresses: [ianire.taboada@ehu.es](mailto:ianire.taboada@ehu.es) (I. Taboada), [fidel.liberal@ehu.es](mailto:fidel.liberal@ehu.es) (F. Liberal), [peter.jacko@gmail.com](mailto:peter.jacko@gmail.com) (P. Jacko).

In the context of non-anticipating strategies the work [7] of Gittins is relevant, which based on the attained service of jobs proposed an index rule that minimizes the mean holding cost when channel capacity is constant. [8,9] propose some heuristics using Gittins approach for the case of time-varying capacity, the first work for deterministic channels. However, to the best of our knowledge, there is no strong analytically founded and well-performing scheduling proposal for randomly time-varying channels that combines this kind of size-awareness with channel-awareness using Gittins approach.

Therefore, in this paper we aim at developing in closed form a simple opportunistic and non-anticipating size-aware scheduler for the problem of minimizing the expected holding cost in random time-varying channels for flows with general size distribution. In order to achieve our goal, the work presented in [3] has been relevant. [3] considers a finite number of channel conditions and exponentially distributed flow sizes, and its associated optimal scheduling problem is formulated as a Markov Decision Process (MDP). Due to the impossibility of solving the general model for being PSPACE-hard [10], the authors of [3] propose a simple Whittle-index-based [11,12] heuristic scheduler, which they show to perform well in several simulation scenarios. Moreover, as shown in [6,13], this Whittle-index-based proposal is maximal stable and fluid-optimal, as well as asymptotically optimal under some assumptions as the number of flows and servers grows to infinity [14]. Hence, so as to achieve our aim, we design a simple Whittle-index-based scheduler extending the framework presented in [3] to the case of general size distribution.

The rest of the paper is structured as follows. In Section 2 we present the problem description. We formulate the problem as a MDP model in Section 3. We design the Whittle-index-based scheduler in Section 4, and we evaluate its performance in Section 5. Finally, Section 6 gathers the main conclusions of the paper. For the sake of readability, some of the proofs are postponed to the Appendix.

## 2. Problem description

We analyze a discrete-time job scheduling problem aimed at minimizing the expected holding cost, in which the feasible transmission rate of each user varies randomly over time. Scheduling decisions are taken at the beginning of time slots  $t \in \mathcal{T} := \{0, 1, \dots\}$ , and are applied during a slot duration.

We consider a system without arrivals with  $K$  jobs waiting for service, which incur a holding cost  $c_k > 0$  per slot while the flow transmission is not completed. We will use terms job/flow/user interchangeably throughout the paper. The job size in bits  $x_k$  follows a general distribution with  $\mathbb{E}[x_k] < \infty$ , characterized by its probability density function  $f_k(x)$ .

The channel of a user  $k$  can take  $N_k$  conditions from a finite set  $\mathcal{N}_k := \{1, 2, \dots, N_k\}$ . These channel conditions are associated to different transmission rates  $r_{k,n}$  (in bits), where  $r_{k,1} \leq r_{k,2} \leq \dots \leq r_{k,N_k}$ . The channel condition of a user  $k$  evolves randomly and independently of other users. We denote the probability of being in state  $n$  by  $q_{k,n}$ , having  $\sum_{n \in \mathcal{N}_k} q_{k,n} = 1$ .

The server makes use of the instantaneous channel information ( $r_{k,n}$  and  $q_{k,n}$ ) and the instantaneous attained service ( $a_k$ ) of each user in order to take decisions. In each decision slot it allows the transmission of a single flow, and it is assumed to be preemptive (the service of a job can be interrupted at the beginning of a slot even if not completed). We refer to job completion or departure probability,  $\mathbb{P}(a_k < X_k \leq a_k + r_{k,n} | X_k > a_k)$ , as  $\mu_{k,(a,n)}$ .

## 3. MDP formulation

In this section we present a MDP formulation of the scheduling problem described in Section 2. First, we provide the MDP model of each job  $k$ . Then, we formulate the optimization problem for the joint MDP model, which takes into account all the jobs in the system.

### 3.1. MDP model of a job

In each time slot  $t$ , a user  $k$  that is in a state  $s_k \in \mathcal{S}_k$  can be allocated either zero or full capacity. We refer to  $\mathcal{B} := \{0, 1\}$  as the action space, in which action 0 means not serving and action 1 serving. Thus, the dynamics of user  $k$  is captured by the action process  $b_k(\cdot)$  and the state process  $s_k(\cdot)$ . As a result of taking action  $b_k(t)$  in state  $s_k(t)$ , the user  $k$  earns a reward, consumes the allocated capacity and evolves its state in the time slot  $t+1$ . In such a way, each user  $k$  is defined independently of other users by tuple  $(\mathcal{S}_k, (\mathbf{R}_{k,s}^b)_{b \in \mathcal{B}}, (\mathbf{W}_{k,s}^b)_{b \in \mathcal{B}}, (\mathbf{P}_{k,s}^b)_{b \in \mathcal{B}})$  as follows:

- $\mathcal{S}_k := (\mathcal{A}_k \times \{1, 2, \dots, N_k\}) \cup \{*\}$  is the state space, where for  $a \in \mathcal{A}_k$  and  $n \in \mathcal{N}_k$  in each state  $(a, n)$  the job is uncompleted, and state  $*$  represents a flow already completed.
- $\mathbf{R}_k^b := (R_{k,s}^b)_{s \in \mathcal{S}_k}$ , where  $R_{k,s}^b$  is the expected one-slot reward earned by user  $k$  at state  $s$  if action  $b$  is decided at the beginning of a slot; it is defined as the expected cost of remaining in the system as:

$$R_{k,(a,n)}^0 = -c_k, \quad R_{k,(a,n)}^1 = -c_k(1 - \mu_{k,(a,n)}), \quad R_{k,*}^b = 0;$$

- $\mathbf{W}_k^b := (W_{k,s}^b)_{s \in \mathcal{S}_k}$ , where  $W_{k,s}^b$  is the expected one-slot capacity consumption or work required by user  $k$  at state  $s$  if action  $b$  is decided at the beginning of a slot, so that

$$W_{k,s}^0 = 0, \quad W_{k,s}^1 = 1;$$

Download English Version:

<https://daneshyari.com/en/article/462965>

Download Persian Version:

<https://daneshyari.com/article/462965>

[Daneshyari.com](https://daneshyari.com)