

# Opportunistic Scheduling of Flows with General Size Distribution in Wireless Time-Varying Channels

I. Taboada<sup>a</sup>, P. Jacko<sup>b</sup>, U. Ayesta<sup>a,c,d,e</sup>, F. Liberal<sup>a</sup>

<sup>a</sup>UPV/EHU, Univ. of the Basque Country, 48013 Bilbao, Spain

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

<sup>c</sup>CNRS, LAAS, 7 avenue du colonel Roche, 31400 Toulouse, France

<sup>d</sup>IKERBASQUE, Basque Foundation for Science, 48011 Bilbao, Spain

<sup>e</sup>Univ. de Toulouse, LAAS, 31400 Toulouse, France

email: {ianire.taboada; urtzi.ayesta; fidel.liberal}@ehu.es; peter.jacko@gmail.com

**Abstract**—In this paper we study how to design an opportunistic scheduler when flow sizes have a general service time distribution with the objective of minimizing the expected holding cost. We allow the channel condition to have two states which in particular covers the important special case of ON/OFF channels. We formulate the problem as a multi-armed restless bandit problem, a particular class of Markov decision processes. Since an exact solution is out of reach, we characterize in closed-form the Whittle index, which allows us to define a heuristic scheduling rule for the problem. We then particularize the index to the important subclass of distributions with a decreasing hazard rate. We finally evaluate the performance of the proposed Whittle-index based scheduler by simulation of a wireless network. The numerical results show that the performance of the proposed scheduler is very satisfactory.

**Index Terms**—Whittle index rule, opportunistic scheduling, size-aware scheduling, mean delay optimization, wireless networks.

## I. INTRODUCTION

Due to fading and interference effects, the quality of a wireless downlink channel, and hence its transmission rate, fluctuates over time. This has given rise to so-called opportunistic schedulers, that is, scheduling disciplines that take advantage of the channel fluctuations by serving a user who has a good channel condition with respect to its own statistical behavior, see for example [1]. In a dynamic scenario, where users arrive and depart upon service completion, and when the objective is to minimize the mean number of users in the system or the mean waiting time, the challenge is how to keep the number of uncompleted flows low while taking advantage of opportunistic gains. Interestingly, greedy or short-term policies that serve the user with best instantaneous capacity, such as the Max Rate policy, are known to perform very poorly in this scenario, see for example [1] or [2].

Over the years the literature on performance evaluation and optimal scheduling of flows in wireless downlink channels has grown tremendously (see [1]–[6]). A particularly relevant paper to our work is [2], where the authors consider a finite

number of channel conditions with exponentially distributed flow sizes, and formulate the optimal scheduling problem in the framework of restless bandit problems [7], which is a fundamental model for resource allocation problems. Notoriously difficult in general, it does not allow to obtain tractable optimal solution, for being PSPACE-hard [8]. Using recent advances in the theory of restless bandits [9], the authors of [2], characterizing the so-called *Whittle index*, develop a simple Whittle index-based heuristic scheduler that they illustrated in simulations to perform well. It was shown in [10], [11] that the Whittle index-based scheduler of [2] has the property of maximal stability and fluid-optimality. Moreover, the Whittle index rule has (in case of no arrivals and departures) the property of being asymptotically optimal as the number of flows and servers grows to infinity, under some technical assumptions [12].

In the present paper we aim at characterizing the optimal scheduling of flows with general size distribution in systems with time-varying service channels. That is, we consider a flow-level model with the objective of minimizing the expected holding cost, which covers as special cases the minimization of the mean number of uncompleted jobs and of the mean delay (waiting time). The extension to general flow size distribution is a major step in comparison with previous literature which assumes exponentially distributed sizes. An exception is [1], where it is shown that the Proportional Fair policy, under some assumptions, with generally distributed flow sizes can be accurately modeled by a state-dependent processor-sharing queue.

Given the difficulty of the problem we restrict ourselves to the case of two channel conditions, *good* and *bad* condition, so as to get fundamental insights into the problem. This *Gilbert-Elliot* channel model has been extensively studied in the context of wireless channels that dates back to the original paper by Gilbert in [13] and covers also the important special case of ON/OFF channels.

In such a way, in order to achieve our goal we extend the

framework introduced in [2] to the case of general flow size distribution. In our main contribution we derive in closed-form the Whittle index, which allows us to define a heuristic scheduling rule for the problem considered.

The rest of the paper is organized as follows. In Section II we present the problem description. Section III formulates the problem as a Markov Decision Process (MDP). We obtain the Whittle index-based solution in Section IV, and its performance is evaluated in Section V. Finally, Section VI gathers the main conclusions of the paper. For the sake of readability, proofs are postponed to the appendix.

## II. PROBLEM DESCRIPTION

We consider a time-slotted system, so that we study a discrete-time job scheduling problem. The decisions are taken in time epochs/instants  $t \in \mathcal{T} := \{0, 1, \dots\}$ , and are applied during the time slots  $t \in \mathcal{T}$ , where slot  $t$  corresponds to the interval between epochs  $[t, t + 1)$ .

### A. Users

Suppose that there are  $K$  users, labeled  $k \in \mathcal{K} := \{1, 2, \dots, K\}$ . Each user is uniquely associated with the flow/job (used interchangeably throughout the paper) it requests to download and with the dedicated wireless channel. For every user  $k$  a holding cost  $c_k > 0$  is paid for every slot while requested download is uncompleted. In the mathematical model we do not consider arrivals of new users. However, in the numerical section we will compare the performance of scheduling disciplines in the presence of arrivals.

*a) Job Sizes:* The (integer-valued) job size  $x_k$  of user  $k$  is measured in *bits* and has a general distribution with  $\mathbb{E}[x_k] < \infty$  for user  $k \in \mathcal{K}$ . The job sizes of users are assumed to be independently distributed, and we denote by  $F_k(x)$  and  $f_k(x)$  the cumulative distribution and density function, respectively.

*b) Channel model:* For user  $k$  the quality of the channel (the *channel condition*) evolves according to a distribution which may depend on  $k$ , independently of all other users present in the system. We assume that for every user the channel can be in two conditions. Different channel conditions correspond to different transmission rates associated with the available modulation and coding schemes. User  $k$  is in channel condition  $n = 1, 2$  with probability  $q_{k,n}$  having  $q_{k,1} + q_{k,2} = 1$ , and  $r_{k,n}$  bits are transmitted to channel in each transmission slot. We assume that  $0 \leq r_{k,1} < r_{k,2}$ .

*c) Departure probability:* We will find it useful to define the *Generalized Hazard Rate* (GHR) function for a distribution  $F_k(x)$  and any positive integer  $r$  as:

$$H_k(x, r) = \frac{F_k(x+r) - F_k(x)}{1 - F_k(x)} \quad (1)$$

Then we can easily see that the departure (or job completion) probability of user  $k$  with attained service  $a$  (the bits that have been transferred of a job) if served in channel condition  $n$  is

$$\mu_{k,(a,n)} = H_k(a, r_n) \quad (2)$$

We will say that a probability distribution belongs to the Decreasing (Increasing) Generalized Hazard Rate (DGHR, IGHR) class if the GHR is decreasing (increasing) for both transmission rates  $r_{k,n}$ .

We further present the job size distributions used in this paper. First, we define the Pareto distribution with shape parameter  $\alpha > 1$  and scale parameter  $\gamma > 0$  whose density function for all  $x \geq 0$  is:

$$f(x) = \frac{\gamma\alpha}{(1+\gamma x)^{\alpha+1}} \quad (3)$$

Second, we define the Weibull distribution with shape parameter  $\alpha > 0$  and scale parameter  $\gamma > 0$  whose density function for all  $x \geq 0$  is:

$$f(x) = \frac{\alpha}{\gamma} \left(\frac{x}{\gamma}\right)^{\alpha-1} \cdot e^{-\left(\frac{x}{\gamma}\right)^\alpha} \quad (4)$$

The Pareto distribution is known to have the property of Decreasing Hazard Rate (DHR). The Weibull distribution with  $\alpha > 1$  belongs to the Increasing Hazard Rate (IHR) class. It can easily be seen that a DHR (IHR) distribution is necessarily DGHR (IGHR). Thus, Pareto belongs to the class DGHR and Weibull with  $\alpha > 1$  to IGHR.

### B. Server

At the beginning of every slot  $t$ , the server (base station) observes the actual channel condition by receiving the Channel Quality Indicator (CQI) and the attained service retrieved from its memory of all the users present in the system, and decides which of them to serve during the slot. We assume that the server is preemptive; that is, at every decision epoch it is permitted to suspend the service of a user whose job is not yet concluded.

## III. MDP APPROACH

In this section we formulate the problem described in Section II by MDP framework. We aim at minimizing the mean holding cost, which covers in particular (when all the holding costs are equal) the minimization of the mean delay and the mean number of uncompleted jobs. The problem studied here fits the multi-armed restless bandit problem adapted to job scheduling (see [14]).

### A. MDP Model of a Job

At the beginning of every time slot, a user  $k$  can only be either served or not. We denote by  $\mathcal{B}$  the action space of user  $k$ ;  $\mathcal{B} := \{0, 1\}$ , where the action 0 means not serving and action 1 serving.

Each user  $k$  is defined by tuple  $\left(\mathcal{S}_k, \left(\mathbf{R}_{k,s}^b\right)_{b \in \mathcal{B}}, \left(\mathbf{W}_{k,s}^b\right)_{b \in \mathcal{B}}, \left(\mathbf{P}_{k,s}^b\right)_{b \in \mathcal{B}}\right)$  as follows:

- State space  $\mathcal{S}_k = (\mathcal{A}_k \times \{1, 2\}) \cup \{*\}$  is the set of all states  $s$  for a user  $k$  such that

- as long as the job is uncompleted: bi-dimensional state  $s = (a, n)$ , with components attained service  $a \in \mathcal{A}_k$ , being  $\mathcal{A}_k$  the space of possible attained service levels, and channel condition  $n \in \{1, 2\}$ ;

– if the job is completed: absorbing state  $s = *$ .

- $\mathbf{R}_k^b := (R_{k,s}^b)_{s \in \mathcal{S}_k}$ , where  $R_{k,s}^b$  is the expected one-slot reward received from user  $k$  at state  $s$  if action  $b$  is decided at the beginning of a slot,

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

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

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

- $\mathbf{P}_k^b := (p_k^b(s, s'))_{s, s' \in \mathcal{S}_k}$ , where  $p_k^b(s, s')$  is the probability for user  $k$  of moving from state  $s$  to state  $s'$  if action  $b$  is decided at the beginning of a slot,

$$p_k^0((a, n), (a, m)) = q_{k,m}, \\ p_k^1((a, n), (a + r_n, m)) = q_{k,m}(1 - \mu_{k,(a,n)}), \\ p_k^1((a, n), *) = \mu_{k,(a,n)}, \quad p_k^b(*, *) = 1;$$

Thus, the dynamics of user  $k$  is captured by state process  $s_k(t) \in \mathcal{S}_k$  and action process  $b_k(t) \in \mathcal{B}$ .

#### B. Optimization Problem, Relaxations and Decomposition

We present now the optimization problem we consider in (5). Let  $\Pi$  be the set of all admissible policies to the studied problem, and for a given discount factor  $\beta$ :

$$\max_{\pi \in \Pi} \mathbb{E}_0^\pi \left[ \sum_{t=0}^{\infty} \sum_{k \in \mathcal{K}} \beta^t R_{k,(s_k(t))}^{b_k(t)} \right] \\ \text{subject to } \sum_{k \in \mathcal{K}} b_k(t) = 1 \quad \forall t \quad (5)$$

The optimization problem formulated in (5) can be relaxed by requiring to serve a job per slot on average as proposed in [7], which is further approached by Lagrangian methods and can be decomposed into a single-job price-based parametrized optimization problem (see [2] for more details). For a price,  $v$ , we will therefore study the user- $k$  subproblem:

$$\max_{\pi_k \in \Pi} \sum_{t=0}^{\infty} \mathbb{E}_0^\pi \beta^t \left[ R_{k,(s_k(t))}^{b_k(t)} - v W_{k,(s_k(t))}^{b_k(t)} \right] \quad (6)$$

#### IV. WHITTLE INDEX-BASED SOLUTION

In this section we study the single-user parametrized problem (6), and we aim at obtaining an optimal solution in terms of the Whittle index. The restless bandit problem and their index-based solution were introduced in [7], generalizing the so-called Gittins index policy that was proved optimal for the multi-armed bandit problem in [15]. The idea of Whittle [7] was to introduce a function to measure a (dynamic) priority for serving, so that a simple scheduling rule appears: *At every slot, serve the user with the highest actual Whittle index value.* Since the Whittle index is a function of state, it is dynamic, and gives rise to an opportunistic scheduler for a wireless

network. Such a Whittle index-based scheduler, however, is only a heuristic for the problem considered. Nevertheless, it has been shown for an increasing amount of models (see [10]–[12]) that the Whittle index rule performs strongly, possesses asymptotically optimal properties, and may even be optimal in some circumstances.

#### A. Whittle Index and Indexability

Let us define the *Whittle index values* and *indexability*, following [16], formalizing the intuitive definition given in [7]. The Whittle index measures the expected efficiency of serving a user for every state, since it is the break-even value of the Lagrangian parameter  $v$ , which can be interpreted as the per-slot cost of serving. From now on we omit the user label  $k$ .

**Definition 1 (Indexability).** *We say that the problem (6) is indexable if there exist values  $v_s^* \in \mathbb{R} \cup \{-\infty, \infty\}$  for all  $s \in \mathcal{S}$  such that*

- it is optimal to serve the user in state  $s$  if  $v_s^* \geq v$ , and*
- it is optimal not to serve the user in state  $s$  if  $v_s^* \leq v$ .*

*Such values  $v_s^*$  are called the (Whittle) index values, and define an optimal index policy for the problem.*

Index policies and indexability has been established for several classes of problems [7], [9], [17]; however, not all problems are indexable. There is an algorithm for computing these index values and verifying sufficient indexability conditions, called *Adaptive-Greedy*, shortly *AG-algorithm*, (see [9] for a survey). If a problem is indexable then the *AG*-algorithm computes the index values.

We now state the Conjecture 1 of indexability, which we have not proved rigorously for our model due to its complexity. Indexability was, however, established for geometric job sizes in [2].

**Conjecture 1.** *Problem (6) is indexable.*

We assume that Conjecture 1 holds throughout the rest of this section. For a given job, we have employed the *AG*-algorithm to compute the Whittle index values numerically. In such a way, we have performed extensive numerical experiments, based on which we further conjecture the following fundamental properties. We illustrate these results in Figure 1 and Figure 2. (All the Whittle index values shown in the figures are normalized, i.e., multiplied by  $1 - \beta$ , to avoid large values obtained for  $\beta \approx 1$ .)

- All the index values for the good condition are greater than those for the bad condition. In particular, for the same attained service the Whittle index value in the good condition is greater:  $v_{(a,2)}^* > v_{(a,1)}^*$ . This property is illustrated in the top-graph of Figure 1.
- For the same channel condition Whittle index values are ordered by GHR as illustrated in Figure 1 and Figure 2. We show that for a DGHR size distribution (in particular, Pareto distribution) index values are decreasing with attained service (see Figure 1) and that for an IGHR size distribution (in particular, Weibull distribution) index

values are increasing with attained service (shown in Figure 2).

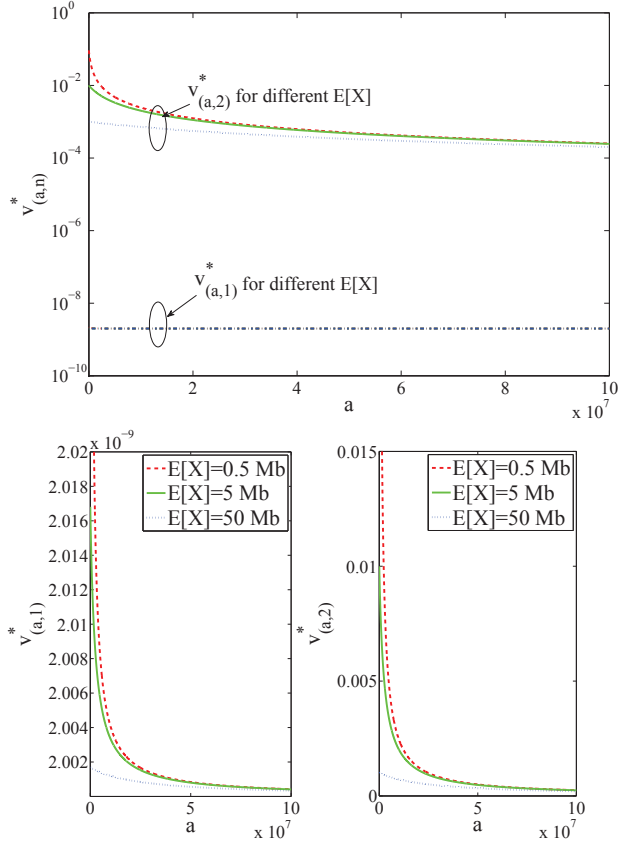


Fig. 1. Normalized Whittle indices for a DGHR size distribution. A Pareto distribution is used with  $\alpha = 1.5$ ,  $q_2 = 0.5$ ,  $r_1 = 8.4$  Kb, and  $r_2 = 16.8$  Kb.

### B. Closed-Form Characterization of Whittle Index

In the previous subsection we illustrated the properties of the Whittle index obtained numerically by applying  $\mathcal{AG}$ -algorithm. The use of the  $\mathcal{AG}$ -algorithm for index value computation becomes time-consuming and sometimes even intractable (it performs  $\mathcal{O}(A^3)$  elementary operations for computing all the index values, where  $A := \max_k |\mathcal{A}_k|$ ). Further, it is numerically unstable for  $\beta \approx 1$ .

Therefore, in this subsection we set out to derive a closed-form characterization of the Whittle index, for which the numerically observed properties have been useful for guessing the structure of the optimal policy. This result improves the undiscounted Whittle index computation both in time and precision. As we present in Proposition 1, the Whittle index admits a rather complicated closed-form expression for general job size distribution, whose computation may still be computationally costly or even intractable.

The methodology to compute this index expression and proofs are presented in the Appendix. The proof of Proposition 1 relies on an analysis of the discounted case and obtaining the undiscounted index values in the limit  $\beta \rightarrow 1$ .

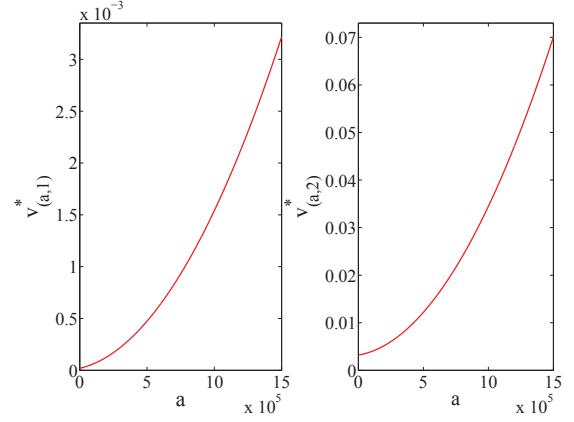


Fig. 2. Normalized Whittle indices for a IGHR size distribution. A Weibull distribution is used with  $\alpha = 3$  and  $\gamma = 5.75 \cdot 10^5$ ,  $q_2 = 0.5$ ,  $r_1 = 4.2$  Kb, and  $r_2 = 8.4$  Kb.

**Proposition 1.** The Whittle index for Problem (6),  $v_{(a,n)}^*$ , in the undiscounted case ( $\beta = 1$ ) is given by:

$$v_{(a,n)}^* = \frac{c\mu_{(a,n)} + (1 - \mu_{(a,n)})R_{a_1} - R_{a_2}}{1 + (1 - \mu_{(a,n)})W_{a_1} - W_{a_2}} \quad (7)$$

where we have defined  $a_1 := a$  and  $a_2 := a + r_n$ , and for  $z \in \{1, 2\}$ ,

if  $t_{z,0} \in \{0, 1\}$  then

$$R_{a_z} = \frac{-c(1 - \sum_{m^{z,(0)} > t_{z,0}} q_{m^{z,(0)}}) + \sum_{m^{z,(0)} > t_{z,0}} q_{m^{z,(0)}} R_{(a_z, m^{z,(0)})}}{\sum_{m^{z,(0)} > t_{z,0}} q_{m^{z,(0)}}} \quad \text{and}$$

$$W_{a_z} = \frac{\sum_{m^{z,(0)} > t_{z,0}} q_{m^{z,(0)}} W_{(a_z, m^{z,(0)})}}{\sum_{m^{z,(0)} > t_{z,0}} q_{m^{z,(0)}}},$$

and if  $t_{z,0} = 2$  then  $R_{a_z} = \lim_{\beta \rightarrow 1} \frac{-c}{1 - \beta}$  and  $W_{a_z} = 0$ , and

$$R_{(a_z, m^{z,(0)})} = -c \left( 1 - \mu_{(a_z, m^{z,(0)})} \right) - \frac{c \left( 1 - \mu_{(a_z, m^{z,(0)})} \right)}{\sum_{m^{z,(1)} > t_{z,1}} q_{m^{z,(1)}}}$$

$$\cdot \sum_{i=0}^{I(a_z)} \left( \prod_{j=1}^i \frac{\sum_{m^{z,(j)} > t_{z,j}} q_{m^{z,(j)}} \left( 1 - \mu_{\left( a_z + \sum_{k=0}^{j-1} r_{m^{z,(k)}, m^{z,(j)}} \right)} \right)}{\sum_{m^{z,(j+1)} > t_{z,j+1}} q_{m^{z,(j+1)}}} \right)$$

$$\cdot \frac{1 - \sum_{m^{z,(i+1)} > t_{z,i+1}} q_{m^{z,(i+1)}} \mu_{\left( a_z + \sum_{k=0}^i r_{m^{z,(k)}, m^{z,(i+1)}} \right)}}{L(a_z, i)}$$



$$W_{(a_z, m^{z,0})} = 1 + \frac{1 - \mu_{(a_z, m^{z,0})}}{\sum_{m^{z,1} > t_{z,1}} q_{m^{z,1}}}$$

$$\cdot \sum_{i=0}^{I(a_z)} \left( \prod_{j=1}^i \frac{\sum_{m^{z,(j)} > t_{z,j}} q_{m^{z,(j)}} \left( 1 - \mu_{\left( a_z + \sum_{k=0}^{j-1} r_{m^{z,(k)}, m^{z,(j)}} \right)} \right)}{\sum_{m^{z,(j+1)} > t_{z,j+1}} q_{m^{z,(j+1)}}} \right)$$

$$\cdot \sum_{m^{z,(i+1)} > t_{z,i+1}} q_{m^{z,(i+1)}}$$

where the thresholds  $t_{z,i} \in \{0, 1, 2\}$  are such that it is optimal to serve in state  $(a_z + \sum_{k=0}^{i-1} r_{m^{z,(k)}, m^{z,(i)}})$  for  $t_{z,i} < m^{(i)} \in \{1, 2\}$ . Further,  $I(a_z)$  is the smallest value of  $i$  that satisfies  $t_{z,i} = 2$ , and  $L(a_z, i)$  equals to 1 when  $i \neq I(a_z)$  and equals either to 1 or  $1 - \beta$  (depending on the type of size distribution) when  $i = I(a_z)$ .

Thus, for a general size distribution the Whittle index rule consists in, at every decision slot, serving the user with the highest value of (7). The values of the thresholds  $t_{z,i}$  and variables  $I(a_z)$  and  $L(a_z, i)$  will be mainly determined by the job size distribution. For a job size distribution which results in alternatively increasing and/or decreasing generalized hazard rate, the value of  $I(a_z)$  will be finite and  $L(a_z, I(a_z)) = 1 - \beta$ . In the cases of both IGHR and DGHR considering infinite attained service levels, we will have  $I(a_z) = +\infty$  and  $L(a_z, I(a_z)) = 1$ . In the next subsection we focus on an important special case of DGHR size distributions in more detail.

### C. Whittle Index for DGHR Size Distributions

In the previous subsection, we have characterized the Whittle index values via a rather complicated analytic expression. Next we focus on providing a simplified expression for certain job size distributions, which allows an easier online index computation and provides further fundamental insights. If the optimal policy possesses a ‘‘nice’’ special structure, we can exploit it to simplify the index expressions.

In the following proposition we consider the case in which size distribution belongs to DGHR class. The proof is straightforward by simplifying the expressions given in Proposition 1, and is, therefore, omitted. Note that as we conjectured in subsection IV-A, for this case the Whittle index values are nonincreasing as a function of attained service for each channel condition as well as all the indices for the good condition are greater than those for the bad condition; consequently, for the bad condition the threshold values will be  $t_{z,i} = 1$  for all  $z \in \{1, 2\}$  and  $0 \leq i \leq I(a_z) = +\infty$ , whereas for the good condition the value of the first threshold will be  $t_{z,0} = 2$ .

**Proposition 2.** Under Conjecture 1, if for a DGHR size distribution the first threshold value  $t_{z,0} = 2$  for the good condition and the threshold values  $t_{z,i} = 1$  for all  $z \in \{1, 2\}$  and  $0 \leq i \leq I(a_z) = +\infty$  for the bad condition, then the Whittle index value in the undiscounted case ( $\beta = 1$ ) is given by:

$$v_{(a,2)}^* = \lim_{\beta \rightarrow 1} \frac{c\mu_{(a,2)}}{1 - \beta} = +\infty \quad (8)$$

$$v_{(a,1)}^* = \frac{c\mu_{(a,1)} + (1 - \mu_{(a,1)})R_1 - R_2}{1 + (1 - \mu_{(a,1)})W_1 - W_2} \quad (9)$$

where, for  $z \in \{1, 2\}$ ,

$$W_z = 1 + \sum_{i=0}^{\infty} \prod_{j=0}^i (1 - \mu_{(a_z + jr_2, 2)})$$

$$R_z = \frac{-c}{q_2} \left( 1 - q_2 \mu_{(a_z, 2)} + \sum_{i=0}^{\infty} \prod_{j=0}^i (1 - \mu_{(a_z + jr_2, 2)}) \cdot (1 - q_2 \mu_{(a_z + (i+1)r_2, 2)}) \right)$$

As we can see in Proposition 2, the Whittle index value for the good condition is infinite, which extends validity of the result of [2] for exponential job sizes. Moreover, the obtained expression of the Whittle index for the bad condition is computationally feasible (up to an acceptable level of precision).

In such a way, the Whittle index-based scheduling rule we propose for DGHR distributions is as follows: at every slot  $t$ ,

- serve the user  $k$  in channel condition 2 with the highest value of  $c_k \mu_{k, (a_k(t), 2)}$ ;
- if there is no user in channel condition 2, then serve the one with highest index value  $v_{k, (a_k(t), 1)}^*$  using (9).

In the next section we will see that the proposed scheduler has a satisfactory performance.

## V. PERFORMANCE EVALUATION

In this section we study the behavior of the proposed Whittle index rule presented in Section IV. To that end, we present several simulation scenarios in which we compare its performance with priority-based schedulers already defined in the literature. Below we define the scheduling algorithms used in the carried out experimental study:

- Max Rate (MR) scheduler:  $v_{k,n}^{\text{MR}} = r_{k,n}$ .
- Proportional Fair (PF) scheduler:  $v_{k, (a,n,d)}^{\text{PF}} = \frac{r_{k,n}}{a_k/d_k}$ ; that is, the ratio of the current transmission rate and the attained throughput, where  $d_k$  is the time already spent in the system.
- $c\mu$ -rule, adapted to the attained service:  $v_{k, (a,n)}^{c\mu} = c_k \cdot \mu_{k, (a,n)}$ .

For all disciplines, in case of ties, these are resolved randomly.

We assume that users are grouped in  $K$  classes. A user from a class  $k$  is characterized by its size distribution, channel characteristics  $r_{k,n}$  and  $q_{k,n}$ , and cost,  $c_k$ . We consider only single-class and two classes of users in order to be able to

CQI	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
r	0	4.2	6.7	8.4	11.2	16.8	21.8	25.2	26.8	33.6	44.6	50.4	53.7	67.2	75.6	80.6

TABLE I  
CQIS AND CORRESPONDING RATES (KB).

Scenario	CQI	$q_2$	Size (Mb)	c
1	{3,5}	0.5	$\{\alpha = 1.5, \gamma = 4 \cdot 10^{-7}, \mathbb{E}[X] = 5\}$	1
2	{3,5}	0.1	$\{\alpha = 1.5, \gamma = 4 \cdot 10^{-7}, \mathbb{E}[X] = 5\}$	1
3	{0,5}	0.5	$\{\alpha = 1.5, \gamma = 4 \cdot 10^{-7}, \mathbb{E}[X] = 5\}$	1
4	(({5,9}), (0.5,0.5))		$\{\alpha_1 = 1.5, \gamma_1 = 4 \cdot 10^{-7}, \mathbb{E}[X_1] = 5\}$ $\{\alpha_2 = 1.5, \gamma_2 = 4 \cdot 10^{-7}, \mathbb{E}[X_2] = 5\}$	(1,5)
5	(({5,9}), (0.5,0.5))		$\{\alpha_1 = 1.5, \gamma_1 = 4 \cdot 10^{-6}, \mathbb{E}[X_1] = 0.5\}$ $\{\alpha_2 = 1.5, \gamma_2 = 4 \cdot 10^{-8}, \mathbb{E}[X_2] = 50\}$	(1,1)
6	{1,3}	0.5	$\{\alpha = 3, \gamma = 5.7 \cdot 10^5, \mathbb{E}[X] = 0.5\}$	1

TABLE II  
PARAMETERS SET IN EXPERIMENTAL STUDY.

easily point out the differences in the performance of the policies studied.

Besides, it is known that mobile Internet traffic flow sizes are properly modeled by means of Pareto distributions [18]. This way, we will particularize our study to Pareto distributed flow sizes defined in (3). For the set of parameters chosen this distribution belongs to DGHR class. Moreover, so as to show the validity of the Whittle-based approach in other type of distributions we consider a scenario with a Weibull (see (4)) size distribution that belongs to IGHR.

Referring to user arrivals in the system, we assume that a new class- $k$  user arrives in a transmission slot according to a Poisson process with rate  $\lambda_k$ . Arrival rate will determine network load,  $\rho$ , where  $\rho_k = \lambda_k \cdot \frac{E[X_k]}{r_{k,2}}$  and  $\rho = \sum_k \rho_k$ . For simplicity, in our simulations we have  $\rho_1 = \rho_2$ .

In order to simulate scenarios as realistic as possible we consider transmission rates employed in 4G networks. A mapping from CQI values to rate values is shown in Table I, which is adapted from [19]. Furthermore, we use a transmission time interval of 1 ms, value employed in current wireless networks.

In this way, we have analyzed six relevant settings, whose parameters are summarized in Table II. Note that we consider two channel conditions of Table I in each scenario and that in the last setting we focus on the Weibull size distribution. In the following we show the results achieved in these scenarios.

#### A. Scenario 1: Basic case

In this first family of simulations we consider a typical setting, which takes into account the equiprobable channel case and medium-sized self-similar flows. As appreciated in the left-graph of Figure 3, by the proposed Whittle-based discipline the distribution of the number of users is more satisfactory. Consequently, if we observe the mean delay employing different scheduling policies for varying  $\rho$  in the right-graph of Figure 3, these results indicate that Whittle behaves better than the rest of disciplines in mean delay terms.

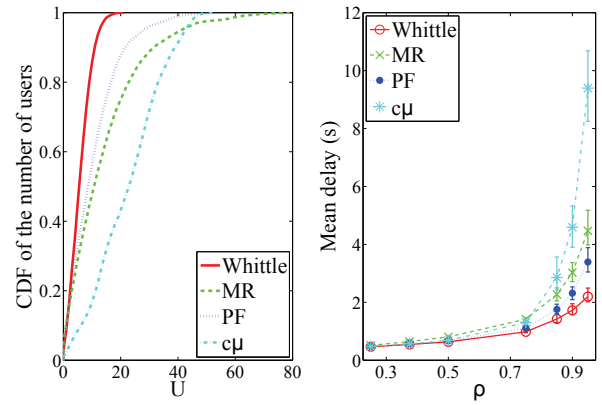


Fig. 3. CDF of the number of users (left) and mean delay (right) for Scenario 1.

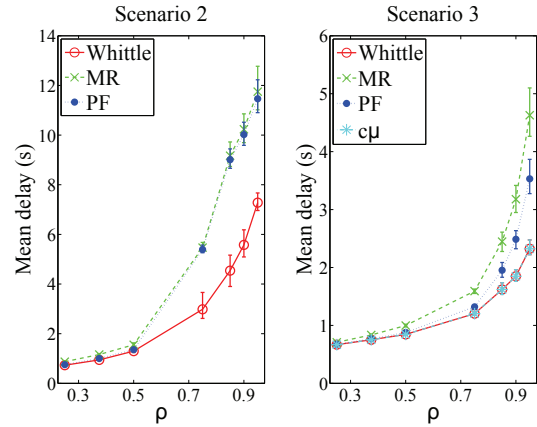


Fig. 4. Mean delay for Scenario 2 (left) and Scenario 3 (right).

#### B. Scenario 2: Low good channel probability case

In these simulations the channel conditions are bad, being the state probability in good channel low (only a 10% of the time in good channel). As depicted in the left graph of Figure 4, Whittle introduces the lowest mean delay for all the network loads considered. We omit results from  $c\mu$  rule since it is unstable for this case.

#### C. Scenario 3: ON-OFF case

We now study the special case of ON/OFF channels, which apart from being applicable to a wireless system covers other areas such as systems with time-varying breakdowns. Results collected in the right plot of Figure 4 show that Whittle and  $c\mu$  are equivalent in this setting, and they achieve the best performance in mean delay terms.

#### D. Scenario 4: Heterogeneous case in cost

In this scenario we would like to analyze the behavior of different disciplines considering that users from class 2 are more important than those belonging to class 1. To deal with this property, we assume that the holding cost in class 2 is five times higher than in class 1. In such a way, as can be seen for the aggregate of classes in Figure 5, Whittle minimizes the mean holding cost among the rest of disciplines considered.

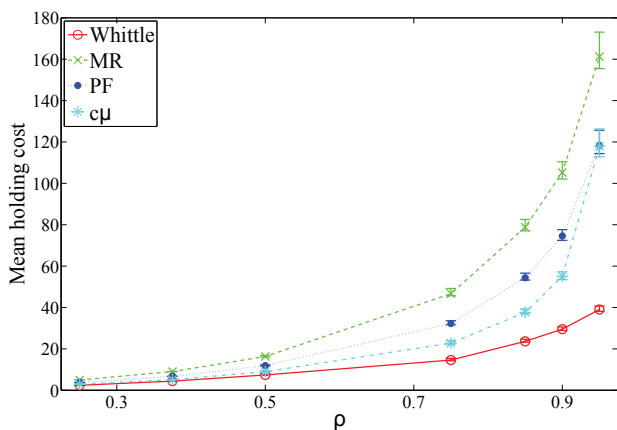


Fig. 5. Mean holding cost for Scenario 4.

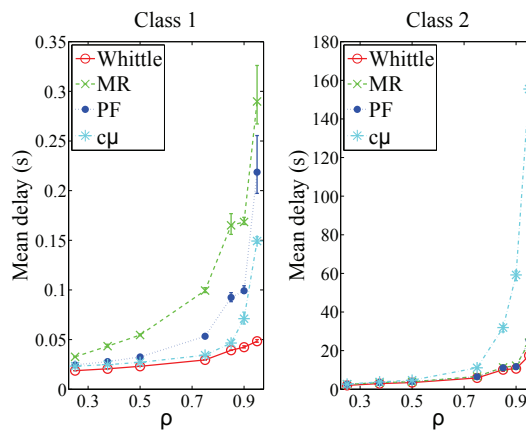


Fig. 7. Mean delay per class for Scenario 5.

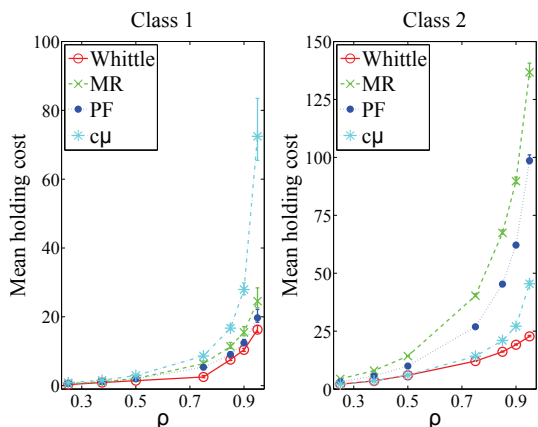


Fig. 6. Mean holding cost per class for Scenario 4.

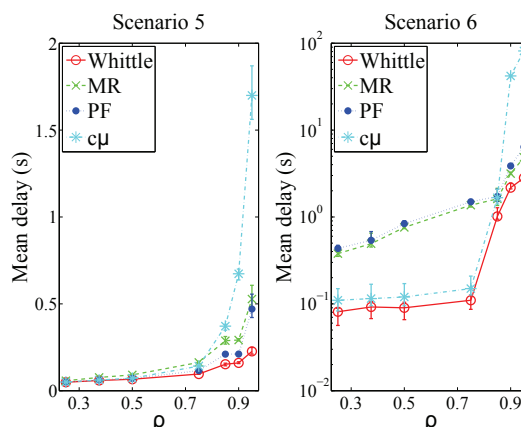


Fig. 8. Mean delay for Scenario 5 (left) and Scenario 6 (right).

Moreover, as illustrated in Figure 6 inside each class Whittle outperforms all the policies considered too.

E. Scenario 5: Heterogeneous case in size

In this setting we consider two classes that differ in the mean job size. We modify the  $\gamma$  parameter from Pareto so the mean size of class 2 is a hundred times the mean size of class 1. We conclude from Figure 7 that the Whittle-based approach is also superior in both classes, and thus, in the mixture of classes as well (see the left graph of Figure 8).

F. Scenario 6: IHR case

All the previous scenarios focus on Pareto flow size distributions. However, in order to show the validity of Whittle index not only for DGHR distributions, in this subsection we study the performance of Whittle index rule when a Weibull distribution with IGHR is considered. As concluded from the right plot in Figure 8, the Whittle index suitably minimizes mean delay respect to the rest of the policies analyzed.

VI. CONCLUSIONS

This paper represents a first attempt on the challenging problem of scheduling flows with general size distribution in a wireless time-varying channel aimed at minimizing the mean holding cost. We have proposed a heuristic scheduling policy

mathematically founded on the Whittle index, which gives an optimal solution for the relaxed problem. Moreover, we have evaluated the performance of the heuristic using simulations, concluding that it outperforms existing disciplines.

There are many interesting research problems that stem from our work. A first important question pertains to stability, since the stability region for this problem is not known. Our numerical simulations indicate that policies that serve a user in his good channel (whenever this is possible) are stable provided  $\rho < 1$ , but no rigorous proof is available. A second important contribution will be to obtain a simpler representation of the Whittle index for the case of general flow size distribution. Indeed the solution obtained in Proposition 1 does not provide insights into the structure of the scheduling rule. Another important extension will be to generalize the results to an arbitrary number of channel conditions.

REFERENCES

- [1] S. Borst, "User-level performance of channel-aware scheduling algorithms in wireless data networks," *IEEE/ACM Transactions on Networking*, vol. 13, no. 3, pp. 636–647, 2005.
- [2] U. Ayesta, M. Eraisquin, and P. Jacko, "A modeling framework for optimizing the flow-level scheduling with time-varying channels," *Performance Evaluation*, vol. 67, pp. 1014–1029, 2010.
- [3] T. Bonald, "A score-based opportunistic scheduler for fading radio channels," in *Proceedings of European Wireless*, 2004, pp. 283–292.

- [4] B. Sadiq and G. de Veciana, "Balancing SRPT prioritization vs opportunistic gain in wireless systems with flow dynamics.," in *ITC 22*, 2010.
- [5] S. Aalto, A. Penttinen, P. Lassila, and P. Osti, "On the optimal trade-off between SRPT and opportunistic scheduling," in *Proceedings of ACM Sigmetrics*, 2011.
- [6] P. Jacko, "Value of information in optimal flow-level scheduling of users with Markovian time-varying channels," *Performance Evaluation*, vol. 68, no. 11, pp. 1022–1036, 2011.
- [7] P. Whittle, "Restless bandits: Activity allocation in a changing world," *A Celebration of Applied Probability*, J. Gani (Ed.), *Journal of Applied Probability*, vol. 25A, pp. 287–298, 1988.
- [8] C. H. Papadimitriou and J. N. Tsitsiklis, "The complexity of optimal queueing network," *Mathematics of Operations Research*, vol. 24, no. 2, pp. 293–305, 1999.
- [9] J. Niño-Mora, "Dynamic priority allocation via restless bandit marginal productivity indices," *TOP*, vol. 15, no. 2, pp. 161–198, 2007.
- [10] U. Ayesta, M. Erausquin, M. Jonckheere, and I. M. Verloop, "Scheduling in a random environment: Stability and asymptotic optimality," *IEEE/ACM Transactions on Networking*, vol. 21, no. 1, pp. 258–271, 2013.
- [11] J. Kim, B. Kim, J. Kim, and Y. H. Bae, "Stability of flow-level scheduling with Markovian time-varying channels," *Performance Evaluation*, vol. 70, no. 2, pp. 148–159, 2013.
- [12] R. Weber and G. Weiss, "On an index policy for restless bandits," *Journal of Applied Probability*, vol. 27, no. 3, pp. 637–648, 1990.
- [13] E. N. Gilbert, "Capacity of a burst-noise channel," *Bell Systems Technical Journal*, vol. 39, pp. 1253–1266, 1960.
- [14] J. C. Gittins, K. Glazebrook, and R. Weber, *Multi-Armed Bandit Allocation Indices*, Wiley-Blackwell, 2011.
- [15] J. C. Gittins and D. M. Jones, "A dynamic allocation index for the sequential design of experiments," in *Progress in Statistics*, J. Gani, Ed., pp. 241–266. North-Holland, Amsterdam, 1974.
- [16] P. Jacko, "Restless bandits approach to the job scheduling problem and its extensions," in *Modern Trends in Controlled Stochastic Processes: Theory and Applications*, A. B. Piunovskiy, Ed., pp. 248–267. Luniver Press, United Kingdom, 2010.
- [17] J. C. Gittins, "Bandit processes and dynamic allocation indices," *Journal of the Royal Statistical Society, Series B*, vol. 41, no. 2, pp. 148–177, 1979.
- [18] D.S. Reginald, "The dynamics of internet traffic: Self-similarity, self-organization, and complex phenomena," *Advances in Complex Systems*, vol. 14, no. 06, pp. 905–949, 2011.
- [19] Stefania Sesia, Issam Toufik, and Matthew Baker, *LTE: the UMTS long term evolution*, Wiley Online Library, 2009.

## APPENDIX

In the discussion below we provide the methodology to obtain a closed-form Whittle index rule for problem (6) (for more details see [9]).

Let us define serving set  $\mathcal{F} \subseteq \mathcal{S}_k$ , which prescribes to serve a user  $k$  if  $(a, n) \in \mathcal{F}$ , while not to serve this user if  $(a, n) \notin \mathcal{F}$ . We will refer to states  $(a, n) \in \mathcal{F}$  as active and  $(a, n) \notin \mathcal{F}$  as passive. The Whittle index,  $v_{(a,n)}^{\mathcal{F}}$ , represents the rate between marginal reward and marginal work, where the marginal reward (work) is the difference of the expected total reward earned (work required) by serving and not serving at the initial state  $(a,n)$  and employing policy  $\mathcal{F}$  afterwards.

**Lemma 1.** *For any state  $(a,n)$  and under any policy  $\mathcal{F}$  we have*

$$v_{(a,n)}^{\mathcal{F}} = \frac{c\mu_{(a,n)} + \beta(1 - \mu_{(a,n)}) \sum_{m \in \mathcal{N}} q_m \mathbb{R}_{(a+r_n,m)}^{\mathcal{F}} - \beta \sum_{m \in \mathcal{N}} q_m \mathbb{R}_{(a,m)}^{\mathcal{F}}}{1 + \beta(1 - \mu_{(a,n)}) \sum_{m \in \mathcal{N}} q_m \mathbb{W}_{(a+r_n,m)}^{\mathcal{F}} - \beta \sum_{m \in \mathcal{N}} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}}} \quad (10)$$

*Proof:* From the definition of reward and work, respectively, we have

$$\mathbb{R}_{(a,n)}^{\mathcal{F}} = \begin{cases} -c(1 - \mu_{(a,n)}) + \beta(1 - \mu_{(a,n)}) \sum_{m \in \mathcal{N}} q_m \mathbb{R}_{(a+r_n,m)}^{\mathcal{F}} & (a, n) \in \mathcal{F} \\ -c + \beta \sum_{m \in \mathcal{N}} q_m \mathbb{R}_{(a,m)}^{\mathcal{F}} & (a, n) \notin \mathcal{F} \end{cases} \quad (11)$$

$$\mathbb{W}_{(a,n)}^{\mathcal{F}} = \begin{cases} 1 + \beta(1 - \mu_{(a,n)}) \sum_{m \in \mathcal{N}} q_m \mathbb{W}_{(a+r_n,m)}^{\mathcal{F}} & (a, n) \in \mathcal{F} \\ \beta \sum_{m \in \mathcal{N}} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}} & (a, n) \notin \mathcal{F} \end{cases} \quad (12)$$

**Lemma 2.** *For a general attained service,  $a$ , the sum of work and reward measures in (10),  $\sum_{m \in \mathcal{N}} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}}$  and  $\sum_{m \in \mathcal{N}} q_m \mathbb{R}_{(a,m)}^{\mathcal{F}}$ , respectively, can be rewritten as:*

$$\sum_{m \in \mathcal{N}} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}} = \begin{cases} \frac{\sum_{m > t} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}}}{1 - \beta + \beta \sum_{m > t} q_m} & (a, m) \in \mathcal{F} \\ 0 & (a, m) \notin \mathcal{F} \end{cases} \quad (13)$$

$$\sum_{m \in \mathcal{N}} q_m \mathbb{R}_{(a,m)}^{\mathcal{F}} = \begin{cases} \frac{-c(1 - \sum_{m > t} q_m) + \sum_{m > t} q_m \mathbb{R}_{(a,m)}^{\mathcal{F}}}{1 - \beta + \beta \sum_{m > t} q_m} & (a, m) \in \mathcal{F} \\ \frac{-c}{1 - \beta} & (a, m) \notin \mathcal{F} \end{cases} \quad (14)$$

*Proof:* Suppose that for attained service  $a$  channel states,  $m$ , higher than a channel state threshold,  $t$ , are active. Thus,  $(a, m) \in \mathcal{F}$  for  $\forall m > t$ , and by (12):

$$\mathbb{W}_{(a,m)}^{\mathcal{F}} = \beta \sum_{m' \in \mathcal{N}} q_{m'} \mathbb{W}_{(a,m')}^{\mathcal{F}}, \quad m \leq t$$

We can solve such a system of linear equations obtaining

$$\mathbb{W}_{(a,m)}^{\mathcal{F}} = \frac{\beta \sum_{m' > t} q_{m'} \mathbb{W}_{(a,m')}^{\mathcal{F}}}{1 - \beta + \beta \sum_{m' > t} q_{m'}}, \quad m \leq t \quad (15)$$

If we define  $X(t)$  as

$$X(t) = \frac{\beta}{1 - \beta + \beta \sum_{m' > t} q_{m'}} \quad (16)$$

We can express (15) as follows

$$\mathbb{W}_{(a,m)}^{\mathcal{F}} = X(t) \sum_{m' > t} q_{m'} \mathbb{W}_{(a,m')}^{\mathcal{F}}, \quad m \leq t \quad (17)$$



On the other hand, using (16) and (17), we can generally express and rewrite

$$\begin{aligned} \sum_{m \in \mathcal{N}} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}} &= (1 - \sum_{m>t} q_m) X(t) \sum_{m>t} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}} + \\ &+ \sum_{m>t} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}} = \frac{\sum_{m>t} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}}}{1 - \beta + \beta \sum_{m>t} q_m} \end{aligned} \quad (18)$$

Besides, if  $\forall m (a, m) \notin \mathcal{F}$ , then  $\sum_{m \in \mathcal{N}} q_m \mathbb{W}_{(a,m)}^{\mathcal{F}} = 0$  (solution to  $\forall m \mathbb{W}_{(a,m)}^{\mathcal{F}} = \beta \sum_{m' \in \mathcal{N}} q_{m'} \mathbb{W}_{(a,m')}^{\mathcal{F}}$ ).

Analogously, we obtain expression (14) for rewards. ■

**Lemma 3.** For a initial channel threshold,  $t = t_0$ , the expressions of  $\mathbb{W}_{(a,m>t_0)}^{\mathcal{F}}$  and  $\mathbb{R}_{(a,m>t_0)}^{\mathcal{F}}$ , in (13) and (14) respectively, are:

$$\begin{aligned} \mathbb{W}_{(a,m>t_0)}^{\mathcal{F}} &= 1 + (1 - \mu_{(a,m)}) X(t_1) \left( \sum_{m^{(1)}>t_1} q_{m^{(1)}} \right. \\ &+ \sum_{i=1}^{I(a)} \left( \prod_{j=1}^i \sum_{m^{(j)}>t_j} q_{m^{(j)}} (1 - \mu_{(a+r_m + \sum_{k=1}^{j-1} r_{m^{(k)}, m^{(j)}})}) X(t_{j+1}) \right) \\ &\cdot \left. \sum_{m^{(i+1)}>t_{i+1}} q_{m^{(i+1)}} \right) \end{aligned} \quad (19)$$

$$\begin{aligned} \mathbb{R}_{(a,m>t_0)}^{\mathcal{F}} &= -c(1 - \mu_{(a,m)}) \left( 1 + X(t_1) \left( 1 - \sum_{m^{(1)}>t_1} q_{m^{(1)}} \right) \right. \\ &+ \sum_{m^{(1)}>t_1} q_{m^{(1)}} (1 - \mu_{(a+r_m, m^{(1)})}) \left. \right) - c(1 - \mu_{(a,m)}) X(t_1) \\ &\cdot \sum_{i=1}^{I(a)} \left( \prod_{j=1}^i \sum_{m^{(j)}>t_j} q_{m^{(j)}} (1 - \mu_{(a+r_m + \sum_{k=1}^{j-1} r_{m^{(k)}, m^{(j)}})}) X(t_{j+1}) \right) \\ &\cdot \frac{1}{L(a, i)} \left( 1 - \sum_{m^{(i+1)}>t_{i+1}} q_{m^{(i+1)}} \right) \\ &+ \sum_{m^{(i+1)}>t_{i+1}} q_{m^{(i+1)}} (1 - \mu_{(a+r_m + \sum_{k=1}^i r_{m^{(k)}, m^{(i+1)}})}) \end{aligned} \quad (20)$$

*Proof:* Let us express  $\mathbb{W}_{(a,m)}^{\mathcal{F}}$  for  $m > t$  components. By (12):

$$\mathbb{W}_{(a,m)}^{\mathcal{F}} = 1 + \beta(1 - \mu_{(a,m)}) \sum_{m' \in \mathcal{N}} q_{m'} \mathbb{W}_{(a+r_m, m')}^{\mathcal{F}}, \quad m > t$$

Using (18)

$$\sum_{m' \in \mathcal{N}} q_{m'} \mathbb{W}_{(a+r_m, m')}^{\mathcal{F}} = \frac{\sum_{m'>t'} q_{m'} \mathbb{W}_{(a+r_m, m')}^{\mathcal{F}}}{1 - \beta + \beta \sum_{m'>t'} q_{m'}}$$

where  $(a + r_m, m') \in \mathcal{F}$  for  $m' > t'$

Therefore,

$$\begin{aligned} \mathbb{W}_{(a,m)}^{\mathcal{F}} &= 1 + \frac{\beta(1 - \mu_{(a,m)}) \sum_{m'>t'} q_{m'} \mathbb{W}_{(a+r_m, m')}^{\mathcal{F}}}{1 - \beta + \beta \sum_{m'>t'} q_{m'}} = \\ &= 1 + (1 - \mu_{(a,m)}) X(t') \sum_{m'>t'} q_{m'} \mathbb{W}_{(a+r_m, m')}^{\mathcal{F}}, \quad m > t \end{aligned}$$

This way,

$$\mathbb{W}_{(a,m>t_0)}^{\mathcal{F}} = 1 + (1 - \mu_{(a,m)}) X(t_1) \sum_{m^{(1)}>t_1} q_{m^{(1)}} \mathbb{W}_{(a+r_m, m^{(1)})}^{\mathcal{F}}$$

where  $(a + r_m, m^{(1)}) \in \mathcal{F}$  for  $\forall m^{(1)} > t_1$

$$\begin{aligned} \mathbb{W}_{(a+r_m, m^{(1)}>t_1)}^{\mathcal{F}} &= 1 + (1 - \mu_{(a+r_m, m^{(1)})}) X(t_2) \\ &\cdot \sum_{m^{(2)}>t_2} q_{m^{(2)}} \mathbb{W}_{(a+r_m+r_{m^{(1)}}, m^{(2)})}^{\mathcal{F}} \end{aligned}$$

where  $(a + r_m + r_{m^{(1)}}, m^{(2)}) \in \mathcal{F}$  for  $\forall m^{(2)} > t_2$

This recursion may happen until at least one of the following statements is true:

- $(a+r_m+r_{m^{(1)}}+\dots+r_{m^{(k+1)}}, m^{(k+2)}) \notin \mathcal{F}$  for  $\forall m^{(k+2)}$
- $(a+r_m+r_{m^{(1)}}+\dots+r_{m^{(k+1)}}, m^{(k+2)}) \notin \mathcal{S}$  for  $\forall m^{(k+2)}$
- $\beta^{k+2} = 0$  when  $k+2 \rightarrow \infty$ , which results in  $X(t_{k+2}) = 0$

and consequently,  $\mathbb{W}_{(a+r_m+r_{m^{(1)}}+\dots+r_{m^{(k)}, m^{(k+1)}})}^{\mathcal{F}} = 1$ .

Applying the previous recursion,  $\mathbb{W}_{(a,m>t_0)}^{\mathcal{F}}$  can be more suitably written as

$$\begin{aligned} \mathbb{W}_{(a,m>t_0)}^{\mathcal{F}} &= 1 + (1 - \mu_{(a,m)}) X(t_1) \sum_{m^{(1)}>t_1} q_{m^{(1)}} \\ &\cdot \left( 1 + (1 - \mu_{(a+r_m, m^{(1)})}) X(t_2) \sum_{m^{(2)}>t_2} q_{m^{(2)}} \right. \\ &\cdot \left( 1 + (1 - \mu_{(a+r_m+r_{m^{(1)}}, m^{(2)})}) X(t_3) \sum_{m^{(3)}>t_3} q_{m^{(3)}} \right. \\ &\cdot \left( 1 + (1 - \mu_{(a+r_m+r_{m^{(1)}}+r_{m^{(2)}}, m^{(3)})}) X(t_4) \sum_{m^{(4)}>t_4} q_{m^{(4)}} \right. \\ &\cdot \dots \left. \left. \left. \left. \sum_{m^{(k+1)}>t_{k+1}} q_{m^{(k+1)}} \dots \right) \right) \right) \end{aligned} \quad (21)$$

Manipulating (21) we can more properly express  $\mathbb{W}_{(a,m>t_0)}^{\mathcal{F}}$  as (19), where  $I(a)$  is the last iteration for which the recursion condition holds. Analogously, using the same procedure for the computation of the expression  $\mathbb{W}_{(a,m>t_0)}^{\mathcal{F}}$ , we achieve expression (20) for reward components, where  $L(a, i) = 1$ , excepting the case of  $I(a)$  is due to for  $\forall m^{(i+2)} (a + r_m + \sum_{k=1}^{I(a)+1} r_{m^{(k)}, m^{(i+2)}}) \notin \mathcal{S}$ , that  $L(a, I(a)) = 1 - \beta$ . ■