

On the Optimal Scheduling of Independent, Symmetric and Time-Sensitive Tasks

Fabio Iannello¹, Osvaldo Simeone² and Umberto Spagnolini³

Abstract

Consider a discrete-time system in which a *centralized controller* (CC) is tasked with assigning at each time interval (or *slot*) K resources (or *servers*) to K out of $M \geq K$ *nodes*. When assigned a server, a node can execute a *task*. The tasks are independently generated at each node by stochastically symmetric and memoryless random processes and stored in a finite-capacity *task queue*. Moreover, they are *time-sensitive* in the sense that within each slot there is a non-zero probability that a task expires before being scheduled. The scheduling problem is tackled with the aim of maximizing the number of tasks completed over time (or the *task-throughput*) under the assumption that the CC has no direct access to the state of the task queues. The scheduling decisions at the CC are based on the outcomes of previous scheduling commands, and on the known statistical properties of the task generation and expiration processes.

Based on a Markovian modeling of the task generation and expiration processes, the CC scheduling problem is formulated as a partially observable Markov decision process (POMDP) that can be cast into the framework of restless multi-armed bandit (RMAB) problems. When the task queues are of capacity one, the optimality of a myopic (or greedy) policy is proved. It is also demonstrated that the MP coincides with the Whittle index policy. For task queues of arbitrary capacity instead, the myopic policy is generally suboptimal, and its performance is compared with an upper bound obtained through a relaxation of the original problem.

Overall, the settings in this paper provide a rare example where a RMAB problem can be explicitly solved, and in which the Whittle index policy is proved to be optimal.

^{1,3}F. Iannello (corresponding author: iannello@elet.polimi.it, phone: +390223993604, fax +390223993413) and U. Spagnolini (spagnoli@elet.polimi.it) are with Politecnico di Milano, P.zza L. da Vinci 32, 20133 Milan, Italy.^{1,2}F. Iannello and O. Simeone (osvaldo.simeone@njit.edu) are with the CWCSPP, New Jersey Institute of Technology, U. Heights, Newark, NJ, 07102, USA

I. INTRODUCTION AND SYSTEM MODEL

The problem of scheduling concurrent tasks under resource constraints finds applications in a variety of fields including communication networks [1], distributed computing [2] and virtual machine scenarios [3]. In this paper we consider a specific instance of this general problem in which a *centralized controller* (CC) is tasked with assigning at each time interval (or *slot*) K resources, referred to as *servers*, to K out of $M \geq K$ nodes as shown in Fig. 1. A server can complete a single task per slot and can be assigned to one node per time interval. The tasks are generated at the M nodes by stochastically symmetric, independent and memoryless random processes. The tasks are stored by each node in a finite-capacity *task queue*, and they are *time-sensitive* in the sense that at each slot there is a non-zero probability that a task expires before being completed successfully. It is assumed that the CC has no direct access to the node queues, and thus it is not fully informed of their actual states. Instead, the scheduling decision is based on the outcomes of previous scheduling commands, and on the statistical knowledge of the task generation and expiration processes. If a server is assigned to a node with an empty queue, it remains idle for the whole slot. The purpose here is thus to pair servers to nodes so as to maximize the average number of successfully completed tasks within either a finite or infinite number of slots (*horizon*), which we refer to as *task-throughput*, or simply throughput.

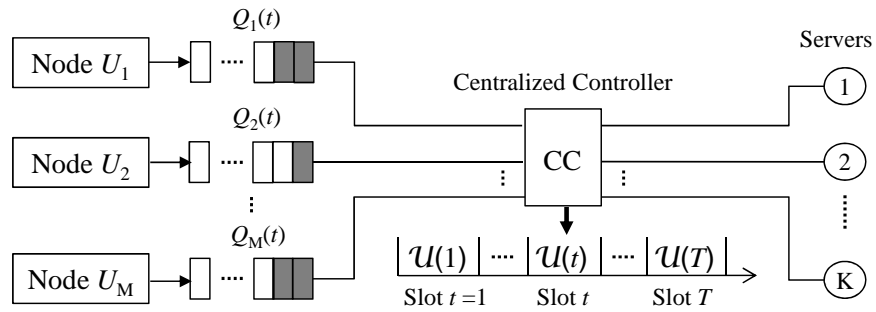


Figure 1. The centralized controller (CC) assigns K resources (servers) to K out of $M \geq K$ nodes to complete their tasks in each slot t . The tasks of node U_i at slot t are stored in a task queue $Q_i(t)$.

A. Markov Formulation

We now introduce the stochastic model that describes the evolution of the task queues across slots. In this section we consider task queues of capacity one (see Sec. V for capacity larger than one), where $Q_i(t) \in \{0, 1\}$ denotes the number of tasks in the queue of node U_i , for $i \in \{1, \dots, M\}$. The stochastic evolution of queue $Q_i(t)$ is shown in Fig. 2 as a function of the scheduling decision $\mathcal{U}(t)$, which consists in the assignment at each slot t of the K servers to a subset $\mathcal{U}(t) \subseteq \{U_1, \dots, U_M\}$ of K nodes, with $|\mathcal{U}(t)| = K$.

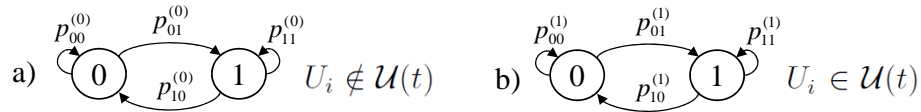


Figure 2. Markov model for the evolution of the state of the task queue $Q_i(t) \in \{0, 1\}$, when the node U_i : a) is not scheduled in slot t (i.e., $U_i \notin \mathcal{U}(t)$); b) is scheduled in slot t (i.e., $U_i \in \mathcal{U}(t)$).

At each slot, node U_i can be either scheduled ($U_i \in \mathcal{U}(t)$) or not ($U_i \notin \mathcal{U}(t)$). If U_i is not scheduled (i.e., $U_i \notin \mathcal{U}(t)$, see Fig. 2-a)) and there is a task in its queue (i.e., $Q_i(t) = 1$), then the task expires with probability (w.p.) $p_{10}^{(0)} = \Pr[Q_i(t+1) = 0 | Q_i(t) = 1, U_i \notin \mathcal{U}(t)]$, while it remains in the queue w.p. $p_{11}^{(0)} = 1 - p_{10}^{(0)}$. Instead, if node U_i is scheduled (i.e., $U_i \in \mathcal{U}(t)$, see Fig. 2-b)) and $Q_i(t) = 1$, its task is completed successfully and its queue in the next slot is either empty or full w.p. $p_{10}^{(1)} = \Pr[Q_i(t+1) = 0 | Q_i(t) = 1, U_i \in \mathcal{U}(t)]$ and $p_{11}^{(1)} = 1 - p_{10}^{(1)}$, respectively. Probability $p_{11}^{(1)}$ accounts for the possible arrival of a new task. If $Q_i(t) = 0$ the probabilities of receiving a new task when U_i is not scheduled and scheduled are $p_{01}^{(0)} = \Pr[Q_i(t+1) = 1 | Q_i(t) = 0, U_i \notin \mathcal{U}(t)]$ and $p_{01}^{(1)} = \Pr[Q_i(t+1) = 1 | Q_i(t) = 0, U_i \in \mathcal{U}(t)]$, respectively, while the probabilities of receiving no task are $p_{00}^{(0)} = 1 - p_{01}^{(0)}$ and $p_{00}^{(1)} = 1 - p_{01}^{(1)}$, respectively.

B. Related Work and Contributions

In this work we assume that the CC has no direct access to the state of the task queues $Q_1(t), \dots, Q_M(t)$, while it knows the transitions probabilities $p_{xy}^{(u)}$, with $x, y, u \in \{0, 1\}$, and the outcomes of previously scheduled tasks. The scheduling problem is thus formalized as a

partially observable Markov decision process (POMDP) [4], and then cast into a restless multi-armed bandit (RMAB) problem [5]. A RMAB is constituted by a set of arms (the queues in our model), a subset of which needs to be activated (or scheduled) in each slot by the controller.

To elaborate, we assume that the transition probabilities of the Markov chains in Fig. 2, the number of nodes M and servers K are such that

$$m = M/K, \text{ is integer, and} \quad (1a)$$

$$p_{11}^{(1)} \leq p_{01}^{(1)} \leq p_{01}^{(0)} \leq p_{11}^{(0)}. \quad (1b)$$

Assumption (1a) states that the ratio $m = M/K$ between the numbers M of nodes and K of servers is an integer, generalizing the single-server case ($K = 1$). Proving the results provided later in this paper for the case of non-integer m remains an open problem. Assumption (1b) is motivated as follows. The inequality $p_{11}^{(1)} \leq p_{01}^{(1)}$ imposes that the probability that a new task arrives when the task queue is full and the node is scheduled ($p_{11}^{(1)}$) is no larger than when the task queue is empty ($p_{01}^{(1)}$). This applies, e.g., when the arrival of a new task is independent on the queue's state and scheduling decisions (i.e., $p_{11}^{(1)} = p_{01}^{(1)}$), or when a new task is not accepted when the queue is full, i.e., $p_{11}^{(1)} = 0$. Inequality $p_{01}^{(1)} \leq p_{01}^{(0)}$ applies, e.g., when the task generation process does not depend on the queue's state and on the scheduling decisions, so that $p_{01}^{(1)} = p_{01}^{(0)}$, or when a new task cannot be accepted while the node is scheduled even if the queue is empty ($p_{01}^{(1)} = 0$). Inequality $p_{01}^{(0)} \leq p_{11}^{(0)}$ indicates that, when a node is not scheduled, the probability $p_{01}^{(0)}$ that its task queue is full in the next slot, given that it is currently empty, is smaller than the probability $p_{11}^{(0)}$ that the task queue is full in the next slot given that it is currently full. This applies, e.g., when the task generation and expiration processes are independent of each other.

Main Contributions: When the task queues are of capacity one, and under assumptions (1), we first show that the *myopic policy* (MP) for the RMAB at hand is a round robin (RR) strategy that: *i*) re-numbers the nodes in a decreasing order according to the initial probability that their respective task queue is full; and then *ii*) schedules the nodes periodically in group of K by

exploiting the initial ordering. The MP is then proved to be throughput-optimal. We then show that, for the special case in which $p_{01}^{(0)} = p_{01}^{(1)}$ and $p_{10}^{(0)} = p_{11}^{(1)} = 0$, the MP coincides with the Whittle index policy, which is a generally suboptimal index strategy for RMAB problems [6]. Finally, we extend the model of Sec. I-A to queues with an arbitrary capacity C . Characterizing optimal policies for $C > 1$ is significantly more complicated than the case of $C = 1$. Hence, inspired by the optimality of the MP for $C = 1$, we compare the performance of the MP for $C > 1$, with an upper bound based on a relaxation of the scheduling constraints of the original RMAB problem [6]. It is recalled that the results in this paper represent a rare case in which the optimal policy for a RMAB can be found explicitly [5].

Related Work: The work in this paper is related to the works [7], [8], in which a RMAB problem similar to the one in this paper is addressed. However, the main difference between our RMAB and the one in [7], [8] is the evolution of the arms across slots. In particular, in our RMAB, each arm evolves across a slot depending on the scheduling decision taken by the controller, while in [7], [8], the evolution of the arms does not depend on the scheduling decision. The transition probabilities for the RMAB in [7], [8] are thus equivalent to setting $p_{01}^{(0)} = p_{01}^{(1)}$ and $p_{11}^{(0)} = p_{11}^{(1)}$ in the Markov chains of Fig. 2. For instance, our model applies to scenarios in which the arms are, e.g., data queues, where each arm draws a data packet from its queue only when scheduled. Instead, the model in [7], [8] applies to scenarios in which the arms are, e.g., communication channels, whose quality evolves across slots regardless whether they are selected for transmission or not.

In [7] it is shown that the MP is optimal for $p_{01}^{(0)} = p_{01}^{(1)} \leq p_{11}^{(0)} = p_{11}^{(1)}$ with $K = 1$, while [8] extends this result to an arbitrary K . The work [7] also demonstrates that the MP is not generally optimal in the case $p_{01}^{(0)} = p_{01}^{(1)} \geq p_{11}^{(0)} = p_{11}^{(1)}$. Finally, paper [9] proves the optimality of the Whittle index policy for $p_{11}^{(0)} = p_{11}^{(1)} \leq p_{01}^{(0)} = p_{01}^{(1)}$. We emphasize that neither our model nor the one considered in [7], [8] subsumes the other, and the results here and in the mentioned previous works should be considered as complementary.

Notation: Vectors are denoted in bold, while the corresponding non-bold letters denote the vectors components. Given a vector $\mathbf{x} = [x_1, \dots, x_M]$ and a set $\mathcal{S} = \{i_1, \dots, i_K\} \subseteq \{1, \dots, M\}$ of cardinality $K \leq M$, we define vector $\mathbf{x}_{\mathcal{S}} = [x_{i_1}, \dots, x_{i_K}]$, where $i_1 \leq \dots \leq i_K$. A function $f(\mathbf{x})$ of vector \mathbf{x} is also denoted as $f(x_1, \dots, x_M)$ or as $f(x_1, \dots, x_l, \mathbf{x}_{\{l+1, \dots, M\}})$ for some $1 \leq l \leq M$, or similar notations depending on the context. Given a set \mathcal{A} and a subset $\mathcal{B} \subseteq \mathcal{A}$, \mathcal{B}^c represents the complement of \mathcal{B} in \mathcal{A} .

II. PROBLEM FORMULATION

Here we formalize the scheduling problem of Sec. I (see Fig. 1), in which the task generation and expiration processes are modeled, independently at each node, by the Markov models of Sec. I-A with queues of capacity one. Extension to task queues of arbitrary capacity is addressed in Sec. V.

A. Problem Definition

The scheduling problem at the CC is addressed in a finite-horizon scenario over slots $t \in \{1, \dots, T\}$. Let $\mathbf{Q}(t) = [Q_1(t), \dots, Q_M(t)]$ be the vector collecting the states of the task queue at slot t . At slot $t = 1$, the CC is only aware of the initial probability distribution $\boldsymbol{\omega}(1) = [\omega_1(1), \dots, \omega_M(1)]$ of $\mathbf{Q}(1)$, whose i th entry is $\omega_i(1) = \Pr[Q_i(1) = 1]$. Thus, the subset $\mathcal{U}(1)$ of $|\mathcal{U}(1)| = K$ nodes scheduled at slot $t = 1$ is chosen as a function of the initial distribution $\boldsymbol{\omega}(1)$ only. For any node $U_i \in \mathcal{U}(t)$ scheduled at slot t , an *observation* is made available to the CC at the end of the slot, while no observations are available for non-scheduled nodes $U_i \notin \mathcal{U}(t)$. Specifically, if $Q_i(t) = 1$ and $U_i \in \mathcal{U}(t)$, the task of U_i is served within slot t , and the CC observes that $Q_i(t) = 1$. Conversely, if $Q_i(t) = 0$ and $U_i \in \mathcal{U}(t)$, no tasks are completed and the CC observes that $Q_i(t) = 0$. We define $\mathcal{O}(t) = \{Q_i(t) : U_i \in \mathcal{U}(t)\}$ as the set of (new) observations available at the CC at the end of slot t . At time t , the CC hence knows the history of all decisions and previous observations and the initial distribution $\boldsymbol{\omega}(1)$, namely $\mathcal{H}(t) = \{\mathcal{U}(1), \dots, \mathcal{U}(t-1), \mathcal{O}(1), \dots, \mathcal{O}(t-1), \boldsymbol{\omega}(1)\}$, with $\mathcal{H}(1) = \{\boldsymbol{\omega}(1)\}$.

Since the CC has only partial information about the system state $\mathbf{Q}(t)$, through $\mathcal{O}(t)$, the scheduling problem at hand can be modeled as a POMDP. It is well-known that a sufficient statistics for taking decisions in such POMDP is given by the probability distribution of $\mathbf{Q}(t)$ conditioned on the history $\mathcal{H}(t)$ [10], referred to as *belief*, and represented by the vector $\boldsymbol{\omega}(t) = [\omega_1(t), \dots, \omega_M(t)]$, with i th entry given by

$$\omega_i(t) = \Pr [Q_i(t) = 1 | \mathcal{H}(t)]. \quad (2)$$

Since the belief $\boldsymbol{\omega}(t)$ fully summarizes the entire history $\mathcal{H}(t)$ of past actions and observations [10], a *scheduling policy* $\pi = [\mathcal{U}^\pi(1), \dots, \mathcal{U}^\pi(T)]$ is defined as a collection of functions $\mathcal{U}^\pi(t)$ that map the belief $\boldsymbol{\omega}(t)$ to a subset $\mathcal{U}(t)$ of $|\mathcal{U}(t)| = K$ nodes, i.e., $\mathcal{U}^\pi(t): \boldsymbol{\omega}(t) \rightarrow \mathcal{U}(t)$. We will refer to $\mathcal{U}^\pi(t)$ as the subset of scheduled nodes, even though, strictly speaking, it is the mapping function defined above. The transition probabilities over the belief space are derived in Sec. II-B.

The *immediate reward* $R(\boldsymbol{\omega}, \mathcal{U})$, accrued by the CC when the belief vector is $\boldsymbol{\omega}$ and action \mathcal{U} is taken, measures the average number of tasks completed within the current slot, and it is

$$R(\boldsymbol{\omega}, \mathcal{U}) = \sum_{U_i \in \mathcal{U}} \omega_i. \quad (3)$$

Notice that $R(\boldsymbol{\omega}, \mathcal{U}) \leq K$ since there are only K servers.

The *throughput* measures the average number of tasks completed over the slots $\{1, \dots, T\}$ that, by exploiting (3) and under policy π , is given by

$$V_1^\pi(\boldsymbol{\omega}(1)) = \sum_{t=1}^T \beta^{t-1} \mathbb{E}^\pi [R(\boldsymbol{\omega}(t), \mathcal{U}^\pi(t)) | \boldsymbol{\omega}(1)]. \quad (4)$$

In (4), the expectation $\mathbb{E}^\pi[\cdot | \boldsymbol{\omega}(1)]$, under policy π for initial belief $\boldsymbol{\omega}(1)$, is intended with respect to the distribution of the Markov process $\boldsymbol{\omega}(t)$, as obtained from the transition probabilities to be derived in Sec. II-B. For generality, the definition (4) includes a discount factor $0 \leq \beta \leq 1$ [7], while the infinite horizon scenario (i.e., $T \rightarrow \infty$) will be discussed in Sec. III-C.

The goal is to find a policy $\pi^* = [\mathcal{U}^*(1), \dots, \mathcal{U}^*(T)]$ that maximizes the throughput (4) so that

$$V_1^*(\boldsymbol{\omega}(1)) = V_1^{\pi^*}(\boldsymbol{\omega}(1)) = \max_{\pi} V_1^{\pi}(\boldsymbol{\omega}(1)), \quad \text{with } \pi^* = \arg \max_{\pi} V_1^{\pi}(\boldsymbol{\omega}(1)) \quad (5)$$

B. Transition Probabilities

The belief transition probabilities, given decision $\mathcal{U}(t) = \mathcal{U}$ and $\boldsymbol{\omega}(t) = \boldsymbol{\omega} = [\omega_1, \dots, \omega_M]$, are

$$p_{\boldsymbol{\omega}\boldsymbol{\omega}'}^{(\mathcal{U})} = \Pr[\boldsymbol{\omega}(t+1) = \boldsymbol{\omega}' | \boldsymbol{\omega}(t) = \boldsymbol{\omega}, \mathcal{U}(t) = \mathcal{U}] = \prod_{i=1}^M \Pr[\omega_i(t+1) = \omega'_i | \omega_i(t) = \omega_i, \mathcal{U}(t) = \mathcal{U}], \quad (6)$$

where $\boldsymbol{\omega}(t+1) = \boldsymbol{\omega}' = [\omega'_1, \dots, \omega'_M]$, while the distribution of entry $\omega_i(t+1)$ is (see Fig. 2)

$$\Pr[\omega_i(t+1) = \omega'_i | \omega_i(t) = \omega_i, \mathcal{U}(t) = \mathcal{U}] = \begin{cases} \omega_i & \text{if } \omega'_i = p_{11}^{(1)} \text{ and } U_i \in \mathcal{U} \\ (1 - \omega_i) & \text{if } \omega'_i = p_{01}^{(1)} \text{ and } U_i \in \mathcal{U} \\ 1 & \text{if } \omega'_i = \tau_0^{(1)}(\omega_i) \text{ and } U_i \notin \mathcal{U} \end{cases}, \quad (7)$$

where we have defined the deterministic function

$$\tau_0^{(1)}(\omega) = \Pr[Q_i(t+1) = 1 | \omega_i(t) = \omega, U_i \notin \mathcal{U}(t)] = \omega p_{11}^{(0)} + (1 - \omega) p_{01}^{(0)} = \omega \delta_0 + p_{01}^{(0)} \quad (8)$$

to indicate the next slot's belief when U_i is not scheduled ($U_i \notin \mathcal{U}(t)$), with $\delta_0 = p_{11}^{(0)} - p_{01}^{(0)} \geq 0$ due to inequalities (1b). Eq. (8) follows from Fig. 2-a), since the next slot's belief is either $p_{11}^{(0)}$ if $Q_i(t) = 1$ (w.p. ω) or $p_{01}^{(0)}$ if $Q_i(t) = 0$ (w.p. $(1 - \omega)$). A generalization of function $\tau_0^{(1)}(\omega)$ that computes the belief $\omega_i(t+k)$ of node U_i when it is not scheduled for k successive slots, e.g., slots $\{t, \dots, t+k-1\}$, and $\omega_i(t) = \omega$, can be obtained as

$$\tau_0^{(k)}(\omega) = \Pr[B_i(t+k) = 1 | \omega_i(t) = \omega, U_i \notin \mathcal{U}(t), \dots, U_i \notin \mathcal{U}(t+k-1)] = \omega \delta_0^k + p_{01}^{(0)} \frac{1 - \delta_0^k}{1 - \delta_0}. \quad (9)$$

Eq. (9) can be obtained recursively from (8) as $\tau_0^{(k)}(\omega) = \tau_0^{(1)}(\tau_0^{(k-1)}(\omega))$, for all $k \geq 1$, with $\tau_0^{(0)}(\omega) = \omega$.

Under assumptions (1b), it is easy to verify that function (8) satisfies the inequalities

$$p_{11}^{(1)} \leq p_{01}^{(1)} \leq \tau_0^{(1)}(\omega), \text{ for all } 0 \leq \omega \leq 1, \text{ and} \quad (10)$$

$$\tau_0^{(1)}(\omega) \leq \tau_0^{(1)}(\omega'), \text{ for all } \omega \leq \omega' \text{ with } 0 \leq \omega, \omega' \leq 1. \quad (11)$$

The inequalities in (10) guarantee that the belief of a non-scheduled node is always larger than that of a scheduled one, while the inequality (11) says that the belief ordering of two non-scheduled nodes is maintained across a slot. These inequalities play a crucial role in the analysis below.

C. Optimality Equations

The dynamic programming (DP) formulation of problem (5) (see e.g., [11]) allows to express the throughput recursively over the horizon $\{t, \dots, T\}$, under policy π and initial belief ω , as

$$V_t^\pi(\omega) = \sum_{j=t}^T \beta^{j-t} \mathbb{E}^\pi [R(\omega(j), \mathcal{U}^\pi(j)) | \omega(t) = \omega] = R(\omega, \mathcal{U}^\pi(t)) + \beta \sum_{\omega'} p_{\omega\omega'}^{(\mathcal{U}^\pi)} V_{t+1}^\pi(\omega'), \quad (12)$$

where $V_t^\pi(\cdot) = 0$ for $t > T$. The DP optimality conditions (or *Bellman equations*) are then expressed in terms of the *value function* $V_t^*(\omega) = \max_{\pi} V_t^\pi(\omega)$, which represents the optimal throughput in the interval $\{t, \dots, T\}$, and it is given by

$$V_t^*(\omega) = \max_{\mathcal{U}(t) = \mathcal{U} \subseteq \{U_1, \dots, U_M\}} \left\{ R(\omega, \mathcal{U}) + \beta \sum_{\omega'} p_{\omega\omega'}^{(\mathcal{U})} V_{t+1}^*(\omega') \right\}. \quad (13)$$

Note that, since the nodes are stochastically equivalent, the value function (13) only depends on the numerical values of the entries of the belief vector ω regardless of the way it is ordered. Finally, an *optimal policy* $\pi^* = [\mathcal{U}^*(1), \dots, \mathcal{U}^*(T)]$ (see (5)) is such that $\mathcal{U}^*(t)$ attains the maximum in the condition (13) for $t \in \{1, 2, \dots, T\}$.

III. OPTIMALITY OF THE MYOPIC POLICY

We now define the myopic policy (MP) and show that, under assumptions (1), it is a round-robin (RR) policy that schedules the nodes periodically and that it is optimal for problem (5).

A. The Myopic Policy is Round-Robin

The MP $\pi^{MP} = \{\mathcal{U}^{MP}(1), \dots, \mathcal{U}^{MP}(T)\}$, with throughput $V_t^{MP}(\cdot)$, is the greedy policy that schedules at each slot the K nodes with the largest beliefs so as to maximize the immediate reward (3), that is, we have

$$\mathcal{U}^{MP}(t) = \arg \max_{\mathcal{U}} R(\boldsymbol{\omega}(t), \mathcal{U}) = \arg \max_{\mathcal{U}} \sum_{U_i \in \mathcal{U}} \omega_i(t). \quad (14)$$

Proposition 1. Under assumptions (1), the MP (14), given an initial belief $\boldsymbol{\omega}'(1)$, is a RR policy that operates as follows: **1)** Sort vector $\boldsymbol{\omega}'(1)$ in a decreasing order to obtain $\boldsymbol{\omega}(1) = [\omega_1(1), \dots, \omega_M(1)]$ such that $\omega_1(1) \geq \dots \geq \omega_M(1)$. Re-number the nodes so that U_i has belief $\omega_i(1)$; **2)** Divide the nodes into m groups of K nodes each, so that the g th group \mathcal{G}_g , $g \in \{1, \dots, m\}$, contains all nodes U_i such that $g = \lfloor \frac{i-1}{K} \rfloor + 1$, namely: $\mathcal{G}_1 = \{U_1, \dots, U_K\}$, $\mathcal{G}_2 = \{U_{K+1}, \dots, U_{2K}\}$, and so on; **3)** Schedule the groups in a RR fashion with period m slots, so that groups $\mathcal{G}_1, \dots, \mathcal{G}_m, \mathcal{G}_1, \dots$ are sequentially scheduled at slot $t = 1, \dots, m, m+1, \dots$ and so on.

Proof: According to (14), the first scheduled set is $\mathcal{U}^{MP}(1) = \mathcal{G}_1 = \{U_1, U_2, \dots, U_K\}$. The beliefs are then updated through (7). Recalling (10), the scheduled nodes, in \mathcal{G}_1 , have their belief updated to either $p_{11}^{(1)}$ or $p_{01}^{(1)}$, which are both smaller than the belief of any non-scheduled node in $\{U_1, \dots, U_M\} \setminus \mathcal{G}_1$. Moreover, the ordering of the non-scheduled nodes' beliefs is preserved due to (11). Hence, the second scheduled group is $\mathcal{U}^{MP}(2) = \mathcal{G}_2$, the third is $\mathcal{U}^{MP}(3) = \mathcal{G}_3$, and so on. This proves that the MP, upon an initial ordering of the beliefs, is a RR policy. ■

We emphasize that the MP sorts the beliefs of the nodes only at the first slot in which it is operated, and then it keeps scheduling the groups of nodes according to their initial ordering, without requiring to recalculate the beliefs.

B. Optimality of the Myopic Policy

We now prove the optimality of the MP by showing that it satisfies the Bellman equations (13). To start with, let us consider a RR policy π^{RR} that operates according to steps **2)** and

3) of Proposition 1 (i.e., without re-ordering the initial belief), and let its throughput (12) be denoted by $V_t^{RR}(\boldsymbol{\omega})$. Note that, when the initial belief $\boldsymbol{\omega}$ is ordered so that $\omega_1 \geq \dots \geq \omega_M$, then $V_t^{RR}(\boldsymbol{\omega}) = V_t^{MP}(\boldsymbol{\omega})$. Based on backward induction arguments similarly to [7], [8], the following lemma establishes a sufficient condition for the optimality of the MP.

Lemma 2. Assume that the MP is optimal at slot $t + 1, \dots, T$, i.e., it satisfies (13). To show that the MP is optimal also at slot t it is sufficient to prove the inequality

$$V_t^{RR}(\boldsymbol{\omega}_S, \boldsymbol{\omega}_{S^c}) \leq V_t^{MP}(\boldsymbol{\omega}_S, \boldsymbol{\omega}_{S^c}) = V_t^{RR}(\omega_1, \omega_2, \dots, \omega_M), \quad \text{for all } \omega_1 \geq \omega_2 \geq \dots \geq \omega_M \quad (15)$$

and all sets $\mathcal{S} \subseteq \{1, \dots, M\}$ of K elements, with the elements in $\boldsymbol{\omega}_{S^c}$ decreasingly ordered.

Proof: Since the MP is optimal from $t + 1$ onward by assumption, it is sufficient to show that scheduling K nodes with arbitrary beliefs at slot t and then following the MP from slot $t + 1$ on is no better than following the MP immediately at slot t . The performance of the former policy is given by the left-hand side (LHS) of (15). In fact $V_t^{RR}(\boldsymbol{\omega}_S, \boldsymbol{\omega}_{S^c})$, for any set \mathcal{S} , represents the throughput of a policy that schedules the K nodes with beliefs $\boldsymbol{\omega}_S$ at slot t , and then operates as the MP from $t + 1$ onward, since beliefs in $\boldsymbol{\omega}_{S^c}$ are decreasingly ordered. The MP's performance is instead given by the right-hand side (RHS) of (15). Note that, for $t = T$, it is immediate to verify that the MP is optimal. This concludes the proof. ■

Theorem 3. Under assumptions (1) the MP is optimal for problem (5), so that $\pi^{MP} = \pi^*$.

Proof: To start with, we first prove in Appendix A that the inequality

$$V_t^{RR}(\omega_1, \dots, \omega_j, y, x, \dots, \omega_M) \leq V_t^{RR}(\omega_1, \dots, \omega_j, x, y, \dots, \omega_M) \quad (16)$$

holds for any $x \geq y$, with $0 \leq j \leq M - 2$, and for all $t \in \{1, \dots, T\}$ and beliefs ω_k (not necessarily ordered), with $k \in \{1, \dots, M\}$. Inequality (16) for $j = 0$ is intended as $V_t^{RR}(y, x, \dots, \omega_M) \leq V_t^{RR}(x, y, \dots, \omega_M)$. If (16) holds, then inequality (15) is satisfied for all $\omega_1 \geq \dots \geq \omega_M$ and all subsets $\mathcal{S} \subseteq \{1, \dots, M\}$ of K elements. In fact, (16) states that the throughput of the RR policy

never increases when, for any pair of adjacent nodes, the one with the smallest belief of the pair is scheduled first. Hence, by starting from the RHS of (15) (i.e., $V_t^{RR}(\omega_1, \omega_2, \dots, \omega_M)$) and by applying a convenient number of successive switchings between pair of adjacent elements of vector $[\omega_1, \omega_2, \dots, \omega_M]$ to achieve $[\omega_{\mathcal{S}}, \omega_{\mathcal{S}^c}]$, for any \mathcal{S} , we can obtain a cascade of inequalities through (16) (one for each switching), which guarantees that (15) holds. By Lemma 2 this is sufficient to prove that the MP is optimal, since the inequality (15) holds for any arbitrary t . ■

C. Extension to the Infinite-Horizon Case

We now briefly describe the extension of problem (5) to the infinite-horizon case, for which the throughput under policy π and its optimal value are given by (see e.g., [7])

$$V^\pi(\boldsymbol{\omega}(1)) = \sum_{t=1}^{\infty} \beta^{t-1} \mathbb{E}^\pi [R(\boldsymbol{\omega}(t), \mathcal{U}^\pi(t)) | \boldsymbol{\omega}(1)], \quad \text{and} \quad V^*(\boldsymbol{\omega}(1)) = \max_{\pi} V^\pi(\boldsymbol{\omega}(1)), \quad (17)$$

where the optimal policy is $\pi^* = \arg \max_{\pi} V^\pi(\boldsymbol{\omega}(1))$ and $0 \leq \beta < 1$. From standard DP theory [11], the optimal policy π^* is stationary, so that the optimal decision $\mathcal{U}^*(t)$ is a function of the current state $\boldsymbol{\omega}(t)$ only, independently of slot t [11]. By following the same reasoning as in [7, Theorem 3], it can be shown that the optimality of the MP for the finite-horizon setting implies the optimality also for the infinite-horizon scenario. Moreover, by following [7, Theorem 4] it can be shown that the MP is optimal also for the undiscounted average reward criterion (i.e., $V_{avg}^\pi(\boldsymbol{\omega}(1)) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^{\infty} \mathbb{E}^\pi [R(\boldsymbol{\omega}(t), \mathcal{U}^\pi(t)) | \boldsymbol{\omega}(1)]$).

IV. OPTIMALITY OF THE WHITTLE INDEX POLICY

In this section, we briefly review the Whittle index policy for RMAB problems [5], and then focus on the infinite-horizon scenario of Sec. III-C, when conditions (1b) are specialized to

$$0 = p_{11}^{(1)} \leq p_{01}^{(1)} = p_{01}^{(0)} = p_{01} \leq p_{11}^{(0)} = 1, \quad (18)$$

and where the task queues are of capacity one. We show that under the assumption (18) (see Sec. I-B for a discussion on these conditions), the RMAB at hand is indexable and we calculate

its Whittle index in closed-form. We then show that the Whittle index policy is equivalent to the MP, and thus optimal for the problem (17).

We emphasize that, our results provide a rare example [5] in which, as in [9], not only indexability is established, but also the Whittle index is obtained in closed form and the Whittle policy proved to be optimal. It is finally remarked that our proof technique is inspired by [9], but the different system model poses new challenges that require significant work.

A. Whittle Index

The Whittle index policy assigns a numerical value $W(\omega_i)$ to each state ω_i of node U_i , referred to as *index*, to measure how rewarding it is to schedule node U_i in the current slot. The K nodes with the largest index are then scheduled in each slot. As detailed below, the Whittle index is calculated independently for each node, and thus the Whittle index policy is not generally optimal for RMAB problems. Moreover, even the existence of a well-defined Whittle index is not guaranteed [5]. To study the indexability and the Whittle index for the RMAB at hand, we can focus on a restless single-armed bandit (RSAB) model, as defined below [5]. A RSAB is a RMAB with a single arm, in which the only decision that needs to be taken by the CC is whether activating the (single) arm or not (i.e., keep it passive).

1) *RSAB with Subsidy for Passivity*: The Whittle index is based on the concept of *subsidy for passivity*, whereby the CC is given a subsidy $m \in \mathbb{R}$ when the arm is not scheduled. At each slot t , the CC, based on the state $\omega(t)$ of the arm, can decide to activate (or schedule) it, i.e., to set $u(t) = 1$, obtaining an immediate reward $R_m(\omega(t), 1) = \omega(t)$. If, instead, the arm is kept passive, i.e., $u(t) = 0$, a reward $R_m(\omega(t), 0) = m$ equal to the subsidy is accrued. The state $\omega(t)$ evolves through (7), which under (18) and adapted to the simplified notation used here becomes

$$\omega(t+1) = \begin{cases} 0 & \text{w.p. } \omega(t) & \text{if } u(t) = 1 \\ p_{01} & \text{w.p. } (1 - \omega(t)) & \text{if } u(t) = 1 \\ \tau_0^{(1)}(\omega(t)) & \text{w.p. } 1 & \text{if } u(t) = 0 \end{cases} . \quad (19)$$

The throughput, given policy $\pi = \{u^\pi(1), u^\pi(2), \dots\}$ and initial belief $\omega(1)$, is

$$V_m^\pi(\omega(1)) = \sum_{t=1}^{\infty} \beta^{t-1} \mathbb{E}^\pi [R_m(\omega(t), u^\pi(t)) | \omega(1)]. \quad (20)$$

The optimal throughput is $V_m^*(\omega(1)) = \max_\pi V_m^\pi(\omega(1))$, while the optimal policy $\pi^* = \arg \max_\pi V_m^\pi(\omega(1))$ is stationary in the sense that the optimal decisions $u_m^*(\omega) \in \{0, 1\}$ are functions of the belief ω only, independently of slot t [9]. Removing the slot index from the initial belief, the optimal throughput $V_m^*(\omega)$ and the optimal decision $u_m^*(\omega)$ satisfy the following DP optimality equations for the infinite-horizon scenario (see [9])

$$V_m^*(\omega) = \max_{u \in \{0, 1\}} \{V_m(\omega|u)\}, \quad (21)$$

$$\text{and } u_m^*(\omega) = \arg \max_{u \in \{0, 1\}} \{V_m(\omega|u)\}. \quad (22)$$

In (21)-(22) we defined $V_m(\omega|u)$, $u \in \{0, 1\}$, as the throughput (20) of a policy that takes action u at the current slot and then uses the optimal policy $u_m^*(\omega)$ onward, we have

$$V_m(\omega|0) = m + \beta V_m^*(\tau_0^{(1)}(\omega)), \text{ and} \quad (23)$$

$$V_m(\omega|1) = \omega + \beta [\omega V_m^*(0) + (1 - \omega) V_m^*(p_{01})]. \quad (24)$$

2) *Indexability and Whittle Index:* We use the notation of [9] to define indexability and Whittle index for the RSAB at hand. We first define the so called *passive set*

$$\mathcal{P}(m) = \{\omega: 0 \leq \omega \leq 1 \text{ and } u_m^*(\omega) = 0\}, \quad (25)$$

as the set that contains all the beliefs ω for which the passive action is optimal (i.e., all $0 \leq \omega \leq 1$ such that $V_m(\omega|0) \geq V_m(\omega|1)$, see (23)-(24)) under the given subsidy for passivity $m \in \mathbb{R}$. The RMAB at hand is said to be *indexable* if the passive set $\mathcal{P}(m)$, for the associated RSAB problem¹, is monotonically increasing as m increases within the interval $(-\infty, +\infty)$, in the

¹Note that in a RMAB with arms characterized by different statistics this condition must be checked for all arms.

sense that $\mathcal{P}(m') \subseteq \mathcal{P}(m)$ if $m' \leq m$ and $\mathcal{P}(-\infty) = \emptyset$ and $\mathcal{P}(+\infty) = [0, 1]$.

If the RMAB is indexable, the Whittle index $W(\omega)$ for each arm with state ω is the infimum subsidy m such that it is optimal to make the arm passive. Equivalently, the Whittle index $W(\omega)$ is the infimum subsidy m that makes passive and active actions equally rewarding, i.e.,

$$W(\omega) = \inf \{m: u_m^*(\omega) = 0\} = \inf \{m: V_m(\omega|0) = V_m(\omega|1)\}. \quad (26)$$

B. Optimality of the Threshold Policy

Here, we show that the optimal policy $u_m^*(\omega)$ for the RSAB of Sec. IV-A1 is a threshold policy over the belief ω . This is crucial in our proof of indexability of the RMAB at hand given in Sec. IV-D. To this end, we observe that: *i*) function $V_m(\omega|1)$ in (24) is linear over the belief ω ; *ii*) function $V_m(\omega|0) = m + \beta V_m^*(\tau_0^{(1)}(\omega))$ in (23) is convex over ω , since the value function $V_m^*(\omega)$ is convex for the problem at hand (see [9], [10]). We need the following lemma.

Lemma 4. The following inequalities hold:

$$a) \text{ For } 0 \leq m < 1: \quad a.1) V_m(0|1) \leq V_m(0|0) \leq V_m(1|1); \quad a.2) V_m(1|0) \leq V_m(1|1); \quad (27a)$$

$$b) \text{ For } m < 0: \quad b.1) V_m(0|0) \leq V_m(0|1) \leq V_m(1|1); \quad b.2) V_m(1|0) \leq V_m(1|1); \quad (27b)$$

$$c) \text{ For } m \geq 1: \quad c.1) V_m(0|0) \leq V_m(1|1) \leq V_m(0|1); \quad c.2) V_m(1|1) \leq V_m(1|0). \quad (27c)$$

Proof: See Appendix B. ■

Leveraging Lemma 4, we can now establish the optimality of a threshold policy $u_m^*(\omega)$.

Proposition 5. The optimal policy $u_m^*(\omega)$ in (22) for subsidy $m \in \mathbb{R}$ is given by

$$u_m^*(\omega) = \begin{cases} 1, & \text{if } \omega > \omega^*(m) \\ 0, & \text{if } \omega \leq \omega^*(m) \end{cases}, \quad (28)$$

where $\omega^*(m) \in \mathbb{R}$ is the optimal threshold for a given subsidy m . The optimal threshold $\omega^*(m)$ is $0 \leq \omega^*(m) \leq 1$ if $0 \leq m < 1$, while it is arbitrary negative for $m < 0$ and arbitrary greater

than unity for $m \geq 1$. In other words we have $u_m^*(\omega) = 1$ if $m < 0$ and $u_m^*(\omega) = 0$ if $m \geq 1$.

Proof: We start by showing that (28), for $0 \leq m < 1$, satisfies (22) and is thus an optimal policy. To see this, we refer to Fig. 3, where we sketch functions $V_m(\omega|1)$ and $V_m(\omega|0)$ for different values of the subsidy m . From (22), we have that $u_m^*(\omega) = 1$ for all ω such that $V_m(\omega|1) > V_m(\omega|0)$ and $u_m^*(\omega) = 0$ otherwise. For $0 \leq m < 1$, from the inequalities of Lemma 4-a), the linearity of $V_m(\omega|1)$ and the convexity of $V_m(\omega|0)$, it follows that there is only one intersection $\omega^*(m)$ between $V_m(\omega|1)$ and $V_m(\omega|0)$ with $0 \leq \omega^*(m) \leq 1$, as shown in Fig. 3-a). Instead, when $m < 0$, by Lemma 4-b), arm activation is always optimal, that is, $u_m^*(\omega) = 1$, since $V_m(\omega|1) > V_m(\omega|0)$ for any $0 \leq \omega \leq 1$ as shown in Fig. 3-b). Conversely, when $m \geq 1$, by Lemma 4-c), it follows that passivity is always optimal, that is, $u_m^*(\omega) = 0$, since $V_m(\omega|0) \geq V_m(\omega|1)$ for any $0 \leq \omega \leq 1$ as shown in Fig. 3-c). ■

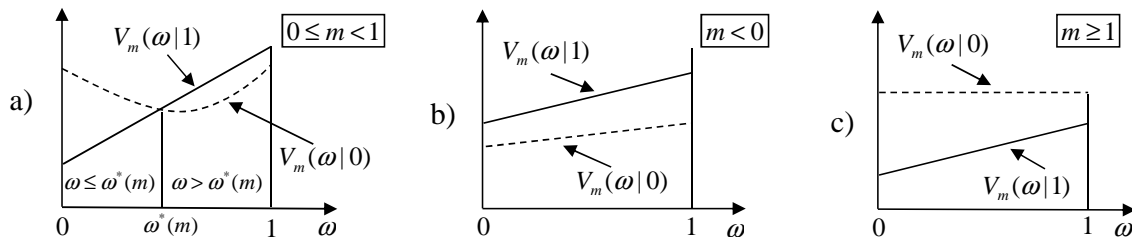


Figure 3. Illustration of the optimality of a threshold policy for different values of the subsidy for passivity m : a) $0 \leq m < 1$; b) $m < 0$; c) $m \geq 1$.

C. Closed-Form Expression of the Value Function

By leveraging the optimality of the threshold policy (28) we derive a closed-form expression of $V_m^*(\omega)$ in (21), being a key step in establishing the RMAB's indexability in Sec. IV-D.

Notice that function $\tau_0^{(k)}(\omega)$ in (9), when specialized to conditions (18), becomes

$$\tau_0^{(k)}(\omega) = 1 - (1 - p_{01})^k (1 - \omega), \quad (29)$$

which is a monotonically increasing function of k , so that $\tau_0^{(k)}(\omega) \geq \tau_0^{(i)}(\omega)$ for any $k \geq i$.

Based on such monotonicity, we can define the average number $L(\omega, \omega')$ of slots it takes for the

belief to become larger than ω' when starting from ω while the arm is kept passive, as

$$L(\omega, \omega') = \min \left\{ k: \tau_0^{(k)}(\omega) > \omega' \right\} = \begin{cases} 0 & \omega > \omega' \\ \left\lfloor \frac{\ln\left(\frac{1-\omega'}{1-\omega}\right)}{\ln(1-p_{01})} \right\rfloor + 1 & \omega \leq \omega' \\ \infty & \omega \leq 1 \leq \omega' \end{cases}. \quad (30)$$

From (30) we have $L(\omega, \omega') = 1$ for $\omega = \omega'$ since, without loss of optimality, we assumed that the passive action is optimal (i.e., $u_m^*(\omega) = 0$) when $V_m(\omega|0) = V_m(\omega|1)$. For $\omega' \geq 1$ instead (according to Proposition 5), the arm is always kept passive and thus $L(\omega, \omega') = \infty$.

Lemma 6. The optimal throughput $V_m^*(\omega)$ in (21) can be written as

$$V_m^*(\omega) = \frac{1 - \beta^{L(\omega, \omega^*(m))}}{1 - \beta} m + \beta^{L(\omega, \omega^*(m))} V_m(\tau_0^{(L(\omega, \omega^*(m)))}(\omega)|1), \quad (31)$$

where $\omega^*(m)$ is the optimal threshold obtained from Proposition 5.

Proof: According to Proposition 5, the optimal policy $u_m^*(\omega)$ keeps the arm passive as long as the current belief is $\omega \leq \omega^*(m)$. Therefore, the arm is kept passive for $L(\omega, \omega^*(m))$ slots, during which a reward $R_m(\omega, 0) = m$ is accrued in each slot. This leads to a total reward within the passivity time given by the following geometric series $\sum_{k=0}^{L(\omega, \omega^*(m))-1} \beta^k m = \frac{1 - \beta^{L(\omega, \omega^*(m))}}{1 - \beta} m$, which corresponds to the first term in the RHS of (31). After $L(\omega, \omega^*(m))$ slots of passivity, the belief becomes larger than the threshold $\omega^*(m)$ and the arm is activated. The contribution to the value function $V(\omega)$ thus becomes $\beta^{L(\omega, \omega^*(m))} V_m(\tau_0^{(L(\omega, \omega^*(m)))}(\omega)|1)$, which is the second term in the RHS of (31). Note that, when $\omega > \omega^*(m)$, activation is optimal, and $V^*(\omega) = V(\omega|1)$. ■

To evaluate $V_m^*(\omega)$ from (31), we only need to calculate $V_m(\omega|1)$ since the other terms, thanks to (30) are explicitly given once $\omega^*(m)$ is obtained from Proposition 5. However, from (24), evaluating $V_m(\omega|1)$ only requires $V_m^*(0)$ and $V_m^*(p_{01})$, which are calculated in the lemma below.

Lemma 7. We have

$$V_m^*(0) = \frac{(m - 2m\beta^{L_m^*} + \beta^{L_m^*}v_m^* - \beta^{L_m^*+1}v_m^* + m\beta^{L_m^*+1} + m\beta^{L_m^*}v_m^* - m\beta^{L_m^*+1}v_m^*)}{(\beta - 1)(\beta^{L_m^*} - \beta^{L_m^*}v_m^* + \beta^{L_m^*+1}v_m^* - 1)} \quad (32a)$$

$$V_m^*(p_{01}) = \frac{(m\beta - m\beta^{L_m^*} + \beta^{L_m^*}v_m^* - \beta^{L_m^*+1}v_m^* + m\beta^{L_m^*+1}v_m^* - m\beta^{L_m^*+2}v_m^*)}{\beta(\beta - 1)(\beta^{L_m^*} - \beta^{L_m^*}v_m^* + \beta^{L_m^*+1}v_m^* - 1)} \quad (32b)$$

where we have defined $L_m^* = L(0, \omega^*(m))$ and $v_m^* = \tau_0^{(L(0, \omega^*(m)))}(0)$.

Proof: By plugging (24) into (31), and evaluating (31) for $\omega = 0$ and $\omega = p_{01}$, we get a linear system in the two unknowns $V_m^*(0)$ and $V_m^*(p_{01})$, which can be solved leading to (32). ■

D. Indexability and Whittle Index

Here, we prove that the RMAB at hand is indexable, we derive the Whittle index in closed form and show that it is equivalent to the MP and thus optimal for the RMAB problem (17).

Theorem 8. a) The RMAB at hand is indexable and b) its Whittle index is

$$W(\omega) = \frac{\left(1 - \beta^{L(0, \omega)} \left(1 - \beta \tau_0^{L(0, \omega)}(0) (1 - \beta) (1 - h)\right)\right) \omega + \beta^{L(0, \omega)} \tau_0^{L(0, \omega)}(0) (1 - \beta) (h\beta + 1)}{(\beta - 1) \left(\beta^{L(0, \omega)} (1 - \beta(1 - h)) \omega - \left(1 + \beta^{L(0, \omega)} \left(\tau_0^{L(0, \omega)}(0) (1 - \beta) + h\beta\right)\right)\right)} \quad (33)$$

Proof: Part **a)**. See Appendix C. Part **b)**. By (26), the Whittle index $W(\omega)$ of state ω is the value of the subsidy m for which activating or not the arm is equally rewarding so that $V_m(\omega|0) = V_m(\omega|1)$. By using (23)-(24) this becomes $\omega + \beta[\omega V_m^*(0) + (1 - \omega)V_m^*(p_{01})] = m + \beta V_m^*(\tau_0^{(1)}(\omega))$. Moreover, since the threshold policy is optimal and $\tau_0^{(1)}(\omega) > \omega$, it follows that, when the belief becomes $\tau_0^{(1)}(\omega)$, it is optimal to activate the arm and thus $V_m^*(\tau_0^{(1)}(\omega)) = V_m(\tau_0^{(1)}(\omega)|1) = \beta \tau_0^{(1)}(\omega) V_m^*(0) + \beta(1 - \tau_0^{(1)}(\omega)) V_m^*(p_{01})$. Plugging this result into $V_m(\omega|0) = V_m(\omega|1)$, along with (32a) and (32b), leads to (33), which concludes the proof. ■

It can be show that the Whittle index $W(\omega)$ in (33) is an increasing function of ω . Therefore, since the Whittle policy selects the K arms with the largest index at each slot, we have:

Corollary 9. The Whittle index policy is equivalent to the MP and is thus optimal.

V. EXTENSION TO TASK QUEUES OF ARBITRARY CAPACITY $C > 1$

The problem of characterizing the optimal policies when $C > 1$ is significantly more complicated than for $C = 1$ and is left open by this work. Moreover, since the dimension of the state space of the belief MDP grows with C , even the numerical computation of the optimal policies is quite cumbersome. Due to these difficulties, here we compare the performance of the MP, inspired by its optimality for $C = 1$, with a performance upper bound obtained following the relaxation approach of [6].

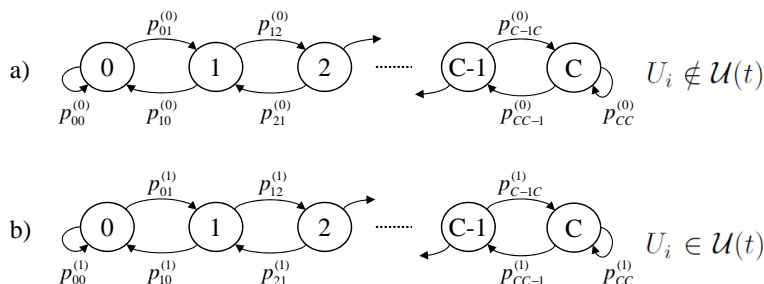


Figure 4. Markov model for the evolution of the queue $Q_i(t)$, of arbitrary capacity C , when the node U_i : a) is not scheduled in slot t (i.e., $U_i \notin \mathcal{U}(t)$); b) is scheduled in slot t (i.e., $U_i \in \mathcal{U}(t)$).

A. System Model and Myopic Policy

Each node U_i has a task queue $Q_i(t) \in \{0, 1, \dots, C\}$ of capacity C . We consider the Markov model of Fig. 4 for the task generation and expiration processes at each node (cf. Sec. I-A). The transition probabilities between queue states when node U_i is not scheduled are $p_{xy}^{(0)} = \Pr[Q_i(t+1) = y | Q_i(t) = x, U_i \notin \mathcal{U}(t)]$, whereas when U_i is scheduled we have $p_{xy}^{(1)} = \Pr[Q_i(t+1) = y | Q_i(t) = x, U_i \in \mathcal{U}(t)]$, for $x, y \in \{0, 1, \dots, C\}$. When node U_i is scheduled at slot t , and $Q_i(t) \geq 1$, one of its task is executed and it also informs the CC about the number of tasks left in the queue (observation). We assume that at most one task can be generated (or dropped) in a slot, so that $p_{xy}^{(u)} = 0$ for $y < x - 1$ and $y > x + 1$, with $u \in \{0, 1\}$ as shown in Fig. 4.

The belief of each i th node is represented by a $(C \times 1)$ vector $\omega_i = [\omega_{i,0}, \dots, \omega_{i,C-1}]$ whose k th entry $\omega_{i,k}$, for $k \in \{0, 1, \dots, C-1\}$, is given by (cf. (2)) $\omega_{i,k} = \Pr[Q_i(t) = k | \mathcal{H}(t)]$. The

immediate reward (3), given the initial belief vectors $\omega_1(t), \dots, \omega_M(t)$ and action \mathcal{U} , becomes

$$R(\omega_1(t), \dots, \omega_M(t), \mathcal{U}) = \sum_{i=1}^M \Pr [Q_i(t) > 0 | \mathcal{H}(t)] 1(U_i \in \mathcal{U}) = K - \sum_{i \in \mathcal{U}} \omega_{i,0}(t). \quad (34)$$

The performance of interest is the infinite-horizon throughput (17).

1) *Myopic Policy*: The MP (14), specialized to the immediate reward (34), becomes

$$\mathcal{U}^{MP}(t) = \operatorname{argmax}_{\mathcal{U}} R(\omega_1(t), \dots, \omega_M(t), \mathcal{U}) = \operatorname{argmin}_{\mathcal{U}} \sum_{i \in \mathcal{U}} \omega_{i,0}(t). \quad (35)$$

Note that, unlike Sec. III-A, when $C > 1$ the MP does not generally have a RR structure.

B. Upper Bound

Here we derive an upper bound to the throughput (17) by following the approach for general RMAB problems proposed in [6]. The upper bound relaxes the constraint that exactly K nodes must be scheduled in each slot. Specifically, it allows a variable number $K^\pi(t)$ of scheduled nodes in each t th slot under policy π , with the only constraint that its discounted average satisfies

$$E^\pi \left[\sum_{t=1}^{\infty} \beta^{t-1} K^\pi(t) \right] = \frac{K}{1-\beta}. \quad (36)$$

The advantage of this relaxed version of the scheduling problem is that it can be tackled by focusing on each single arm independently from the others [6], [12]. This is because, by the symmetry of the nodes, the constraint (36) can be equivalently handled by imposing that each node is active on average for a discounted time $E^\pi [\sum_{t=1}^{\infty} \beta^{t-1} 1(U_i \in \mathcal{U}^\pi(t))] = \frac{K}{M(1-\beta)}$. We can thus calculate the optimal solution of the relaxed problem by solving a single RSAB problem.

We now elaborate on such a RSAB by dropping the node index. Here, the immediate reward when the arm is in state ω (a vector since $C > 1$, see Sec. V-A), and action $u \in \{0, 1\}$ is chosen, is $R(\omega, u) = 1 - \omega_0$ if $u = 1$ and $R(\omega, u) = 0$ if $u = 0$, while the Markov evolution of the belief follows from Fig. 4 and similarly to Sec. I-A. The problem consists in optimizing the throughput under the constraint $E^\pi [\sum_{t=1}^{\infty} \beta^{t-1} 1(U_i \in \mathcal{U}^\pi(t))] = \sum_{t=1}^{\infty} \beta^{t-1} E^\pi [u^\pi(t)] = K/(M(1-\beta))$, as

introduced above. Under the assumption that the state ω belongs to a finite state space \mathcal{W} (to be discussed below), this optimization can be done by resorting to a linear programming (LP) formulation [12]. Specifically, let $z_{\omega}^{(u)}$ be the probability of being in state ω and selecting action $u \in \{0, 1\}$ under a given policy. The optimization at hand leads to the following LP

$$\text{maximize } \sum_{\omega, u} R(\omega, u) z_{\omega}^{(u)}, \quad (37a)$$

$$\text{subject to : } \sum_{\omega, u} z_{\omega}^{(u)} = 1, \quad (37b)$$

$$\sum_{\omega} z_{\omega}^{(1)} = \frac{K}{M(1 - \beta)}, \quad (37c)$$

$$z_{\omega}^{(0)} + z_{\omega}^{(1)} = \delta(\omega - \omega(1)) + \beta \sum_{\omega', u} z_{\omega'}^{(u)} p_{\omega\omega'}^{(u)}, \text{ for all } \omega \in \mathcal{W}, \quad (37d)$$

where (37c) is the constraint on the average time in which the node is scheduled, while (37d) guarantees that $z_{\omega}^{(u)}$ is the stationary distribution [12], in which $\delta(\omega - \omega(1)) = 1$ if $\omega = \omega(1)$ and $\delta(\omega - \omega(1)) = 0$ if $\omega \neq \omega(1)$. Note that, as discussed in Sec. II, the term $p_{\omega\omega'}^{(u)}$ is the probability that the next state is ω' given that action u is taken in state ω .

We are left to discuss the cardinality of the set \mathcal{W} . While the belief ω can generally assume any value in the C -dimensional probability simplex, the number of states actually assumed by ω during any *limited* horizon of time is finite due to the finiteness of the action space [10]. In our problem, since the time horizon is unlimited, this fact alone is not sufficient to conclude that the set \mathcal{W} is finite. However, after each t th slot in which the arm is activated, the belief at the $(t + 1)$ th slot can only take C values given that the queue state is learned by the CC. Therefore, the evolution of the belief is reset after each activation, and in practice, the time between two activations is finite since the node must be kept active for a discounted fraction of time $K / (M(1 - \beta))$. Hence, by constraining the maximum time interval between two activations to a sufficiently large value, the state space \mathcal{W} remains finite and the optimal performance is not affected. We used this approach for the numerical evaluation of the upper bound in Sec. V-C.

C. Numerical Results

We now present some numerical results to compare the performance of the MP with the upper bound of Sec. V-B. The performance is the throughput (17) normalized by its ideal value $K/(1 - \beta)$ that is obtained if the nodes always have a task to be completed when scheduled.

In Fig. 5 we show the normalized throughput versus the queue capacity C for different ratio M/K between the number M of nodes and the number K of nodes scheduled in each slot. We keep $K = 3$ fixed and vary M . We assume a uniform distribution for the initial number of tasks in the queues $Q_i(1)$ for all the nodes, so that $\omega_{i,k}(1) = 1/(C + 1)$ for all i, k . The probabilities that a new task is generated when the arm is kept passive are $p_{01}^{(0)} = 0.15$ and $p_{kk+1}^{(0)} = 0.1$, for $k \in \{1, C - 1\}$, while under activation they are $p_{01}^{(1)} = 0.05$ and $p_{kk+1}^{(1)} = 0$. The probability that a task expires when the arm is kept passive and activated are $p_{kk-1}^{(0)} = 0.05$ and $p_{kk-1}^{(1)} = 0.95$ respectively. The remaining transitions probabilities are $p_{CC}^{(0)} = 0.9$, $p_{CC}^{(1)} = 0.05$, while $\beta = 0.95$.

From Fig. 5 it can be seen that when C and/or M/K are small the MP's performance is close to the upper bound. In fact, for small M/K , most of the nodes are scheduled in each slot and the relaxed system in Sec. V-B approaches the original one, while for small C we get closer to the optimality of the MP for $C = 1$. For moderate to large values of M/K and/or C instead, the more flexibility in the relaxed system enables larger gains over the MP.

VI. CONCLUSIONS

This paper considers a centralized scheduling problem for independent, symmetric and time-sensitive tasks under resources constraints. The problem is to assign a finite number of resources to a larger number of *nodes* that may have tasks to be completed in their *task queue*. It is assumed that the *central controller* has no direct access to the queue of each node. Based on a Markovian modeling of the task generation and expiration processes, the scheduling problem is formulated as a partially observable Markov decision process (POMDP) and then cast into the framework of restless multi-armed bandit (RMAB) problems. Under the assumption that the task queues are of capacity one, a greedy, or *myopic policy* (MP), operating in the space of the a posteriori

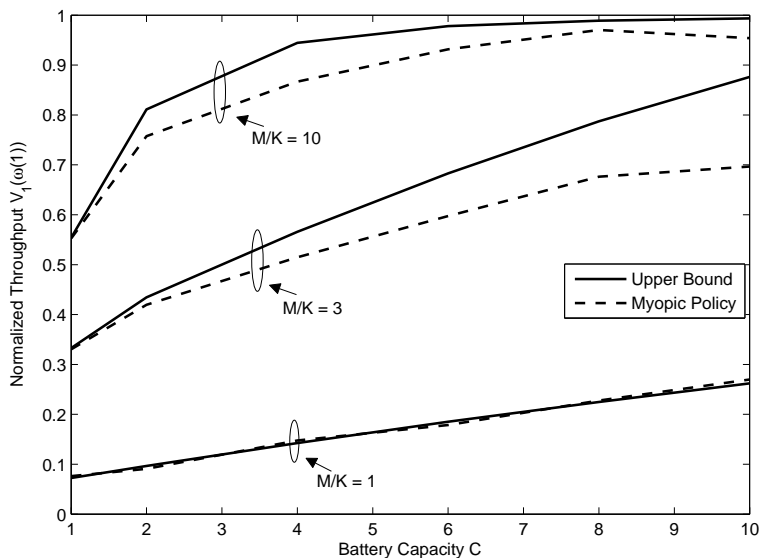


Figure 5. Normalized optimal throughput of the MP in (35) as compared to the upper bound versus the queue capacity C for different ratios $M/K \in \{1, 3, 10\}$ (system parameters are $K = 3$, $\beta = 0.95$, $\omega_{i,k}(1) = 1/(C + 1)$ for all i, k , $p_{01}^{(0)} = 0.15$, $p_{01}^{(1)} = 0.05$, $p_{CC}^{(0)} = 0.9$, $p_{CC}^{(1)} = 0.05$, $p_{kk-1}^{(0)} = 0.05$, $p_{kk-1}^{(1)} = 0.95$, $p_{kk+1}^{(0)} = 0.1$, $p_{kk+1}^{(1)} = 0$, for $k \in \{1, C - 1\}$).

probabilities (beliefs) of the number of tasks in the queues, is proved to be optimal, under appropriate assumptions, for both finite and infinite-horizon throughput criteria. The MP selects at each slot the nodes with the largest probability of having a task to be completed. It is shown that the MP is round-robin since it schedules the nodes periodically. We have also established that the RMAB problem at hand is indexable, derived the Whittle index in closed form and shown that the Whittle index policy is equivalent to the MP and thus it is optimal.

Systems in which the task queues have arbitrary capacities have been investigated as well by comparing the performance of the MP, which is generally suboptimal, with an upper bound based on a relaxation of the scheduling constraint.

Overall, this paper proposes a general framework for resource allocation that finds applications in several areas of current interest including communication networks and distributed computing.

APPENDIX A

PROOF OF THEOREM 3

The proof is divided into two steps. In the first step we derive the throughput of the RR policy in closed form, and then we show that inequality (16) holds.

As for the first step, the throughput for the RR policy (and thus of the MP) can be calculated as the sum of the contribution of each node separately (due to the round robin structure). To elaborate, let us focus on node U_i , with initial belief $\omega_i(1)$, and assume that $U_i \in \mathcal{G}_1$. Nodes in group \mathcal{G}_1 are scheduled at slots $t \in \{1+(j-1)m\}$, for $j \in \{1, 2, \dots\}$. Let $r_j(\omega_i(1)) = \mathbb{E}^{\text{RR}}[\omega_i(1+(j-1)m)|\omega_i(1)]$ be the average reward accrued by the CC from node U_i only, when scheduling it for the j th time at slot $t = 1+(j-1)m$ (see the RHS of (3)) (i.e., when operating the RR policy). At slot $t = 1$ we have $r_1(\omega_i(1)) = \omega_i(1)$. To calculate $r_2(\omega_i(1))$ we first derive the average value of the belief (see (7)) after the slot of activity in $t = 1$ as $\mathbb{E}^{\text{RR}}[\omega_i(2)|\omega_i(1)] = \tau_1^{(1)}(\omega_i(1))$, where $\tau_1^{(1)} = \omega\delta_1 + p_{01}^{(1)}$ with $\delta_u = (p_{11}^{(u)} - p_{01}^{(u)})$ (cf. (8)). We then account for the $(m-1)$ slots of passivity by exploiting (8), so that $r_2(\omega_i(1)) = \mathbb{E}^{\text{RR}}[\omega_i(1+m)|\omega_i(1)] = \phi^{(1)}(\omega_i(1))$, where we have set $\phi^{(1)}(\omega) = \tau_0^{(m-1)}(\tau_1^{(1)}(\omega)) = \omega\alpha_m + \psi_m$ with $\alpha_m = \delta_1\delta_0^{m-1}$ and $\psi_m = p_{01}^{(1)}\delta_0^{m-1} + p_{01}^{(0)}\frac{1-\delta_0^{m-1}}{1-\delta_0}$, and where $\tau_0^{(k)}(\omega) = \tau_0^{(1)}(\tau_0^{(k-1)}(\omega))$ indicates the belief of a node after k slots of passivity when the initial belief is ω (i.e., $\tau_0^{(k)}(\omega)$ is obtained recursively by applying $\tau_0^{(1)}(\omega)$ to itself k times). In general, we can obtain $r_j(\omega_i(1)) = \mathbb{E}^{\text{RR}}[\omega_i(1+(j-1)m)|\omega_i(1)]$, for $j \geq 2$, by iterating the procedure above by applying $\phi^{(1)}(\omega)$ to itself $(j-1)$ times. After a little algebra we get $\phi^{(j-1)}(\omega) = \phi^{(1)}(\phi^{(j-2)}(\omega)) = \omega\alpha_m^{j-1} + \psi_m\frac{1-\alpha_m^{j-1}}{1-\alpha_m}$, so that $r_j(\omega_i(1)) = \phi^{(j-1)}(\omega_i(1))$, where we set $\phi^{(0)}(\omega) = \omega$. The reasoning above can be applied when starting from any arbitrary slot t .

Finally, the total reward accrued by the CC from a node that is scheduled H times, when its belief at the first slot in which it is scheduled is ω , can be calculated by summing up the average reward $r_j(\cdot)$ during each slot in which the node is scheduled (see definition above), as

$$\theta^{(H)}(\omega) = \sum_{j=1}^H \beta^{(j-1)m} r_j(\omega) = \frac{\psi_m}{1-\alpha_m} \left(\frac{1-\beta^{mH}}{1-\beta^m} - \frac{1-(\beta^m\alpha_m)^H}{1-\beta^m\alpha_m} \right) + \frac{1-(\beta^m\alpha_m)^H}{1-\beta^m\alpha_m} \omega. \quad (38)$$

Note that, for a node $U_i \in \mathcal{G}_g$, for $g \geq 1$ and with belief equal to ω at $t = 1$, the first slot in which the node is scheduled is $t = g$, and thus its belief at time $t = g$ becomes $\tau_0^{(g-1)}(\omega)$ (i.e., after $(g-1)$ slots of passivity while other groups are scheduled). Therefore, for a node $U_i \in \mathcal{G}_g$, with initial belief ω , the total contribution to the throughput is given by $\beta^{g-1}\theta^{(H)}(\tau_0^{(g-1)}(\omega))$.

Let us now focus on the second step, i.e., proving the inequality (16). At $t = T$, it is easily seen to hold due to (3) and (12). We then need to show that (16) also holds at t . To do so, let us denote as \mathcal{L} and \mathcal{R} the RR policies whose throughputs are given by the LHS and RHS of (16) respectively. The differences between \mathcal{L} and \mathcal{R} are the positions of the nodes with belief x and y in the initial belief vectors. Therefore, some of the m groups created by the two policies might have different nodes (see the RR operations in Proposition 1). To simplify, we refer to the node with belief x (y) as node x (y). Let us assume that nodes x and y belong to groups $\mathcal{G}_{g'}$ and $\mathcal{G}_{g''}$ under policy \mathcal{R} , respectively, while they belong to groups $\mathcal{G}_{g''}$ and $\mathcal{G}_{g'}$ under policy \mathcal{L} , respectively, with $g'' \geq g'$, and $g', g'' \in \{1, \dots, m\}$. If $g'' = g'$, then the two policies coincide and (16) holds with equality. If $g'' = g' + 1$ (nodes are adjacent but do not belong to the same group), the only difference between policies \mathcal{L} and \mathcal{R} is the scheduling order of nodes x and y .

To verify that inequality (16) holds, we need to prove that scheduling node y in group $\mathcal{G}_{g'}$ and node x in group $\mathcal{G}_{g''}$ is no better than doing the opposite for any $x \geq y$. To elaborate, let $H_x^{\mathcal{R}}(t) = H_y^{\mathcal{L}}(t)$ and $H_y^{\mathcal{R}}(t) = H_x^{\mathcal{L}}(t)$ be the number of times that node x (or y) is scheduled under policy \mathcal{R} (or \mathcal{L}) and node y (or x) is scheduled under policy \mathcal{L} (or \mathcal{R}) in the horizon $\{t, t+1, \dots, T\}$, respectively. By recalling (38) and the discount factor β , the contribution generated by node x and y under policy \mathcal{R} is $\beta^{g'-1}\theta^{(H_x^{\mathcal{R}}(t))}(\tau_0^{(g'-1)}(x))$ and $\beta^{g''-1}\theta^{(H_y^{\mathcal{R}}(t))}(\tau_0^{(g''-1)}(y))$ respectively, and similarly under policy \mathcal{L} we have $\beta^{g''-1}\theta^{(H_x^{\mathcal{L}}(t))}(\tau_0^{(g''-1)}(x))$ and $\beta^{g'-1}\theta^{(H_y^{\mathcal{L}}(t))}(\tau_0^{(g'-1)}(y))$. Note that, in the argument of function $\theta^{(\cdot)}$, we have considered that the nodes in group $\mathcal{G}_{g'}$ are scheduled for the first time at slot $g' - 1$, and thus the belief must be updated through function $\tau_0^{(g'-1)}(\cdot)$, and similarly for nodes in $\mathcal{G}_{g''}$ the first slot is $g'' - 1$. Moreover, the discount factor is $\beta^{g'-1}$ is common to all the nodes in group $\mathcal{G}_{g'}$, and so is $\beta^{g''-1}$ for group $\mathcal{G}_{g''}$.

By recalling that all the nodes, except x and y , are scheduled at the same slot under the two policies \mathcal{R} and \mathcal{L} (thus giving the same contribution to the throughput), the inequality (16) can thus be reduced to $\beta^{g'-1}\theta^{(H_x^{\mathcal{R}}(t))}(\tau_0^{(g'-1)}(x)) + \beta^{g''-1}\theta^{(H_y^{\mathcal{R}}(t))}(\tau_0^{(g''-1)}(y)) - \beta^{g''-1}\theta^{(H_x^{\mathcal{L}}(t))}(\tau_0^{(g''-1)}(x)) - \beta^{g'-1}\theta^{(H_y^{\mathcal{L}}(t))}(\tau_0^{(g'-1)}(y)) \geq 0$, which must hold for all admissible $H_x^{\mathcal{R}}(t) = H_y^{\mathcal{L}}(t)$ and $H_y^{\mathcal{R}}(t) = H_x^{\mathcal{L}}(t)$.

$= H_x^{\mathcal{L}}(t)$ and all $g'' \geq g'$, with $g', g'' \in \{1, \dots, m\}$. There are two cases: **1)** $H_x^{\mathcal{R}}(t) = H_y^{\mathcal{L}}(t) = H_y^{\mathcal{R}}(t) = H_x^{\mathcal{L}}(t) = H \geq 1$, that is, nodes x and y are scheduled the same number of times within the horizon of interest under the two policies \mathcal{R} and \mathcal{L} ; **2)** $H_x^{\mathcal{R}}(t) = H_y^{\mathcal{L}}(t) = H$, and $H_y^{\mathcal{R}}(t) = H_x^{\mathcal{L}}(t) = H - 1$, for $H \geq 1$, namely, node x (or y) is scheduled one time more than node y (or x) under policy \mathcal{R} (or \mathcal{L}). By exploiting the RHS of (38), after a little algebra, one can verify that the inequality above holds in both cases, which concludes the proof of Theorem 3.

APPENDIX B

PROOF OF LEMMA4

Proof of case a). From (23)-(24), and recalling that $\tau_0^{(1)}(0) = p_{01}$ from (29), the leftmost inequality in (27a.1) follows immediately as it becomes $V_m(0|1) = \beta V_m^*(p_{01}) \leq m + \beta V_m^*(p_{01}) = V_m(0|0)$. For the rightmost inequality in (27a.1), we have $V_m(1|1) = 1 + \beta V_m^*(0)$, while from (21) and the fact that $V_m(0|1) \leq V_m(0|0)$ we have $V_m^*(0) = \max\{V_m(0|0), V_m(0|1)\} = V_m(0|0)$. Therefore, we have $V_m(1|1) = 1 + \beta V_m^*(0) \geq 1 + \beta V_m(0|0) \geq V_m(0|0)$, which holds as $1 + \beta V_m(0|0) \geq V_m(0|0)$ implies $V_m(0|0) \leq \frac{1}{1-\beta}$. The latter bound always holds, since for $m < 1$ the infinite horizon throughput is upper bounded as $V_m^*(\omega) \leq \sum_{t=0}^{\infty} \beta = \frac{1}{1-\beta}$ given that we can get at most a reward of $R_m(\omega, u) \leq 1$ in each slot. Hence, inequalities (27a.1) are proved. Inequality (27a.2) can be proved by contradiction. Specifically, let us assume that: *hp.I*) $V_m(1|0) \geq V_m(1|1)$. From (21) we would have $V_m^*(1) = \max\{V_m(1|0), V_m(1|1)\} = V_m(1|0)$, i.e., the passive action would be optimal when $\omega = 1$. Moreover, from (23) we would have $V_m(1|0) = m + \beta V_m^*(1) = m + \beta V_m(1|0)$, which can be solved with respect to $V_m(1|0)$ to get $V_m(1|0) = \frac{m}{1-\beta} = V_m^*(1)$. Therefore, if hypothesis *hp.I*) holds, we also have that $V_m(1|1) = 1 + \beta V_m^*(0) \leq V_m(1|0) = V_m^*(1) = \frac{m}{1-\beta}$. However, the value function $V_m^*(\omega)$ is bounded $\frac{m}{1-\beta} \leq V_m^*(\omega) \leq \frac{1}{1-\beta}$, where the lower bound is obtained considering a policy that always chooses the passive action for any belief ω . The boundedness of the value function, thus implies that if *hp.I*) holds then $1 + \beta \frac{m}{1-\beta} \leq 1 + \beta V_m(0) = V_m(1|1) \leq V_m(1|0) = \frac{m}{1-\beta}$, which yields $1 + \beta \frac{m}{1-\beta} \leq \frac{m}{1-\beta}$ and thus $(1 - \beta)(1 - m) \leq 0$. But this is clearly impossible as $m, \beta < 1$.

Consequently, we have proved that $V_m(1|1) \geq V_m(1|0)$.

Proof of case b) Inequality $V_m(0|0) \leq V_m(0|1)$ follows immediately since $m + \beta V_m^*(p_{01}) \leq \beta V_m^*(p_{01})$ holds for $m < 0$. The second inequality $V_m(0|1) \leq V_m(1|1)$ becomes $V_m(0|1) \leq 1 + \beta V_m^*(0)1 + \beta V_m(0|1)$, which leads to $V_m(0|1) \leq \frac{1}{1-\beta}$, which always holds as discussed above. Inequality $V_m(1|0) \leq V_m(1|1)$ holds since an active action is always optimal when $m < 0$.

Proof of case c) The inequality holds since a passive action is always optimal for any $m \geq 1$.

APPENDIX C

PROOF OF THEOREM 8

Following the discussion in Sec. IV-A2, to prove indexability it is sufficient to show that the threshold $\omega^*(m)$ is monotonically increasing with the subsidy m , for $0 \leq m < 1$. In fact, from Proposition 5 the passive set (25) for $m < 0$ is $\mathcal{P}(m) = \emptyset$, while for $m \geq 1$, we have $\mathcal{P}(m) = [0, 1]$. We then only need to prove the monotonicity of $\omega^*(m)$ for $0 \leq m < 1$, which has been shown to hold in [9, Lemma 9] if

$$\left. \frac{dV_m(\omega|1)}{dm} \right|_{\omega=\omega^*(m)} < \left. \frac{dV_m(\omega|0)}{dm} \right|_{\omega=\omega^*(m)}. \quad (39)$$

To check if (39) holds, we differentiate (23)-(24) at the optimal threshold $\omega = \omega^*(m)$ as

$$V_m(\omega^*(m)|1) = \omega^*(m) + \beta\omega^*(m)V_m^*(0) + \beta(1 - \omega^*(m))V_m^*(p_{01}), \text{ and} \quad (40)$$

$$V_m(\omega^*(m)|0) = m + \beta \left[\tau_0^{(1)}(\omega^*(m)) (1 + \beta V_m^*(0)) + \beta(1 - \tau_0^{(1)}(\omega^*(m)))V_m^*(p_{01}) \right], \quad (41)$$

where (41) follows from (24) and from the fact that $\tau_0^{(1)}(\omega) \geq \omega$, for any ω (see (29)), and hence $V_m^*(\tau_0^{(1)}(\omega^*(m))) = V_m(\tau_0^{(1)}(\omega^*(m))|1)$, since arm activation is optimal for any $\omega > \omega^*(m)$.

By letting $D_m(\omega) = \frac{dV_m^*(\omega)}{dm}$, then from (40) we have $\left. \frac{dV_m(\omega|1)}{dm} \right|_{\omega=\omega^*(m)} = \beta\omega^*(m)D_m(0) + \beta(1 - \omega^*(m))D_m(p_{01})$, while from (41) we get $\left. \frac{dV_m(\omega|0)}{dm} \right|_{\omega=\omega^*(m)} = 1 + \beta^2\tau_0^{(1)}(\omega^*)D_m(0) + \beta^2(1 - \tau_0^{(1)}(\omega^*))D_m(p_{01})$. Finally, after some algebraic manipulations, and recalling that $D_m(0) = \frac{dV_m^*(0)}{dm} = \frac{d(m + \beta V_m^*(p_{01}))}{dm} = 1 + \text{recursively } \beta D_m(p_{01})$, we can rewrite (39) as

$D_m(p_{01})\beta(1-\beta)[1-\omega(1-\beta(1-p_{01}))] + \beta[\omega(1-\beta(1-p_{01})) - \beta p_{01}] < 1$. To show that the last inequality holds when $0 \leq m < 1$, we first upper bound the derivative of the value function as $D_m(\omega) \leq \frac{1}{1-\beta}$, since $\frac{d}{dm}R_m(\omega) \leq 1$. Finally, using this upper bound $D_m(p_{01}) \leq \frac{1}{1-\beta}$ after a little algebra (39) reduces to $\beta(1-\beta p_{01}) < 1$, which clearly holds for any $\beta \in [0, 1)$ as $0 \leq p_{01} \leq 1$. This concludes the proof of Theorem 8.

REFERENCES

- [1] D. Bertsekas, R. G. Gallager, *Data Networks*. Englewood Cliffs, NJ: Prentice Hall, 1992.
- [2] A. Benoit, L. Marchal, J.-F. Pineau, Y. Robert, F. Vivien, "Scheduling concurrent bag-of-tasks applications on heterogeneous platforms," *IEEE Trans. Computers*, vol. 59, no. 2, pp. 202-217, Feb. 2010.
- [3] Y. Bai, C. Xu and Z. Li, "Task-aware based co-scheduling for virtual machine system," in *Proc. ACM Symp. On Applied Comp.*, Sierre, Switzerland, pp. 181-188, Mar. 2010.
- [4] G. E. Monahan, "A survey of partially observable Markov decision processes: Theory, models, and algorithms," *Manag. Sci.*, vol. 28, no. 1, pp. 1-16, 1982.
- [5] J. Gittins, K. Glazerbrook, R. Weber, *Multi-armed Bandit Allocation Indices*. West Sussex, UK: Wiley, 2011.
- [6] P. Whittle, "Restless bandits: Activity allocation in a changing world," *J. Appl. Probab.*, vol. 25, pp. 287-298, 1988.
- [7] S. H. A. Ahmad, M. Liu, T. Javidi, Q. Zhao and B. Krishnamachari, "Optimality of myopic sensing in multi-channel opportunistic access," *IEEE Trans. Inf. Theory*, vol. 55, No. 9, pp. 4040-4050, Sept. 2009.
- [8] S. H. A. Ahmad, M. Liu, "Multi-channel opportunistic access: A case of restless bandits with multiple plays," in *Proc. 47th Ann. Allerton Conf. Commun., Contr., Comput.*, Monticello, IL, pp. 1361-1368, Sept. 2009.
- [9] K. Liu and Q. Zhao, "Indexability of restless bandit problems and optimality of Whittle index for dynamic multichannel access," *IEEE Trans. Inf. Theory*, vol. 56, no. 11, pp. 5547-5567, Nov. 2010.
- [10] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra, "Planning and acting in partially observable stochastic domains," *Artif. Intell.*, vol. 101, pp. 99-134, May 1998.
- [11] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Hoboken, NJ: Wiley, 2005.
- [12] D. Bertsimas and J. E. Niño-Mora, "Restless bandits, linear programming relaxations, and a primal-dual heuristic," *Oper. Res.*, vol. 48, no. 1, pp. 80-90, Jan. 2000.