# Resource-on-Demand Schemes in 802.11 WLANs With Non-Zero Start-Up Times

Jorge Ortín, Carlos Donato, Pablo Serrano, Senior Member, IEEE, and Albert Banchs, Senior Member, IEEE

Abstract-Increasing the density of access points is one of the most effective mechanisms to cope with the growing traffic demand in wireless networks. To prevent energy wastage at low loads, a resource-on-demand (RoD) scheme is required to opportunistically (de)activate access points as network traffic varies. While previous publications have analytically modeled these schemes in the past, they have assumed that resources are immediately available when activated, an assumption that leads to inaccurate results and might result in inappropriate configurations of the RoD scheme. In this paper, we analyze a general RoD scenario with N access points and non-zero start-up times. We first present an exact analytical model that accurately predicts performance but has a high computational complexity, and then derive a simplified analysis that sacrifices some accuracy in exchange for a much lower computational cost. To illustrate the practicality of this model, we present the design of a simple configuration algorithm for RoD. Simulation results confirm the validity of the analyses, and the effectiveness of the configuration algorithm.

*Index Terms*—WLAN, 802.11, resource on demand, energy consumption, infrastructure on demand.

#### I. INTRODUCTION

**T**O COPE with the growing demand of wireless traffic, one of the most effective approaches is to increase the density of access points (AP) in the network. The side effect of this strategy, though, is the increase of the power consumed, which can result in energy wastage if all the infrastructure is kept powered on when the load is low [1], [2]. Techniques to "green" the operation of the network include the design of more energy efficient hardware [3], the optimization of the radio transmission chain [4], or the implementation of

Manuscript received January 31, 2016; revised May 14, 2016 and September 18, 2016; accepted October 22, 2016. Date of publication November 1, 2016; date of current version December 29, 2016. The work of J. Ortín was supported in part by the Gobierno de Aragón (research group T98) and the European Social Fund, in part by the EU H2020 Wi-5 Project under Grant 644262, and in part by the Centro Universitario de la Defensa under Project CUD2013-05. The work of C. Donato, P. Serrano, and A. Banchs was supported by the European Commission in the framework of the H2020-ICT-2014-2 Project Flex5Gware under Grant 671563 and in part by the Madrid Regional Government through the TIGRE5-CM program under Grant S2013/ICE-2919. (*Corresponding author: Jorge Ortín.*)

J. Ortín is with the Centro Universitario de la Defensa, Zaragoza 50090 Spain, and also with the Aragón Institute of Engineering Research, Universidad de Zaragoza, Zaragoza 5018, Spain (e-mail: jortin@unizar.es).

C. Donato and A. Banchs are with the Department of Telematics Engineering, Universidad Carlos III de Madrid, Madrid 28911, Spain, and also with the Institute IMDEA Networks, Madrid, Spain (e-mail: carlos.donato@imdea.org; banchs@it.uc3m.es).

P. Serrano is with the Department of Telematics Engineering, Universidad Carlos III de Madrid, Madrid, Spain (e-mail: pablo@it.uc3m.es).

Color versions of one or more of the figures in this paper are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/JSAC.2016.2624158

resource-on-demand (RoD) strategies that dynamically adapt to the network load, activating resources as it grows and deactivating them when it shrinks [5].

RoD schemes are relatively easy to deploy, as they do not require the introduction of major changes in the network, and have been proposed to decrease the energy consumption of base stations in "traditional" mobile networks (GSM, UMTS) [6], [7], as these devices account for up to 60% of the total energy consumed [3]. Following [8], RoD policies can be divided into static and dynamic, depending on whether the switching on/off of the devices is scheduled or it follows real-time traffic patterns. In general, dynamic approaches are more efficient than static solutions, although they require higher switching on/off rates. In [5], a number of dynamic approaches are classified according to the wireless technology, performance metric, reaction time of the algorithm and control scheme (centralized or distributed).

Regarding WLAN networks, it has been shown that RoD techniques can potentially provide substantial gains in energy savings when the number of considered APs increases (gains of approx. 37%, or 26 kW, can be achieved for a university campus [9] or even higher [10]). Several publications have shown the feasibility of RoD policies in campus networks [11], [12]. The first analytical model for these techniques [13] focuses on the case of "clusters" of APs covering the same area, and studies the impact of the strategy used to (de)activate APs on parameters such as the energy savings and the switchoff rate of the devices. In [14], the work is extended to account for the case when APs do not completely overlap their coverage areas. Following this interest in RoD schemes for WLANs, new publications analyse the performance when some assumptions are relaxed,<sup>1</sup> e.g., Garroppo *et al.* [15] analyse the impact of using an accurate energy consumption model on performance.

In this paper, we analyse the impact of start-up times on the performance of a RoD scheme. By "start-up time" we mean the time it takes between the AP is activated until the WLAN is announced. According to the seminal work of [16], typical start-up times of an AP range between 12 and 35 seconds, yet they have not been considered in previous analytical models. However, in our previous work [17], we already showed that for the simple case of 2 APs, start-up times have a notable impact on performance, both qualitatively and quantitatively, as compared to the ideal case of immediate boot times. In that

0733-8716 © 2016 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications\_standards/publications/rights/index.html for more information.

<sup>&</sup>lt;sup>1</sup>In fact, in both [5] and [8] it is noted that current RoD policies are made using over-simplified models.

work, we also confirmed that the time to power on an AP is is on the same order of magnitude –approx. 45 s– and practically constant.

We analyse now the general case of a RoD scenario consisting of N overlapping APs with non-zero start-up times. We present an analytic *exact* model that accurately predicts performance in terms of energy consumption and service time, but with a high computational complexity. Because of this complexity, we then present a *simplified* model that sacrifices some numerical accuracy in exchange for more affordable computational times. Finally, we present one possible use of this simplified model, namely, the design of a simple configuration algorithm for RoD, based on the minimization of the average service time. As the results show, the simplified model supports the design of optimization policies that tradeoff performance for significant gains in energy efficiency.

# II. SYSTEM MODEL

# A. Scenario

We consider a *cluster model* like the one analysed in [13], consisting of N identical APs serving the same area but using non-overlapping channels. Although in typical high-density deployments the APs may not be located exactly at the same position, the high level of overlapping allows making this assumption, which simplifies the theoretical analysis. Indeed, as will be seen in Section V, this assumption does not impact the validity of the RoD strategy.

The need for a dense deployment such as the one addressed in this paper is motivated by the current trends in the increase of traffic demand. This trend has been forecasted by a number of sources. According to [18], the number of devices and connections per user is steadily growing, which increases user densities; in addition, the throughput required per user is also increasing, as new services such as HD video streaming are becoming ubiquitous. Along the same lines, the forecasts for future 5G networks,<sup>2</sup> predict data rates 100 times higher than today's. Even now, a recent research estimates that typical densities in the deployment of APs may exceed 4000 APs per square kilometre [10].

The scenario considered could be mapped to a very-dense 802.11a setup, where there are many available channels in the 5 GHz band (the specific number depending on the country). One of the APs is always on, in order to maintain the WLAN coverage, while the other APs are opportunistically powered on (off) as users arrive (leave) the system. Powering an AP takes a deterministic time  $T_{on}$  and, during this time, the AP is not available, so arriving requests are served by any of the other APs. We neglect the time required to power off an AP.

Each AP consumes  $P_{AP}$  units of power when active (i.e., during start-up and when powered on) and zero otherwise. Although commodity hardware can support an intermediate state (i.e., switching on/off the wireless card), this does not bring as much savings as powering on/off the complete device [19]. A "user" is a new connection generated by a wireless client. Following [12] and [20], these are generated according to a Poisson process at rate  $\lambda$  and are always

<sup>2</sup>http://5g-ppp.eu/.

served by the less loaded AP. Also following [12], we further assume that users' demands are exponentially distributed (i.e., each user downloads an amount of data that is exponentially distributed) and that the AP bandwidth is evenly shared among all the users.<sup>3</sup>

Based on the above assumptions, service times (i.e., the time elapsed since a user arrives to the WLAN until it has fully downloaded its demanded data) would be exponentially distributed (with mean  $1/\mu$ ) if every user got all the bandwidth of an AP, and the service rate (i.e., the inverse of the average service time) is  $\mu$  when there is only one serving AP,  $2\mu$  when there are two APs serving, etc. (i.e., we neglect the impact of channel sharing via contention). The total load is given by  $\rho = \lambda/N\mu$ .

We also assume a load-balancing algorithm such that users (re)associate while they are being served, and that this (re)association time is negligible –note that this can be achieved with the recent 802.11v and 802.11r amendments [23], which support triggering re-associations and performing fast transitions, respectively, with minor disruption of the service.

#### B. Resource on Demand Policy

In order to power on/off the APs we assume there is a "target" number M > 1 of users per AP, i.e., the system will opportunistically power on/off APs in order to keep that "target" number across resources (except for one AP that will be always on, to guarantee coverage). Based on this, we assume a threshold-based policy with hysteresis, namely:

- An AP will be powered on when the number of user per AP is ρ<sub>h</sub> higher than this target value.
- An AP will be powered off when the number of user per AP is ρ<sub>l</sub> lower than this target value.

In this way, with K APs powered on, the K + 1-th AP will be powered on when the number of users reaches

Threshold to power on another AP  $(N_K)$ :  $\lceil (1 + \rho_h) K M \rceil$ 

while with K APs powered on, one AP will be powered off when the number of user reaches

Threshold to power off an AP  $(n_K)$ :  $\lfloor (1 - \rho_l) K M \rfloor$ 

We next impose some conditions on these thresholds to support an efficient operation. On the one hand, with K APs powered on we impose that there are at least K associated users, so all APs are serving traffic. This results in that the threshold to power off an AP with K users has to be at least K, i.e.,

$$n_K = \lfloor (1 - \rho_l) K M \rfloor \ge K,\tag{1}$$

which results on the following condition for  $\rho_l$ 

$$\rho_l < 1 - \frac{1}{M}.\tag{2}$$

<sup>3</sup>The assumption on the Poissonian nature of user arrivals is aligned with the characterization driven by measurements provided by [21] and [22]. Furthermore, [12] shows that, while the distribution of the duration of user connections is not a memory-less process, it can still be approximated by an exponential distribution with reasonable accuracy. In the numerical evaluation, we will rely on a more accurate traffic model in order to assess the impact of the simplifying assumptions upon which our analysis relies. On the other hand, to prevent (or, at least, reduce) "flipflop" effects in the WLAN (i.e., to power on an AP and, once active, immediately power it off), we assume that the  $\rho_h$  and  $\rho_l$ thresholds are set such that the number of users to power on an AP when K of them are already serving traffic is larger than the number of users required to power off an AP when K + 1 are serving traffic, i.e.,

$$\lceil (1+\rho_h)KM \rceil > \lfloor (1-\rho_l)(K+1)M \rfloor.$$
(3)

Based on the above condition, and neglecting the rounding operations, to prevent the flip-flop effects the following condition between  $\rho_h$  and  $\rho_l$  should hold

$$\rho_h > \frac{1 - \rho_l}{K} - \rho_l,\tag{4}$$

where the rhs of (4) is maximum for K = 1, i.e., the case of one AP, and therefore the  $\rho_h$  threshold should be set to at least<sup>4</sup>

$$\rho_h > 1 - 2\rho_l. \tag{5}$$

In addition to the above, for analytical tractability we introduce the following restriction on the RoD policy: at any point in time there will be at most one AP being powered on. More specifically, while one AP is powering on there will be no decisions taken w.r.t. powering on or off other resources, and only once the AP is available the system will decide on the amount of resources needed.

# III. EXACT ANALYSIS A. Model Overview

We model the system with the semi-Markov process illustrated in Fig. 1. The label in each arrow corresponds to the number of users (or range of users) in the system that triggers the transition between the corresponding stages. There are four types of *stages*, depending on the transitions that could happen between them:

- Stage 1, which is the initial situation with only one AP active. The only possible transition is to stage 1\* (another AP is powered on), that is triggered when the number of users reaches *N*<sub>1</sub>.
- Stage N, when all APs are active and serving traffic. The only possible transition is to stage N 1 (one AP is powered off), what happens when the number of users is  $n_N$ .
- Stages K (with 1 < K < N), where there are K active APs. In this case there are two possible transitions: one to stage K \*, triggered when the number of users reaches  $N_K$  and another AP is powered on (label  $N_K$  in Fig. 1); and other to stage K 1, triggered when the number of users in the system is  $n_K$  and one AP is powered off (label  $n_K$  in Fig. 1).
- Stages K\* (with  $1 \le K < N$ ), where in addition to the *K* active APs there is another AP booting up. For this type of stage there is a larger number of possible transitions, which are determined by the number of users in the system after  $T_{on}$ :

<sup>4</sup>Note that for simplicity our policy is set on fixed values of  $\rho_h$  and  $\rho_l$ , i.e., they do not change with the number of active APs.



Fig. 1. Semi-Markov process for an IoD scheme with N = 4 APs.

- If there are  $N_{K+1}$  or more users, the system will move to stage K + 1\*, as the number of users is already above the threshold to switch on an additional AP. These transitions are marked with the label  $\geq N_{K+1}$  in Fig. 1.
- If there are between  $n_{K+1}+1$  and  $N_{K+1}-1$  users, the system will move to stage K + 1. These transitions are marked with the label  $(n_{K+1}, N_{K+1})$  in Fig. 1.
- If there are  $n_{K+1}$  users or less, the next stage will depend on whether the number of users is also less than or equal to  $n_2$  (and therefore the next stage will be '1'), between  $n_2 + 1$  and  $n_3$  (the next stage will be '2'), and so on. These transitions are marked with the labels  $\leq n_2$ ,  $(n_2, n_3]$ , ... in Fig. 1.

With the above, we have introduced the different stages of the semi-Markov process. We next analyse each type of stage, their holding times, and the transition probabilities between them.

#### B. Modelling the Stages of the Semi-Markov Model

1) Stage 1 ( $S_1$ ): For the initial situation with only one AP active, following our assumptions the system can be modelled with the continuous-time Markov chain (CTMC) illustrated in Fig. 2, where each state represents the number of users being served and therefore reaching the absorbing state  $N_1$  corresponds to the case when another AP will be powered on (and stage 1 will be left).

$$\underbrace{0}_{\mu} \underbrace{1}_{\mu} \underbrace{2}_{\mu} \underbrace{\lambda}_{\mu} \underbrace{\lambda}_{\mu} \underbrace{N_{1}-1}_{\mu} \underbrace{N_{1}}_{\mu} \underbrace{N_{1}-1}_{\mu} \underbrace{N_{1}}_{\mu} \underbrace{N_{1}-1}_{\mu} \underbrace{N_{1}}_{\mu} \underbrace{N_{1}-1}_{\mu} \underbrace{N_{1}}_{\mu} \underbrace{N_{1}-1}_{\mu} \underbrace{N_{1}-1}_$$

Fig. 2. CTMC for stage 1: one active AP and no AP powering on.

The average time in this stage  $H_1$  corresponds to the *time* until absorption of the Markov chain, i.e., the time since the system arrived to the chain until it reaches the absorbing state  $N_1$ . If we define  $L_i(t)$  as the expected total time that a CTMC spends in state *i* during the interval [0, t),  $H_1$  can be expressed as the sum of the terms  $L_i(t)$  for all the nonabsorbing states of the CTMC when  $t \to \infty$  [24]

$$H_1 = \sum_{i=0}^{N_1 - 1} L_i^{(1)}(\infty).$$
(6)

The values of  $L_i^{(1)}(\infty)$  (the superscript (1) indicates that we are referring to the CTMC modelling stage 1) can be obtained solving the following system of equations:

$$\mathbf{L}^{(1)}(\infty)\mathbf{Q}^{(1)} = -\boldsymbol{\pi}^{(1)}(0), \qquad (7)$$

where

$$\mathbf{L}^{(1)}(\infty) = \left[ L_0^{(1)}(\infty), L_1^{(1)}(\infty), \dots L_{N_1 - 1}^{(1)}(\infty) \right], \quad (8)$$

$$\boldsymbol{\pi}^{(1)}(0) = \left[\pi_0^{(1)}(0), \pi_1^{(1)}(0), \dots, \pi_{N_1-1}^{(1)}(0)\right], \tag{9}$$

with  $\pi_i^{(1)}(0)$  the initial probability of state *i*, and  $\mathbf{Q}^{(1)}$  a  $N_1 \times N_1$  matrix with the following non-zero elements:

$$q_{ij} = \begin{cases} -\lambda & i = 1, j = 1 \\ -\lambda - \mu & i = 2, \dots, N_1, j = i \\ \lambda & i = 1, \dots, N_1 - 1, j = i + 1 \\ \mu & i = 2, \dots, N_1, j = i - 1 \end{cases}$$
(10)

The computation of  $\pi^{(1)}(0)$  is not straightforward, as it depends on the stage the system was before arriving to stage 1, which could be stage 2 or any other stage  $K^*$ , with  $K \ge 1$ . We detail how to compute  $\pi^{(1)}(0)$  for this and the other cases in the next section, after we present the modelling of the other stages of the semi-Markov process.

Finally, let  $P(S_T | S_{T'})$  denote the transition probability from stage T' to stage T, with T and T' referring indistinctly to any stage K, including stages 0 and N, or K\*. For the case of stage 1, we have that

$$P(S_{1*} \mid S_1) = 1. \tag{11}$$

2) Stages K ( $S_K$ ), 1 < K < N: For these stages, the resulting CTMC is illustrated in Fig. 3. In this case, while the arrival rate is also  $\lambda$ , the service rate accounts for the total number of powered-on APs, which is constant and equal to  $K \cdot \mu$  for all states.<sup>5</sup> As described above, there are two absorbing states: one corresponding to the powering on of another AP (when the system reaches  $N_K$  users), and another corresponding to the de-activation of one AP (when the number of users is  $n_K$ ).

<sup>5</sup>Note that we have imposed  $n_K > K$  with (1).

$$(n_{K}) \underbrace{n_{K}+1}_{K\mu} \underbrace{n_{K}+1}_{K\mu} \underbrace{\dots}_{K\mu} \underbrace{N_{K}-1}_{K\mu} \underbrace{N_{K}}_{K} \underbrace{N_{K}}_{K}$$

Fig. 3. CTMC for Stage K: K active APs and no AP powering on.

Similarly to the previous case, the average time in a stage *K* can be computed as

$$H_K = \sum_{i=n_K+1}^{N_K-1} L_i^{(K)}(\infty).$$
 (12)

In order to compute  $L_i^{(K)}(\infty)$  we use the same expression as in the previous case

$$\mathbf{L}^{(K)}(\infty)\mathbf{Q}^{(K)} = -\boldsymbol{\pi}^{(K)}(0), \qquad (13)$$

where

$$\mathbf{L}^{(K)}(\infty) = \left[ L_{n_{K}+1}^{(K)}(\infty), L_{n_{K}+2}^{(K)}(\infty), \dots L_{N_{K}-1}^{(K)}(\infty) \right], \quad (14)$$
$$\boldsymbol{\pi}^{(K)}(0) = \left[ \pi_{n_{K}+1}^{(K)}(0), \pi_{n_{K}+2}^{(K)}(0), \dots \pi_{N_{K}-1}^{(K)}(0) \right], \quad (15)$$

and  $\mathbf{Q}^{(K)}$  is a  $(N_K - n_K - 1) \times (N_K - n_K - 1)$  matrix, whose non-zero elements are

$$q_{ij} = \begin{cases} -\lambda - K\mu & i = 1, \dots, N_K - n_K - 1, j = i \\ \lambda & i = 1, \dots, N_K - n_K - 2, j = i + 1 \\ K\mu & i = 2, \dots, N_K - n_K - 1, j = i - 1 \end{cases}$$
(16)

Again, the computation of  $\pi^{(K)}(0)$  requires the knowledge of the stage of the system before entering stage K, which we will address in the next section.

To finalise the analysis of this stage, we have to compute the two transition probabilities from this stage to stage K \* (when the chain ends in the absorbing state  $N_K$ ) and to stage K - 1 (when the chain falls into the absorbing state  $n_K$ ), denoted as  $P(N_K)$  and  $P(n_K)$  respectively,

$$P(S_{K*} \mid S_K) = P(N_K),$$
(17)

$$P(S_{K-1} | S_K) = P(n_K).$$
(18)

These can be computed as [25]

$$[P(n_K) \ P(N_K)] = \pi^{(K)}(0)\mathbf{B}^{(K)}, \tag{19}$$

with  $\mathbf{B}^{(K)}$  a  $(N_K - n_K - 1) \times 2$  matrix whose element  $b_{ij}$  is the probability of ending in the absorbing state j, given that the chain starts in the transient state i. This matrix can be computed as

$$\mathbf{B}^{(K)} = \left[\mathbf{I} - \mathbf{T}^{(K)}\right]^{-1} \mathbf{R}^{(K)}, \qquad (20)$$

where **I** is the identity matrix,  $\mathbf{T}^{(K)}$  is a  $(N_K - n_K - 1) \times (N_K - n_K - 1)$  matrix with the transition probabilities between nonabsorbing states, and  $\mathbf{R}^{(K)}$  is a  $(N_K - n_K - 1) \times 2$  matrix denoting the transition probabilities from non-absorbing to absorbing states. Both matrices are obtained from the associated discrete-time Markov chain (DTMC) of the CTMC, and

$$(n_N)$$
  $(n_N+1)$   $(n_N+2)$   $(n_N+2$ 

Fig. 4. Markov chain when all APs are active.

their non-zero elements are

$$t_{ij} = \begin{cases} \frac{\lambda}{\lambda + K\mu} & i = 1, \dots, N_K - n_K - 2, \, j = i+1\\ \frac{K\mu}{\lambda + K\mu} & i = 2, \dots, N_K - n_K - 1, \, j = i-1 \end{cases}$$
(21)

$$r_{ij} = \begin{cases} \frac{K\mu}{\lambda + K\mu} & i = 1, j = 1\\ \frac{\lambda}{\lambda + K\mu} & i = N_K - n_K - 1, j = 2 \end{cases}$$
(22)

3) Stage N ( $S_N$ ): When all APs are active and serving traffic the resulting CTMC is the one depicted in Fig. 4, where the arrival rate is  $\lambda$  and the service rate is  $N \cdot \mu$ . As in the case of stage 1, there is only one absorbing state, the one corresponding to the switching off of one AP when there are  $n_N$  users in the system and all the APs are on, but now the chain has an infinite number of states.

To compute the holding time in this stage, we assume that the system is stable (i.e.,  $\lambda < N\mu$ ), so there is a state  $n_D$  with  $n_D > n_N$  such that

$$\sum_{n=D+1}^{\infty} L_i^{(N)}(\infty) \approx 0, \qquad (23)$$

and therefore the holding time is

$$H_N \approx \sum_{i=n_N+1}^{n_D} L_i^{(N)}(\infty), \qquad (24)$$

where  $\mathbf{L}^{(N)}(\infty)$  is obtained from

$$\mathbf{L}^{(N)}(\infty)\mathbf{Q}^{(N)} = -\boldsymbol{\pi}^{(N)}(0), \qquad (25)$$

with

$$\mathbf{L}^{(N)}(\infty) = \left[ L_{n_N+1}^{(N)}(\infty), L_{n_N+2}^{(N)}(\infty), \dots L_{N_D}^{(N)}(\infty) \right], \quad (26)$$

$$\boldsymbol{\pi}^{(N)}(0) = \left[ \pi_{n_N+1}^{(N)}(0), \pi_{n_N+2}^{(N)}(0), \dots, \pi_{n_D}^{(N)}(0) \right],$$
(27)

and  $\mathbf{Q}^{(N)}$  is a  $(n_D - n_N) \times (n_D - n_N)$  matrix with the following non-zero elements

$$q_{ij} = \begin{cases} -\lambda - N\mu & i = 1, \dots, n_D - n_N, j = i \\ \lambda & i = 1, \dots, n_D - n_N - 1, j = i + 1 \\ N\mu & i = 2, \dots, n_D - n_N, j = i - 1 \end{cases}$$
(28)

The computation of  $\pi^{(N)}(0)$  is described in the next section. From this stage, the only possible transition is to stage N - 1, i.e.,

$$P(S_{N-1} \mid S_N) = 1.$$
<sup>(29)</sup>

4) Stages  $K * (S_{K*})$ : For the stages with K active APs and one AP being powered on, the resulting Markov chain is illustrated in Fig. 5. In these stages there are no absorbing states that trigger the transition to other stages, since this happens when the amount of time spent in the stage is  $T_{on}$ . Because of this, the number of users that can be in the system during this stage varies between zero an infinity. Additionally, the service rate depends on the number of users



Fig. 5. Markov chain for the case of K active APs and one AP powering on.

as there should be at least one user per active AP for the total rate to be  $K\mu$ .<sup>6</sup>

For completeness, the time spent in a stage K \* is given by

$$H_{K*} = T_{on}.$$
 (30)

In this case, we need to obtain the expected time that the system spends in each state *i* during the  $T_{on}$  seconds that a stage K \* lasts,  $L_i^{(K*)}(T_{on})$ , and the probability of each state after  $T_{on}$ ,  $\pi_i^{(K*)}(T_{on})$ . These terms are required to compute the transition probabilities from stage K \* to the other stages and to obtain the performance figures of the system. The values for  $L_i^{(K*)}(T_{on})$  and  $\pi_i^{(K*)}(T_{on})$  can be obtained with the expressions of the transient analysis of an M/M/K queue, which are

$$L_i^{(K*)}(t) = \int_0^t \pi_i^{(K*)}(u) du, \qquad (31)$$

where  $\pi_i^{(K*)}(t)$  is the probability of being in state *i* at time *t*, which is determined by the fundamental equations of the CTMC

$$\frac{d\pi^{(K*)}(t)}{dt} = \pi^{(K*)}(t) \mathbf{Q}^{(K*)},$$
(32)

with  $\pi^{(K*)}(t) = [\pi_i^{(K*)}(t)]_i$  the transient state probability vector and  $\mathbf{Q}^{(K*)}$  the infinitesimal generator matrix of the CTMC. The non-zero elements of  $\mathbf{Q}^{(K*)}$  are

$$q_{ij} = \begin{cases} -\lambda - (i-1)\mu & i = 1, \dots, K, j = i \\ -\lambda - K\mu & i = K+1, \dots, j = i+1 \\ \lambda & i = 1, \dots, j = i+1 \\ (i-1)\mu & i = 2, \dots, K, j = i-1 \\ K\mu & i = K+1, \dots, j = i-1 \end{cases}$$
(33)

Note that to solve (31) and (32), we need again the vector of initial state probabilities  $\pi^{(K*)}(0)$ . On the other hand, as there are no closed expressions for the transient behaviour of an M/M/K queue, we need to use approximate methods (such as *uniformization* [24]) to solve it and compute  $L_i(T_{on})$  and  $\pi_i(T_{on})$ . Like in the previous cases, the computation of  $\pi^{(S_K*)}(0)$  is described in the next section.

Finally, as noted before, from stage  $K^*$  the system can go to any other stage K', with  $K' \leq K + 1$ , and to stage  $K + 1^*$ . In this way, after  $T_{on}$  the system can have K or less APs powered on, with the following probabilities

$$P(S_{K'} \mid S_{K*}) = \begin{cases} \sum_{i=0}^{n_2} \pi_i^{(K*)}(T_{on}), & K' = 1\\ \sum_{i=n_{K'}+1}^{n_{K'+1}} \pi_i^{(K*)}(T_{on}), & 1 < K' \le K \end{cases}$$
(34)

while the probability of having more APs on (or being powered on) after  $T_{on}$  depends on whether there are more APs to be

<sup>&</sup>lt;sup>6</sup>In fact, the CTMC corresponds to the classic M/M/K queue

former case, we have that

$$P(S_{K+1} \mid S_{K*}) = \sum_{i=n_{K+1}+1}^{N_{K+1}-1} \pi_i^{(K*)}(T_{on}), \qquad (35)$$

$$P(S_{K+1*} \mid S_{K*}) = \sum_{i=N_{k+1}}^{\infty} \pi_i^{(K*)}(T_{on}),$$
(36)

while for the case of K = N - 1 there are no more APs to activate, and therefore

$$P(S_N \mid S_{N-1*}) = \sum_{i=n_N+1}^{\infty} \pi_i^{(N-1*)}(T_{on}).$$
(37)

### C. Computing the Steady-State Distribution

To complete the analysis of the steady-state distribution of the semi-Markov process, we have to express the set of initial conditions for every stage in terms of the final state probabilities of the other stages. To this end, we denote  $\pi_i^{(T)}(0)$  as the probability that the initial state is i for the stage T (again we use T for generalization purposes, with stage T we refer indistinctly to any stage K, including stages 0 and N, or K\*). This probability can be computed with the law of total probability as

$$\pi_i^{(T)}(0) = \sum_{S_{T'} \in \mathcal{T}_T} \pi_i^{(T|T')}(0) P_T(T'), \qquad (38)$$

where

- $T_T$  is the set of stages that can reach stage T in one transition between stages,
- $\pi_{i}^{(T|T')}(0)$  is the probability that the initial state of the CTMC modelling stage T is i, given that the system was in stage T' and transitioned to state T,
- $P_T(T')$  is the probability that the system was in stage T' before the stage transition, given that it is now in stage T.

The set T can be easily derived for each stage from the Semi-Markov model described in Section III-A. Specifically we have,

$$\mathcal{T}_1 = \{S_2, S_{1*}, \dots S_{N-1*}\},\tag{39}$$

$$\mathcal{T}_K = \{S_{K+1}, S_{K-1*}, \dots S_{N-1*}\}$$
 for  $1 < K < N$ , (40)

$$\mathcal{T}_N = \{S_{N-1*}\},\tag{41}$$

$$\mathcal{T}_{1*} = \{S_1\}, \tag{42}$$

$$\mathcal{T}_{K*} = \{S_K, S_{K-1*}\} \text{ for } 1 < K < N-1.$$
 (43)

As an example, for the case of Fig. 1 with 4 APs, we have  $T_1 = \{S_2, S_{1*}, S_{2*}, S_{3*}\}, T_2 = \{S_3, S_{1*}, S_{2*}, S_{3*}\}, T_3 =$  $\{S_4, S_{2*}, S_{3*}\}, \ \mathcal{T}_4 = \{S_{3*}\}, \ \mathcal{T}_{1*} = \{S_1\}, \ \mathcal{T}_{2*} = \{S_2, S_{1*}\},$  $\mathcal{T}_{3*} = \{S_3, S_{2*}\}.$ 

The computation of  $\pi_i^{(T|T')}(0)$  depends on whether stage T' corresponds to a stage with an AP being powered on or not. For the latter case, the transition is triggered because the number of stations reached a (de)activation threshold (i.e., an absorbing state), and therefore we have

$$\pi_i^{(K*|K)}(0) = \begin{cases} 1, & i = N_K \\ 0, & \text{otherwise} \end{cases}$$
(44)

powered on, i.e., if K < N - 1, or not (K = N - 1). For the for 1 < K < N (note that we have included here the transition from stage 1 to stage 1\* as well), and

$$\pi_i^{(K-1|K)}(0) = \begin{cases} 1, & i = n_K \\ 0, & \text{otherwise} \end{cases}$$
(45)

for  $1 < K \leq N$  (we have included the transition from stage N to stage N - 1 as well). On the other hand, when stage T' is a K\* stage, there are multiple states that can result in a transition to a stage, which results in the following cases: (i) If the transition is to stage 1 ( $S_T = S_1$ ), then

$$\pi_i^{(1|K*)}(0) = \begin{cases} \frac{\pi_i^{(K*)}(T_{on})}{\sum_{j=0}^{n_2} \pi_j^{(K*)}(T_{on})}, & 0 \le i \le n_2 \\ 0, & n_2 < i < N_1 \end{cases}$$
(46)

(ii) If the transition is to a stage  $1 < K' \leq K$ , then

$$\pi_i^{(K'|K*)}(0) = \begin{cases} \frac{\pi_i^{(K*)}(T_{on})}{\sum_{j=n_{K'}+1}^{n_{K'+1}} \pi_j^{(K*)}(T_{on})}, & n_{K'} < i \le n_{K'+1} \\ 0, & n_{K'+1} < i < N_{K'} \end{cases}$$
(47)

(iii) If K < N - 1 (i.e.  $S_{T'} \neq S_{N-1}$ ) and the transition is to stage K + 1, then

$$\pi_i^{(K+1|K*)}(0) = \frac{\pi_i^{(K*)}(T_{on})}{\sum_{j=n_{K+1}+1}^{N_{K+1}-1} \pi_j^{(K*)}(T_{on})}.$$
(48)

(iv) If K < N - 1 (again  $S_{T'} \neq S_{N-1}$ ) and the transition is to stage K + 1\*, then

$$\pi_i^{(K+1|K*)}(0) = \begin{cases} 0, & 0 \le i \le N_{K+1} \\ \frac{\pi_i^{(K*)}(T_{on})}{\sum_{j=N_{K+1}}^{\infty} \pi_j^{(K*)}(T_{on})}, & i > N_{K+1} \end{cases}$$
(49)

(v) If K = N - 1 and the transition is to stage N, then

$$\pi_i^{(N|N-1*)}(0) = \frac{\pi_i^{(N-1*)}(T_{on})}{\sum_{j=n_N+1}^{\infty} \pi_j^{(K*)}(T_{on})}.$$
 (50)

Finally, the computation of  $P_T(T')$  can be done with the law of total probability again

$$P_T(T') = \frac{P(S_T | S_{T'})\phi_{T'}}{\sum_{S_Q \in \mathcal{I}_T} P(S_T | S_Q)\phi_Q},$$
(51)

where  $P(S_T \mid S_{T'})$  denotes the stage transition probability computed in (11), (17), (18), (29), (34)-(37), and  $\phi_T$  is the stationary probability of stage T in the embedded Markov chain of the semi-Markov process. The computation of  $\phi_T$  is done via the system

$$\boldsymbol{\phi} = \boldsymbol{\phi} \mathbf{P},\tag{52}$$

where  $\phi$  is a row vector whose components are the values of  $\phi_T$ , and **P** is a matrix composed of the stage transition probabilities of the embedded Markov chain.

With the above, we have completed the analysis that enables in the next section the computation of the steady state probabilities of the semi-Markov model. We also address there how to compute performance figures based on these probabilities.

#### D. Performance Figures

We characterise the performance of the system with two figures:

- The average power consumed by the infrastructure P.
- The average service time of a user  $T_s$ .

The average power consumed by the infrastructure can be expressed in terms of the average number of APs that are powered on,  $N_{AP}$ , as follows

$$P = N_{AP} P_{AP}. (53)$$

 $N_{AP}$  is computed as the weighted sum of the number of APs powered on in each stage times the probability of being in that stage

$$N_{AP} = \sum_{K=1}^{N} K \mathcal{P}_{K} + \sum_{K=1}^{N-1} (K+1) \mathcal{P}_{K*},$$
 (54)

where  $\mathcal{P}_K$  and  $\mathcal{P}_{K*}$  are the stationary probabilities of the stages of the semi-Markov process, i. e., the probability of being in stage K (including stages 0 and N) or K\* at a specific moment. These probabilities are related to the stage probabilities of the embedded Markov chain as follows

$$\mathcal{P}_T = \frac{H_T \phi_T}{\sum_{K=1}^N \phi_K H_K + \sum_{K=1}^{N-1} \phi_{K*} H_{K*}}.$$
 (55)

The average service time  $T_s$ , which corresponds to the time between the instant when a user generates a service request and when this request is completely served, can be obtained via Little's formula

$$T_s = \frac{N_u}{\lambda},\tag{56}$$

with  $N_u$  the average number of users in the system. This can be computed with the law of total probability as follows

$$N_{u} = \sum_{i=1}^{\infty} i \left( \sum_{K=1}^{N} \pi_{i}^{(K)} \mathcal{P}_{K} + \sum_{K=1}^{N-1} \pi_{i}^{(K*)} \mathcal{P}_{K*} \right), \quad (57)$$

where  $\pi_i^{(K)}$  and  $\pi_i^{(K*)}$  are the average probabilities of having *i* users, given that the system is in stage K or K\*, respectively. This can be computed, for each type of stage, as

$$\pi_i^{(K)} = \frac{L_i^{(K)}(\infty)}{H_K},$$
(58)

$$\pi_i^{(K*)} = \frac{L_i^{(K*)}(T_{on})}{T_{on}}.$$
(59)

As can be seen, all the performance metrics depend on the variables  $\phi_T$ ,  $H_T$ ,  $L_i^{(K)}(\infty)$  and  $L_i^{(K*)}(T_{on})$ , whose relationships have been described through Sections III-B to III-C. In order to obtain an exact solution for them, we should solve a system of non-linear equations with the additional problem that there are no closed expressions for the transient analysis of the CTMC modelling stages K\*. To solve this, we propose the iterative algorithm described in Algorithm 1. In this algorithm, the initial values of  $\pi^{(K)}(0)$  can be set assuming that all the states with non-zero probabilities according to (46)-(50) have the same initial probability. Regarding  $\pi^{(K*)}(0)$ , a good starting guess is to assume that  $\pi_i^{(K*)}(0) = 1$  for  $i = N_K$  and 0

Algorithm 1 Solution to the Exact Model

1: Set initial estimations of  $\pi^{(K)}(0)$  and  $\pi^{(K*)}(0)$ 

- 2: repeat Compute  $\mathbf{L}^{(K)}(\infty)$  with (7), (13) and (25)
- 3: Obtain  $H_K$  with (6), (12) and (24) 4:
- Solve (31)-(32) to obtain  $\mathbf{L}^{(K*)}(T_{on})$  and  $\boldsymbol{\pi}^{(K*)}(T_{on})$ 5:
- Compute P(T | T') with (17)-(20) and (34) (37) 6:
- Solve (52) to obtain  $\phi$ 7:
- 8:
- 9:
- Obtain  $P_T(T')$  with (51) Compute  $\pi_i^{(T|T')}(0)$  with (44)-(50) Update  $\pi^{(K)}(0)$  and  $\pi^{(K*)}(0)$  with (38) 10:
- 11: **until** Stopping criterion is met 12: Obtain  $\pi_i^{(K)}$  and  $\pi_i^{(K*)}$  with (58) and (59)
- 13: Compute  $\mathcal{P}_T$  with (55)
- 14: Obtain  $N_u$  with (57) and  $N_{AP}$  with (54)
- 15: Compute  $T_s$  with (56) and P with (53)

otherwise (this is what would happen if  $T_{on} = 0$ ). Finally, a common stopping criterion is that the norm of the vector difference between the old and updated version of vectors  $\pi^{(K)}(0)$  and  $\pi^{(K*)}(0)$  is below a threshold  $\epsilon$ .

#### **IV. SIMPLIFIED ANALYSIS**

#### A. Motivation and Simplification

The main weaknesses of the model derived in the previous section is that the initial probabilities of a stage depends on the "final" probabilities of the rest of stages, which depend in turn of their initial probabilities. This causes a loop that requires the use of an iterative algorithm with non-negligible computational complexity as the one described above. We next describe how to simplify the analytical model of Section III to enable an efficient computation of the performance figures, at the cost of some numerical inaccuracy.

The proposed simplification affects exclusively the transitions from stages K\*. As can be seen in Fig. 1, from these stages the system could go to stage K + 1\* or any other stage K', with  $K' \leq K + 1$ . The direct transitions between stages K\* make that the initial state probabilities for these stages  $\pi_i^{(K*)}(0)$  could be non-zero for  $i \ge N_K$ . To break the coupling between stages K\*, we assume that

the initial state probabilities of stages K \* are fixed and equal to

$$\pi_i^{(S_{K*})}(0) = \begin{cases} 1, & i = N_K \\ 0, & \text{otherwise} \end{cases}$$
(60)

This implies that the system enters into stages K \* always with  $N_K$  users. This assumption holds as long as the transition probability between a stage K \* and the stage K + 1 \* is small, which is true for typical  $T_{on}$  values.

Once this assumption is made, the transition probabilities from stages K \* to other stages are fixed and independent of the initial state probabilities of the rest of stages. Now, we also have to tackle the same apparent coupling for the initial state probabilities of stages K. To solve this, we build a new semi-Markov model derived from the one depicted in Fig. 1 substituting stages K by the embedded DTMC of their



Fig. 6. Simplified model.

corresponding CTMC. The description of this new model is performed in the next Section.

#### B. Model Description

Fig. 6 shows the embedded DTMC of the semi-Markov process described above. The leftmost states, which model stages *K* (including 0 and *N*), are defined by the pair (*i*, *K*), with *i* the number of users in the system and *K* the number of powered-on APs. The holding time of these states is an exponential random variable with mean  $(\lambda + K\mu)^{-1}$ . The rightmost states model the stages  $K^*$  and their holding time is constant and equal to  $T_{on}$ . To keep a uniform notation, we note these states as (\*, *K*). The non-null transitions probabilities are described in (61), as shown at the bottom of the next page.

The first two equations model the transitions between states of Stage 1, the first one corresponds to the departure of a user and the second one its arrival. The third and fourth equations model the transitions between states of stages  $2 \le N \le N-1$  and the fifth and sixth the transitions between states of stage N. The seventh equation corresponds to the switch off of an AP when a user departures and stage K remains with  $n_K$  users, which triggers the transition to stage K - 1. The eighth equation models the switching on of a new AP (i.e. the transition to stage K\*) when the  $N_K$ -th user arrives and K APs are on. Note that in all the cases the transition probabilities only depend on the parameters  $\lambda$ ,  $\mu$  and the number of APs that are serving traffic at the moment of the transition.

The next equations model the transitions from states (\*, K), (i.e., from stages K\*). Now the transition probabilities are of the form  $\pi_i^{(K*)}(T_{on})$  and can be computed solving (31) and (32) assuming the initial state probabilities given in (60). Specifically, the ninth equation corresponds to the transition from stage K \* to stage K + 1 \* because the system reaches  $N_{K+1}$  users during the booting-up of the K + 1-th AP. The tenth equation models the transition from stage K \* to a state where there are K + 1 APs powered on and a number of users ranging between  $n_{K+1} + 1$  and  $N_{K+1} - 1$ . The eleventh equation is similar to the previous one but for stage N - 1\*. In this case, there is no upper limit in the number of users since there is not any remaining AP to boot up. The twelfth equation models the transition from stage K \* to states where the number of APs on is below K + 1. This implies that during the booting up of the K + 1-th AP several users have left forcing the system to switch off some APs. The last equation is similar to the previous one and corresponds to transitions to states where only one AP is on.

With the previous equations, the DTMC can be easily solved to obtain the stationary distribution of the state probabilities, that we name P(i, K) (or P(\*, K)) hereafter. With these, the stationary probability of each state of the semi-Markov process,  $\Phi(i, K)$  (or  $\Phi(*, K)$ ) are

$$\Phi(i,1) = \frac{P(i,1)}{\Omega \cdot (\lambda + \mu)}, \quad 0 \le i < N_1$$
(62)

$$\Phi(i, K) = \frac{P(i, K)}{\Omega \cdot (\lambda + K\mu)}, \quad n_K < i < N_K, \ 1 < K < N$$
(63)

$$\Phi(i, N) = \frac{P(i, N)}{\Omega \cdot (\lambda + N\mu)}, \quad i > n_N \tag{64}$$

$$\Phi(*,K) = \frac{P(*,K)T_{on}}{\Omega}, \quad 1 \le K < N$$
(65)

with

$$\Omega = \sum_{j=0}^{N_1-1} \frac{P(j,1)}{\lambda+\mu} + \sum_{K'=2}^{N_1-1} \sum_{j=n_{K'}+1}^{N_{K'}-1} \frac{P(j,K')}{\lambda+K'\mu} + \sum_{j=n_N+1}^{\infty} \frac{P(j,N)}{\lambda+N\mu} + \sum_{K'=1}^{N-1} P(*,K')T_{on}.$$
 (66)

The stationary probabilities of stages K \* are directly  $\mathcal{P}_{K*} = \Phi(*, K)$ , while for stages K we have

$$\mathcal{P}_{1} = \sum_{i=0}^{N_{1}-1} \Phi(i, 1), \tag{67}$$

# Algorithm 2 Solution to the Approximate Model

- 1: Set  $\pi^{(K*)}(0)$  with (60)
- 2: Solve (31) and (32) to obtain  $\mathbf{L}^{(K*)}(T_{on})$  and  $\boldsymbol{\pi}^{(K*)}(T_{on})$
- 3: Solve the DTMC with transitions given by (61) to obtain P(i, K) and P(\*, K)
- 4: Compute (62)-(66) to obtain  $\Phi(i, K)$  and  $\Phi(*, K)$
- 5: Obtain  $\mathcal{P}_1$ ,  $\mathcal{P}_K$  and  $\mathcal{P}_N$  with (67)-(69)
- 6: Obtain  $N_u$  with (70) and  $N_{AP}$  with (54)
- 7: Compute  $T_s$  with (56) and P with (53)

$$\mathcal{P}_K = \sum_{i=n_K+1}^{N_K-1} \Phi(i, K), \tag{68}$$

and

$$\mathcal{P}_N = \sum_{i=n_N+1}^{\infty} \Phi(i, N).$$
(69)

Once these terms are known, we can compute the average power P with (55) and (54). The average service time  $T_s$  is obtained with (56) as well, but in this case  $N_u$  is

$$N_{u} = \sum_{i=1}^{\infty} i \left( \sum_{K=1}^{N} \Phi(i, K) + \sum_{K=1}^{N-1} \pi_{i}^{(K*)} \mathcal{P}_{K*} \right), \quad (70)$$

with  $\pi_i^{(K*)}$  the same as in (59).

To end this Section, we present in *Algorithm 2* the different steps required to obtain the performance figures of the system. As can be seen, in this case we avoid the presence of loops.

## V. NUMERICAL RESULTS

We next present a numerical evaluation of a RoD system in terms of the performance figures considered, namely, the average service time  $T_s$  and the power consumed by the infrastructure P. To this end, we compute these two variables for a variety of scenarios, these being defined in terms of the network load or the configuration of the RoD scheme (given by the parameters M,  $\rho_h$ ,  $\rho_l$ ). ted with Octave. In the simulation results presented, we compare the results of our approximate model against the ones obtained via simulation,<sup>7</sup> while in Section V-C we assess the computational complexity of this model against the accurate one.

Throughout all simulations, we consider the following scenario<sup>8</sup>: (i) various APs can be simultaneously activated (instead of only one, as assumed in the analysis); (ii) there is no complete overlap of the coverage areas: we assume a deployment centred around one AP with a 10 m coverage radio that is always on, and N-1 APs with the same coverage radios that are randomly deployed within a 4 m circle centred around the first AP and that will be opportunistically (de)activated; and (iii) users are not static but follow the classical random waypoint model [26], selecting a novel destination at random after reaching the previous one, and moving at a speed that is randomly chosen between 0.3 and 0.7 m/s. We further assume that there are up to N = 10 APs available,<sup>9</sup> that a single AP consumes 3.5 W when active (which corresponds to the average power consumed by a Linksys device [19]) and zero otherwise, and that  $\mu = 0.1 \text{ s}^{-1}$ .

# A. Impact of Network Load

We first analyse the power consumption as the network load  $\rho = \lambda/(N\mu)$  varies. To this end, we fix a target distribution of M = 5 users per AP and the following two configurations of the (de)activation thresholds { $\rho_h$ ,  $\rho_l$ }: {100%, 30%} and {50%, 25%}, the former being more "reluctant" to increase the number of APs when the network load increases. To understand the impact of  $T_{on}$  on performance, we consider the cases of zero and 30 s start-up times. We plot the computed

<sup>7</sup>Our approximate model is solved numerically using Octave (https://www.gnu.org/software/octave/), while simulation results are obtained from a discrete event simulator written in C++.

 $^{8}$ Note that this scenario relaxes some of the simplifying assumptions behind our model, and thus allows to assess the impact of such assumptions on the results.

<sup>9</sup>This is a reasonable number for dense scenarios: for instance, an auditorium with 360 users, each of them demanding 3 Mbps for HD video, would require 31 802.11n APs with a throughput of 35 Mbps (data taken from [18]). Results of the same order of magnitude are obtained in [27] and [28].

$$\begin{cases} P(i-1, 0 \mid i, 0) = \frac{\mu}{\lambda + \mu} & i = \{1, \dots, N_1 - 1\} \\ P(i+1, 0 \mid i, 0) = \frac{\lambda}{\lambda + \mu} & i = \{0, \dots, N_1 - 2\} \\ P(i-1, K \mid i, K) = \frac{K\mu}{\lambda + K\mu} & i = \{n_K + 2, \dots, N_K - 1\}, K = \{2, \dots, N - 1\} \\ P(i+1, K \mid i, K) = \frac{\lambda}{\lambda + K\mu} & i = \{n_K + 1, \dots, N_K - 2\}, K = \{2, \dots, N - 1\} \\ P(i-1, N \mid i, N) = \frac{N\mu}{\lambda + N\mu} & i = \{n_N + 2, \dots\} \\ P(i+1, N \mid i, N) = \frac{\lambda}{\lambda + N\mu} & i = \{n_N + 1, \dots\} \\ P(n_K, K-1 \mid n_K + 1, K) = \frac{K\mu}{\lambda + K\mu} & K = \{2, \dots, N\} \\ P(*, K \mid N_K - 1, K) = \frac{\lambda}{\lambda + K\mu} & K = \{1, \dots, N - 1\} \\ P(*, K + 1 \mid *, K) = \sum_{i=N_{K+1}}^{\infty} \pi_i^{(K*)}(T_{on}) & K = \{1, \dots, N - 2\} \\ P(i, K+1 \mid *, K) = \pi_i^{(K*)}(T_{on}) & i = \{n_K + 1, \dots, N_{K+1} - 1\}, K = \{1, \dots, N - 2\} \\ P(i, N \mid *, N - 1) = \pi_i^{(N-1*)}(T_{on}) & i = \{n_K + 1, \dots, N_{K+1} - 1\}, K = \{2, \dots, N - 1\}, K' = \{2, \dots, K\} \\ P(i, 1 \mid *, K) = \pi_i^{(K*)}(T_{on}) & i = \{n_K + 1, \dots, n_{K'+1}\}, K = \{2, \dots, N - 1\}, K' = \{2, \dots, K\} \\ P(i, 1 \mid *, K) = \pi_i^{(K*)}(T_{on}) & i = \{0, \dots, n_2\}, K = \{1, \dots, N - 1\} \end{cases}$$



Fig. 7. Average power consumption vs. network load.

figures of P in Fig. 7, where we use squares for the simulation values (average of 10 simulation runs, each consisting of more than 100k users) and lines for the analysis.

According to the results, the power consumption is monotonously increasing with the network load, with the analysis practically coinciding with the simulation values, with some minor deviations (approx. 1.8%) for high loads (we depict a zoomed version of the figure for these values). Considering the relative performance of each configuration, for the case of  $T_{on} = 0$  the results overlap, while for the case of  $T_{on} = 30$  s, the policy that is "more eager" to power APs leads to higher power consumption.

We next analyse the performance in terms of service time and the trade-off with power consumption. To this end, we plot  $T_s$  vs. P in Fig. 8, with each simulation point corresponding to a different value of  $\rho$ , which varies from 0.05 to 0.9 in steps of 0.05. Here we also provide for comparison the "ordinary" case of no RoD scheme (all APs always on), which leads to the smallest service times and the largest power consumptions. As in the previous case, the analysis accurately predicts simulation results, with differences below 2.5%. The figure also illustrates that the service time is a monotonous increasing function of the load: steep for  $\rho \leq 0.3$ , which is caused by the "drastic" impact of powering on an AP when the number of active resources is relatively low, and then more gradual until  $\rho \approx 0.9$ . Concerning the impact of the considered configurations, for the same value of the  $\{\rho_l, \rho_h\}$  parameters, the non-zero start time has an impact of approx. 5 s for the more dynamic configuration, and approx. 2 s for the more "reluctant" configuration, while the impact of the activation policy results in differences of approx. 12 s.

#### B. Impact of RoD Configuration

Next, we consider the case of a fixed value of  $\rho = 0.5$ , and compute the service time and power consumed for the two considered { $\rho_h$ ,  $\rho_l$ } configurations and different values of the target number of users per AP *M*. We plot the service time and the power consumption as a function of *M*, with the results being depicted in Fig. 9.



Fig. 8. Average service time vs. average power consumption.



Fig. 9. Average service time (top) and power consumption (bottom) vs. target number of users per AP.

For the case of the service time (Fig. 9, top), again simulation results practically coincide with the analysis. The larger Mis, the longer the service times are, as users are more likely to share the capacity of a single AP before activating new resources. In fact, the relation is practically linear, e.g., when M changes from 5 to 10, the service time doubles for all considered scenarios: as there are, on average, more users per AP, the service times will be longer.

For the case of the power consumption, the resulting values are depicted in Fig. 9 (bottom). Here, we note that the results for both RoD configurations for  $T_{on} = 0$  overlap, and result in a constant power consumption regardless of the value of M. The reason for this behaviour is that, as M increases, more users per AP are required to power on additional resources, but also longer service times will result, leading to more users in the system. In fact, the power consumption of 17.5 W implies that, on average, 5 out of the 10 available APs are on, which matches the  $\rho = 0.5$  load. When the start-up times are nonzero, there is a small reduction of P as M increases. The reason for this is that, on average, the system is less likely to power on additional APs, which incurs in the overhead of the

TABLE I Relative Differences and Computational Times of the Exact and Simplified Analyses

$T_{on}$ (s)	$ ho_h, ho_l$	ρ	Error		Comp. time (s)	
			$\Delta T_s$	$  \Delta P$	Exact	Simpl.
0	[0.5, 0.75]	0.25	$\approx 0\%$	$\approx 0\%$	104.28	2.35
	$\{0.3, 0.75\}$	0.75	pprox 0%	$\approx 0\%$	102.34	2.35
	{1, 0.7}	0.25	pprox 0%	$\approx 0\%$	102.78	2.40
		0.75	0.06%	0.04%	104.58	2.30
15	{0.5, 0.75}	0.25	0.22%	0.14%	167.69	3.62
		0.75	0.91%	0.58%	269.43	5.77
	{1, 0.7}	0.25	0.02%	0.01%	166.94	3.9
		0.75	0.17%	0.14%	276.20	5.65
30	{0.5, 0.75}	0.25	1.97%	1.05%	320.52	6.75
		0.75	3.09%	1.80%	519.96	11.36
	{1, 0.7}	0.25	0.41%	0.23%	316.62	7.03
		0.75	1.17%	0.66%	541.04	11.72

start-up process. Finally, we also note that simulation are very close to those from the analysis, with relative differences of approx. 4% (the smaller M is, the larger the differences are, as the impact of non-perfect overlap of coverage areas is more noticeable for a small number of users).

# C. Computational Complexity

We next estimate the computational complexity of obtaining the numerical solution for the exact and the simplified analysis. To this end, we assume a scenario with N = 10 APs, fix M = 4, and consider different configurations of  $T_{on}$ ,  $\rho$ , and  $\{\rho_h, \rho_l\}$  parameters. For each set of parameters, we compute the average service time  $T_s$  and power consumption P using the exact and the simplified analysis, as well as the time required to compute these values for each case. We note that we use Octave to compute the numerical solution for these analysis, running over an Intel<sup>®</sup> Xeon<sup>®</sup> X5550 @2.67GHz with 48 GB RAM, and therefore our comparison serves to illustrate the relative differences in complexity, and not absolute values.

We provide in Table I the results of the above computation. More specifically, we provide in the Table, for each considered configuration, the relative difference between the two analyses in terms of service time (denoted as  $\Delta T_s$ ) and power consumption (denoted as  $\Delta P$ ), and the corresponding computation times. There are two main observations from the results: (*i*) on the one hand, for both power and service time figures, the resulting differences between the numerical analyses are at most 3%, and in many cases well below 1%; and (*ii*) on the other hand, for the computational times, there are two orders of magnitude of difference between them in all but for two cases. Finally, it is also worth noting that, for the case of the exact analysis, computational times grow with  $T_{on}$ , which confirms to some extent that the K\* stages are responsible for the computational burden.

# D. Realistic Traffic Model

To analyse the impact of the simplifying assumptions on the traffic model of our analysis, in the following we compare the results obtained with our analysis against those obtained from simulations with a "realistic" traffic model. In particular, we



Fig. 10. Average delay vs. power consumed with non-exponential service demands.

follow [29] and assume that when a station joins the WLAN, it performs a random number of download requests that follows a BiPareto distribution. The length of each download also follows a BiPareto distribution, and the interarrival time of requests follows a lognormal distribution. We fix the average number of requests to 10, with the following parameters of the BiPareto distribution:  $\alpha = 0.06$ ,  $\beta = 1.73$ , c = 6.61 and k = 1; the lognormal distribution is simulated with parameters  $\mu = 0.34$  and  $\sigma = 0.63$ ; and the request lengths are initially modelled with parameters  $\alpha = 0.0$ ,  $\beta = 2.13$ , c = 20.0 and k = 1.5 (which leads to an average download size of 30 MB), while the user arrival rate is Poissonian at a rate ranging from 0.05 to 0.9 s<sup>-1</sup>.

We show in Fig. 10 the resulting average service time vs. power consumption for different configurations of the RoD scheme and values of  $T_{on}$ , where (like in Fig. 8) we vary the load from 0.05 to 0.9 in steps of 0.05. We observe that the accuracy of the model worsens as the service rate increases: the deviations are smaller than 5% for  $\rho < 0.8$  but notably higher as the system gets closer to saturation. We conclude from these results that overall the accuracy of the model is reasonable for the range of loads of interest (i.e., sufficiently far from congestion).

# E. Optimal Configuration of a RoD Scheme

While the exact analysis incurs in a notable complexity, we have seen that the simplified analysis is able to compute the performance figures of a RoD scheme in an affordable manner while keeping a notable accuracy. In this way, it can be used, for instance, to compute the optimal configuration of a RoD algorithm, given a set of estimated network conditions, these being expressed in terms of  $\lambda$  and  $\mu$ . In the following, we present one example of such configuration algorithms, although we restrict ourselves for simplicity to the considered RoD policy (although there could be many others) and a simple optimization criterion. Our optimization scheme works as follows. Given an estimation of the network conditions, we set a bound on the maximum service time  $T_{max}$ , and perform

TABLE II Performance an Optimal Configurations of a RoD Scheme

$\mu(s^{-1})$	ρ	$T_{on}(s)$	Μ	$\rho_h$	$\rho_l$	$T_s(s)$	P(W)
0.05	0.25	0	3	1.20	0.55	75.93	8.76
		15	4	0.75	0.30	79.17	8.96
		30	3	1.20	0.30	74.77	9.16
	0.5	0	3	1.20	0.30	76.89	17.33
		15	3	1.15	0.30	76.60	17.55
		30	3	1.15	0.30	78.41	17.81
	0.75	0	3	1.20	0.30	77.44	25.35
		15	3	1.20	0.30	79.75	26.00
		30	2	0.95	0.45	53.00	25.34
0.10	0.25	0	3	1.20	0.55	37.96	8.76
		15	3	1.20	0.30	37.38	9.16
		30	3	1.20	0.30	39.98	9.50
	0.5	0	3	1.20	0.30	38.44	17.33
		15	3	1.15	0.30	39.21	17.81
		30	2	0.95	0.45	28.41	17.88
	0.75	0	3	1.20	0.30	38.72	25.35
		15	2	0.95	0.45	26.50	26.34
		30	2	0.95	0.45	28.20	25.98
0.20		0	3	1.20	0.55	18.98	8.76
	0.25	15	3	1.20	0.30	19.99	9.50
		30	3	0.80	0.30	19.87	10.12
	0.5	0	3	1.20	0.30	19.22	17.33
		15	2	0.95	0.45	14.21	17.88
		30	2	1.00	0.45	16.36	17.61
	0.75	0	3	1.20	0.30	19.36	25.35
		15	2	0.95	0.45	14.10	25.98
		30	2	0.95	0.45	15.90	25.42

a sweep on the configuration space  $\{M, \rho_h, \rho_l\}$  to look for the configuration that minimises power while guaranteeing an average service time  $T_s$  below  $T_{\text{max}}$ . In our search, Mgoes from 2 to 10 in steps of one, while  $\rho_h$  and  $\rho_l$  go from 0.05 to 1.25 in steps of 0.05.

The configuration resulting from this search and the corresponding performance figures are given in Table II for three different service rates  $\mu = \{0.05, 0.1, 0.2\} s^{-1}$  and the corresponding three service time bounds  $T_{\text{max}} = \{80, 40, 20\} s$ , respectively. If we compare the consumed power with a reference scenario of the 10 APs always on (i.e., consuming 35 W), the reduction is quite considerable, ranging between 25% and 75% depending on the network load. Finally, it is also worth remarking that  $T_{on}$  has a non-negligible effect on the resulting configuration parameters.

## VI. CONCLUSIONS

Resource-on-Demand schemes are required in dense networks to adapt to the varying load while maintaining an energy efficient performance. In this paper, we have developed an analytical model of these schemes that, in contrast to previous publications, accounts for the non-zero start-up times of real hardware. We have also presented a simplified model, whose computational times are appox. 50x shorter while maintaining relative errors below 3%. We have illustrated the practicality of this simplified model with a simple algorithm to derive the optimal configuration of a RoD scheme.

#### References

 A. Fehske, G. Fettweis, J. Malmodin, and G. Biczok, "The global footprint of mobile communications: The ecological and economic perspective," *IEEE Commun. Mag.*, vol. 49, no. 8, pp. 55–62, Aug. 2011.

- [2] Y. Chen et al., "Fundamental trade-offs on green wireless networks," IEEE Commun. Mag., vol. 49, no. 6, pp. 30–37, Jun. 2011.
- [3] C. Han et al., "Green radio: Radio techniques to enable energy-efficient wireless networks," *IEEE Commun. Mag.*, vol. 49, no. 6, pp. 46–54, Jun. 2011.
- [4] G. Lim, C. Xiong, L. J. Cimini, and G. Y. Li, "Energy-efficient resource allocation for OFDMA-based multi-RAT networks," *IEEE Trans. Wireless Commun.*, vol. 13, no. 5, pp. 2696–2705, May 2014.
- [5] L. Budzisz et al., "Dynamic resource provisioning for energy efficiency in wireless access networks: A survey and an outlook," *IEEE Commun. Surveys Tut.*, vol. 16, no. 4, pp. 2259–2285, 4th Quart., 2014.
- [6] P. Serrano, A. de la Oliva, P. Patras, V. Mancuso, and A. Banchs, "Greening wireless communications: Status and future directions," *Comput. Commun.*, vol. 35, no. 14, pp. 1651–1661, 2012.
- [7] Y. S. Soh, T. Q. S. Quek, M. Kountouris, and H. Shin, "Energy efficient heterogeneous cellular networks," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 5, pp. 840–850, May 2013.
- [8] J. Wu, Y. Zhang, M. Zukerman, and E. K. N. Yung, "Energy-efficient base-stations sleep-mode techniques in green cellular networks: A survey," *IEEE Commun. Surveys Tut.*, vol. 17, no. 2, pp. 803–826, 2nd Quart., 2015.
- [9] W. Sun, H. Li, and J. Wu, "Study on real energy consumption of large-scale campus wireless network," in *Proc. Int. Conf. Comput. Netw. Commun. (ICNC)*, Jan. 2013, pp. 605–609.
- [10] F. Ganji, L. Budzisz, and A. Wolisz, "Assessment of the power saving potential in dense enterprise WLANs," in *Proc. IEEE 24th Int. Symp. Pers. Indoor Mobile Radio Commun. (PIMRC)*, Sep. 2013, pp. 2835–2840.
- [11] F. Ganji et al., "Greening campus WLANs: Energy-relevant usage and mobility patterns," Comput. Netw., vol. 78, pp. 164–181, Feb. 2015.
- [12] F. G. Debele, M. Meo, D. Renga, M. Ricca, and Y. Zhang, "Designing resource-on-demand strategies for dense WLANs," *IEEE J. Sel. Areas Commun.*, vol. 33, no. 12, pp. 2494–2509, Dec. 2015.
- [13] M. A. Marsan, L. Chiaraviglio, D. Ciullo, and M. Meo, "A simple analytical model for the energy-efficient activation of access points in dense WLANs," in *Proc. 1st Int. Conf. Energy–Efficient Comput. Netw.*, May 2010, pp. 159–168.
- [14] A. P. C. da Silva, M. Meo, and M. A. Marsan, "Energy-performance trade-off in dense WLANs: A queuing study," *Comput. Netw.*, vol. 56, no. 10, pp. 2522–2537, 2012.
- [15] R. G. Garroppo, G. Nencioni, G. Procissi, and L. Tavanti, "The impact of the access point power model on the energy-efficient management of infrastructured wireless LANs," *Comput. Netw.*, vol. 94, pp. 99–111, Jan. 2016.
- [16] A. P. Jardosh, K. Papagiannaki, E. M. Belding, K. C. Almeroth, G. Iannaccone, and B. Vinnakota, "Green WLANs: On-demand WLAN infrastructures," *Mobile Netw. Appl.*, vol. 14, no. 6, pp. 798–814, 2009.
- [17] J. Ortín, P. Serrano, and C. Donato, "Optimal configuration of a resource-on-demand 802.11 WLAN with non-zero start-up times," *Comput. Commun.*, to be published. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S014036641630161X, DOI: /10.1016/j.comcom.2016.04.022.
- [18] J. Florwick, J. Whiteaker, A. C. Amrod, and J. Woodhams, "Wireless LAN design guide for high density client environments in higher education," Cisco, San Jose, CA, USA, Tech. Rep. [Online]. Available: http://www.cisco.com/c/dam/en\_us/solutions/industries/docs/education/ cisco\_wlan\_design\_guide.pdf
- [19] P. Serrano, A. Garcia-Saavedra, G. Bianchi, A. Banchs, and A. Azcorra, "Per-frame energy consumption in 802.11 devices and its implication on modeling and design," *IEEE/ACM Trans. Netw.*, vol. 23, no. 4, pp. 1243–1256, Aug. 2015.
- [20] M. Papadopouli, H. Shen, and M. Spanakis, "Modeling client arrivals at access points in wireless campus-wide networks," in *Proc. 14th IEEE Workshop Local Metropolitan Area Netw. (LANMAN)*, Apr. 2005, pp. 6–12.
- [21] V. Paxson and S. Floyd, "Wide area traffic: The failure of Poisson modeling," *IEEE/ACM Trans. Netw.*, vol. 3, no. 3, pp. 226–244, Jun. 1995.
- [22] A. Feldmann, A. C. Gilbert, W. Willinger, and T. G. Kurtz, "The changing nature of network traffic: Scaling phenomena," *SIGCOMM Comput. Commun. Rev.*, vol. 28, no. 2, pp. 5–29, Apr. 1998.
- [23] G. R. Hiertz, D. Denteneer, L. Stibor, Y. Zang, X. P. Costa, and B. Walke, "The IEEE 802.11 universe," *IEEE Commun. Mag.*, vol. 48, no. 1, pp. 62–70, Jan. 2010.
- [24] G. Bolch, S. Greiner, H. de Meer, and K. S. Trivedi, *Queueing Networks and Markov Chains*. Hoboken, NJ, USA: Wiley, 2006.
- [25] O. C. Ibe, Markov Processes for Stochastic Modeling. New York, NY, USA: Academic, 2009.

- [27] "High Density Wi-Fi Deployment Guide (CVD)," Cisco, San Jose, CA, USA, Tech. Rep. [Online]. Available: https://documentation.meraki.com/MR/WiFi\_Basics\_and\_Best\_Practices /High\_Density\_Wi-Fi\_Deployment\_Guide\_(CVD)
- [28] "High-density wireless networks for auditoriums validated reference design," Aruba Netw., Sunnyvale, CA, USA, Tech. Rep. [Online]. Available: https://community.arubanetworks.com/aruba/attachments /aruba/Aruba-VRDs/21/1/High-Density%20Wireless%20Networks%20 for%20Auditoriums.pdf
- [29] F. Hernández-Campos, M. Karaliopoulos, M. Papadopouli, and H. Shen, "Spatio-temporal modeling of traffic workload in a campus WLAN," in *Proc. 2nd Annu. Int. Workshop Wireless Internet (WICON)*, Boston, MA, USA, Aug. 2006, Art. no. 1.



**Jorge Ortín** received the Telecommunication Engineering and Ph.D. degrees from the Universidad de Zaragoza in 2005 and 2011, respectively. He was a Post-Doctoral Research Fellow with the Universidad Carlos III of Madrid in 2012 and a Visiting Researcher with the Politecnico de Milano in 2016. He has been with the Centro Universitario de la Defensa Zaragoza since 2013, where he is currently an Associate Professor. His research interests include wireless communications systems, specifically the analysis of resource allocation

problems for cellular, WLAN, and sensor networks.



**Carlos Donato** received the B.Sc. degree in telecommunication engineering and the M.Sc. degree in telematics engineering from the Universidad Carlos III de Madrid in 2014 and 2015, respectively, where he is currently pursuing the Ph.D. degree. He is also with the Institute IMDEA Networks. His research focuses on performance evaluation and centralized optimization of wireless networks.



**Pablo Serrano** (M'04–SM'15) received the Telecommunication Engineering and Ph.D. degrees from the Universidad Carlos III de Madrid (UC3M) in 2002 and 2006, respectively. He was a Visiting Researcher with the Computer Network Research Group, University of Massachusetts Amherst, in 2007, under a Jos Castillejo Grant, the Telefonica Research Center, Barcelona, in 2013, the School of Computer Science and Statistics, Trinity College Dublin in 2015, and the Universita degli Studi di Brescia in 2016. He has been with the Telematics

Department, UC3M, since 2002, where he is currently an Associate Professor. He has authored over 70 scientific papers in peer-reviewed international journal and conferences. He serves on the Editorial Board of the IEEE COMMUNICATIONS LETTERS.



Albert Banchs (M'04–SM'12) received the M.Sc. and Ph.D. degrees from the Polytechnic University of Catalonia in 1997 and 2002, respectively. He was with ICSI Berkeley in 1997, Telefonica I+D in 1998, and NEC Europe Ltd. from 1998 to 2003. He is currently a Full Professor with the University Carlos III of Madrid (UC3M) and the Deputy Director of the Institute IMDEA Networks. His research interests include the performance evaluation and algorithm design in wireless and wired networks. He is an Editor of the IEEE TRANSACTIONS ON WIRELESS

COMMUNICATIONS and the IEEE/ACM TRANSACTIONS ON NETWORKING.