

RAN energy efficiency and failure rate through ANN traffic predictions processing

Original

RAN energy efficiency and failure rate through ANN traffic predictions processing / Vallero, Greta; Renga, Daniela; Meo, Michela; Marsan, Marco Ajmone. - In: COMPUTER COMMUNICATIONS. - ISSN 0140-3664. - 183:(2022), pp. 51-63. [10.1016/j.comcom.2021.11.011]

Availability:

This version is available at: 11583/2942572 since: 2021-12-10T14:41:06Z

Publisher:

Elsevier

Published

DOI:10.1016/j.comcom.2021.11.011

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

Elsevier postprint/Author's Accepted Manuscript

© 2022. This manuscript version is made available under the CC-BY-NC-ND 4.0 license
<http://creativecommons.org/licenses/by-nc-nd/4.0/>. The final authenticated version is available online at:
<http://dx.doi.org/10.1016/j.comcom.2021.11.011>

(Article begins on next page)

RAN Energy Efficiency and Failure Rate Through ANN Traffic Predictions Processing

Greta Vallero^{a,1}, Daniela Renga^a, Michela Meo^a, Marco Ajmone Marsan^{a,b}

^a*Politecnico di Torino, Department of Electronics and Telecommunications, Torino, Italy*

^b*IMDEA Networks Institute, Madrid, Spain*

Abstract

In this paper, we focus on the application of ML tools to resource management in a portion of a Radio Access Network (RAN) and, in particular, to Base Station (BS) activation and deactivation, aiming at reducing energy consumption while providing enough capacity to satisfy the variable traffic demand generated by end users. In order to properly decide on BS (de)activation, traffic predictions are needed, and Artificial Neural Networks (ANN) are used for this purpose. Since critical BS (de)activation decisions are not taken in proximity of minima and maxima of the traffic patterns, high accuracy in the traffic estimation is not required at those times, but only close to the times when a decision is taken. This calls for careful processing of the ANN traffic predictions to increase the probability of correct decision. Numerical performance results in terms of energy saving and traffic lost due to incorrect BS deactivations are obtained by simulating algorithms for traffic predictions processing, using real traffic as input. Results suggest that good performance trade-offs can be achieved even in presence of non-negligible traffic prediction errors, if these forecasts are properly processed. The impact of forecast processing for dynamic resource allocation on the BS failure rate is also investigated. Results reveal that conservative approaches better prevent BSs from hardware failure. Nevertheless, the deployment of newer devices, designed for fast dynamic networks,

*Corresponding author

Email address: greta.vallero@polito.it (Greta Vallero)

allows the adoption of approaches which frequently activate and deactivate BSs, thus achieving higher energy saving.

Keywords: Radio access network; base station; energy efficiency; traffic prediction; neural network; base station lifetime

1. Introduction

The growth of computational power, the availability of data, the improvement of learning algorithms are the boosts behind the pervasiveness of new Artificial Intelligence (AI)- and Machine Learning (ML)-based mechanisms to respond to the new challenges of today's networks, which are often too complex to be properly understood, modelled, and managed with traditional approaches. This is the case of decision making for network management and configuration in presence of a huge set of parameters and fast changing scenarios, but also of catching the effect of complex interactions among multitudes of heterogeneous users and network elements (such as macro and small cells in heterogeneous networks), as well as understanding the hidden correlations among systems.

While these approaches are attractive because they are, by their nature, suited to handle problems with very large state spaces and complexity, in practice, effectively exploiting the potential of ML technologies is not easy. A deep domain knowledge is needed, as well as a careful processing of the outputs of ML tools. In this paper, we experiment this in Radio Access Network (RAN) management. We consider a portion of a RAN in which Base Stations (BSs) of macro and small cells can be activated and deactivated based on traffic load, so as to reduce energy consumption while guaranteeing that enough capacity is provided to satisfy the demand. BS activation and deactivation decisions are taken based on traffic predictions that are performed through Artificial Neural Networks (ANN). In order to properly operate the network so as not to deteriorate Quality of Service (QoS), the outputs of the ANNs have to undergo a number of processing steps that, driven by a deep domain knowledge, are carefully tailored for the scope.

Solutions for RAN management and resource on demand provisioning have been formulated in several contexts and with a multitude of different objectives: the trade-off between the opposite needs to on the one hand reduce energy consumption, and on the other provide QoS, is a timely objective, motivated by concerns on sustainability, climate change, and network operational cost increase. The deployment of BS management mechanisms, in its turn, is easier today due to the flexibility of new network architectures and it is effective for energy consumption reduction due to the typical traffic demand profiles at the edge of the network. Traffic demand profiles are characterised by (often short) peaks, followed by (often long - especially during night) valleys, and this makes the installed RAN equipment under-utilised for long periods of time. During these under-utilisation periods, some of the RAN equipment can be put in low power consuming sleep modes. Moreover, in some areas and in some periods (typically right after technology upgrades), the RAN capacity is over-provisioned even with respect to traffic peaks, and this makes BS management even more attractive for energy saving purposes.

As mentioned above, BS activation and deactivation decisions are taken based on traffic predictions. When the traffic is predicted to be small enough, some small cell BSs can be deactivated, and traffic is carried by the small and macro BSs that remain active. Conversely, when traffic grows and additional capacity is needed, some BSs in sleep mode are re-activated. Operating BSs, by activating and deactivating them, has an impact on BS failure rate: on the one side, switching is harmful to BS failure rate; on the other side, the time spent in sleep mode saves the BS from deterioration. The balance between these two phenomena depends on the hardware (HW) components of the BS, and on the RAN management strategy. The deterioration of the failure rate of the BS directly impacts its maintenance and its operation and maintenance cost for the operator. According to [1], this cost accounts for 3 k€ and 1 k€, per year, for each macro and micro cell BS, respectively. Up to 4 billions of BSs were counted worldwide back in 2012, as described in [2], and this number is bound to remarkably increase because of the RAN densification, planned with the 5G

RAN deployment [3, 4]. Thus, in addition to energy saving and QoS, these RAN management strategies have to be designed taking into consideration the impact on the BS failure rate, in order to avoid the explosion of the RAN Operational (OPEX) cost.

In [5], the performance of this BS management strategy was tested using several different ANNs for traffic predictions, performed over an hourly time scale. The results showed a limited sensitivity to the type of ANN. Indeed, critical BS (de)activation decisions are taken in correspondence of specific traffic values, and high accuracy in the estimations is not required in general, but only close to the times when decisions are taken. Hence, to significantly improve performance, traffic predictions need to be carefully processed and the overall pattern understood. In [6], we show that performing traffic predictions through ANN, over a shorter time scale, and their processing is fundamental for improving performance. In this paper, we introduce the BS failure rate as a new variable in the design space of RAN management. Results reveal that energy saving strategies based on conservative processing of the traffic demand forecasts reach significant energy consumption reduction, preserving QoS, as well as the BSs failure rate, also in case the BS HW has not been designed for dynamic switching. When the processing of the traffic demand predictions results in a dynamic BSs activation and deactivation, energy saving is further slightly improved, at the expense of a small loss in QoS and in the BS failure rate, suggesting that these approaches are suitable only in case the BS HW design is optimised for BS switching.

The paper is organised as follows. After the related work discussed in Section 2, in Section 3 we present the scenario and the methodology of our study. The proposed approach models each BS as defined in Section 4 and is based on traffic predictions, obtained with the tools that are presented in Section 5, and on the prediction processing algorithms that are reported in Section 6. After presenting performance indicators in Section 7, results are discussed in Section 8. Section 9 summarises our findings and Section 10 concludes our work.

2. Related work

In the green networking literature, many RAN management solutions have been proposed, based on Resource on Demand (RoD) approaches. An overview of the RoD strategies to dynamically adapt the set of active radio resources to the current traffic demand is presented in [7, 8, 9]. With the purpose of reducing energy demand and limiting the RAN operational cost, the authors of [10] exploit RoD strategies to adapt energy consumption to the actual traffic load. In [11], a framework is proposed to efficiently allocate spectrum resources to users, switching off unneeded BSs, in order to minimise power consumption. A time-varied probabilistic ON/OFF switching algorithm for cellular networks is presented in [12]. In [13, 14], RoD strategies are applied in a green mobile access network, with the objective of improving the interaction with the smart grid in a demand-response scenario, thus reducing the electricity bill and providing ancillary services. Recently, the effects of BS switching on the lifetime of the BSs have been investigated. The authors of [15] showed that putting a BS in sleep mode, besides reducing its energy consumption, increases its lifetime, since the BS operating temperature is reduced. This reduction depends on its HW components, i.e., on the materials used to build the device, and on the time spent in sleep mode. However, the same paper highlights also that power states transitions, which imply transition in the HW operating temperature, negatively affect the BS lifetime. For this reason, the works presented in [16, 17], formalise the optimal BS switching to maximise the RAN lifetime; the problem is solved through an heuristic, which allows to save up to 40% of energy during night, without decreasing the network survival duration. The effects of network device switching are also analysed in optical backbone networks, in [18, 19]. Many of these works, which focus on BSs switching, aim at dynamically allocating resources, under the assumption that the future traffic demand is exactly known. This means that predictions of the amount of traffic demand are necessary in order to make the proposed approaches viable. This aspect is very critical, since errors in the traffic estimation can significantly affect the perfor-

mance of these strategies. If the traffic demand is overestimated, waste of energy occurs; in case of traffic underestimation, incorrect BS deactivations may deteriorate QoS. To overcome the issue related to QoS deterioration due to traffic underestimation, [20] uses deep learning neural network based predictions, employing a customised loss function, to predict the needed network capacity. In particular, in case of underestimation, such function gives higher penalty than when the needed capacity is overestimated.

In the recent literature, many works focus on traffic estimation. In [21], an Auto-Regressive Integrated Moving Average (ARIMA) is used for the prediction of mobile data traffic and a Seasonal ARIMA (SA) model is used in [22]. These works demonstrate that these two methods provide high accuracy, but require slow training and forecasting, which make them impractical for on-line forecasting. The work presented in [23], uses Markovian models, while [22], [24], [25], [26], [27], [28], [29] employ ML approaches. According to [22] and [24], ANNs provide promising results in forecasting the hourly amount of traffic in TCP/IP networks. In [25], very good performance is reached in the forecast of the mobile traffic of an LTE BS, using a Recurrent Neural Network (RNN) and 1 ms resolution data. High accuracy in traffic predictions is achieved with the same approach, in [26]. An hybrid scheme, structured in an ANN and a RNN is discussed in [27]. Moreover, [28] and [29] predict traffic demand with Least Squares Support Vector Machines and a Linear Regression based approach, respectively.

Differently from the previous literature, in this paper, traffic predictions are based on different ANN structures, and RoD strategies are applied based on traffic predictions, after processing, with the objective of reducing the RAN energy consumption without (or with minimal) QoS deterioration. This means that, in our work, high accuracy traffic forecasts are not the primary objective. Predictions are a necessary tool to achieve the proper resource allocation, which allows energy saving without compromising QoS. In particular, since traffic maxima and minima do not trigger BS switch-on/off, we are not interested in careful estimations of such values. Rather, we try to carefully identify the instants when to activate or deactivate BSs. To achieve this goal, we apply processing

techniques on the traffic forecasts, and we show that a careful selection of these techniques can be more impactful on the achieved RAN energy efficiency than
150 the careful selection of the traffic predictor.

3. Scenario

A portion of an heterogeneous LTE RAN is considered, comprising one macro cell BS, and a few small cell BSs, whose coverage overlaps with the macro cell. This is a typical scenario considered for 5G and beyond RAN architectures, that
155 leverage small cell BSs exploiting high frequency bands (tipycally millimeter wave). Small cell BSs are deployed to provide additional capacity during high traffic demand periods. A centralised Management and Orchestration System (MANO) decides the activation of resources (i.e., small cell BSs), according to predictions of the future traffic demand. These predictions are performed
160 on a temporal horizon of 15 minutes (which is the time granularity chosen by the operator whose data we used in this work, and is thus taken as the time slot). Small cell BSs can be switched on and off to reduce the RAN energy consumption, with attention to QoS. This means that, when not all the capacity is needed to satisfy the predicted traffic demand, some small cell
165 BSs are put in sleep mode. On the contrary, all BSs are activated in those periods when all the capacity is required for the traffic demand satisfaction. The activation/deactivation of a BS cannot occur at intervals shorter than one hour, to avoid too frequent switches.

Traffic predictions are performed through ANNs, and are processed before reaching a decision about BS activation and deactivation. The BS management works
170 in two phases.

1. **Training phase.** This phase is performed only once, as a preliminary step of our online management system. In this step, the ML algorithm used to predict the traffic demand is trained using historical data.
- 175 2. **Run-time phase.** The traffic is predicted using the previously trained ANN, and BS activations or deactivations are decided. During this phase,

the following two steps are performed at every time slot, i.e., every 15 minutes:

- (a) **Prediction.** The traffic demand is forecast for the following 4 time slots.
- (b) **Prediction processing and decision.** The four predictions are processed by the MANO, which decides which small cell BS must be active in the next hour.

180

4. Modelling the BS

The input power, in watt, required for the operation of a BS at time slot t , denoted as $P_{in}(t)$, is derived according to the linear model proposed in [30], which has a fixed component, corresponding to the amount of power needed to keep the BS active, and a load-dependent component. It can be expressed as follows:

$$P_{in}(t) = N_{trx} \cdot [P_0 + \Delta_p P_{max} \rho(t)], \quad 0 \leq \rho \leq 1 \quad (1)$$

where N_{trx} is the number of transceivers, P_0 represents the power consumption in watt when the radio frequency output power is null, Δ_p is the slope of the load dependent power consumption, $\rho(t)$ is the traffic load at time slot t , defined as the traffic carried by the BS in b/s and the BS capacity in b/s. P_{max} is the maximum radio frequency output power in watt at maximum load. Table 1 summarises the value of the parameters for macro and small cell BSs [30]. The consumption of the BS in sleep mode is considered negligible.

We adopt the model of the BS failure rate presented in [17] and [31]. The model treats the BS as a whole, i.e. as a single entity, rather than using a model for each single component of the BS, even if this assumption results in a less detailed model, which does not specify the dependencies among the BS components. The failure rate of a BS, and in general, of an arbitrary device, is given by the Arrhenius law [32] and is strictly dependent on the operating temperature. As derived in [17, 31, 32], putting a BS in sleep mode positively impacts its lifetime and failure rate, since its operating temperature is reduced.

Nevertheless, the power state change negatively affects the failure rate, increasing it, since, as explained in [33], the metal is sensitive to temperature variations and, in particular, to state cycling. As a result, the failure rate of a BS, and, in general, of an arbitrary device, expressed in failure/h, is given by two contributions: i) the failure rate in active and sleep modes, weighted by the time the device is in those states, and ii) the frequency of the device's operational state changes. In [16, 17, 31], a failure is every HW event which causes interruption of the BS operation, both when its repair needs some HW substitution or not. According to [16, 17], the failure rate of a BS b can be expressed as:

$$\gamma_b = (1 - \tau_{sleep})\gamma_{on} + \tau_{sleep}\gamma_{sleep} + \frac{f_{tr}}{N_F} \quad (2)$$

where τ_{sleep} is the fraction of time the device spends in sleep mode, γ_{on} and γ_{sleep} , in failure/h, are the failure rates when the BS is active and in sleep mode, respectively, computed with the Arrhenius law [32]. The parameter f_{tr} , in cycle/h, is the frequency of the sleep mode cycle and N_F , in cycle/failure, is the number of cycles supported by the device before a failure occurs. Usually, to measure the impact of the device switching on its lifetime, the Accelerator Factor (AF) is estimated. This indicator provides the mean lifetime increase/decrease with respect to the always on condition, as the ratio between the resulting failure rate and the failure rate of the always on scenario. A value of AF larger than 1 means that the failure rate increases, while a value smaller than 1 indicates that the failure rate decreases. Similar to (2), the resulting AF is given by two contributions: the time spent in sleep mode, which decreases the BS failure rate (and AF), and the frequency of the operating state changes, which deteriorates the BS failure rate (and increases AF). For each BS b , the AF can be computed over a period of duration θ , as follows:

$$AF_{b,\theta} = \frac{\gamma_b}{\gamma_b^{on}} = 1 - \underbrace{(1 - AF_{sleep})\tau_{sleep}}_{\text{Lifetime Increase}} + \underbrace{\chi f_{tr}}_{\text{Lifetime Decrease}} \quad (3)$$

185 where AF_{sleep} is the AF, computed assuming that the device is always kept in sleep mode. According to [32], it is always lower than 1, otherwise putting the device in sleep mode would mean increasing the failure rate of the device.

Table 1: Values of the parameters of the consumption model for macro and small cell BSs.

BS type	N_{trx}	P_{max} (W)	P_0 (W)	Δ_p
Macro	6	20	84	2.8
Small	2	6.3	56	2.6

Then, τ_{sleep} is the fraction of time the BS has spent in sleep mode in the period of duration t . The parameter f_{tr} , in cycle/h, is the frequency of the switching cycle which is measured over t and χ , in h/cycle, is defined as $\frac{1}{\gamma_b^{on} N_F}$ and acts as weight of the frequency f_{tr} . As a result, the drop of the BS failure rate is achieved when $\chi f_{tr} < (1 - AF_{sleep})\tau_{sleep}$. Notice that the parameters χ and AF_{sleep} depend on the HW component used to build the BS, while τ_{sleep} and f_{tr} depend on the switching strategy.

5. Traffic predictions

In this section, we present the traffic prediction tools used during the BS run-time as a preliminary step to the BS management decision. An ANN-based approach is used for this purpose. This is because, as mentioned, in the literature, the potentiality of ANN has been widely demonstrated [22] [24]. Moreover, [5] shows that the performance of the BS activation/deactivation are not significantly affected by the ML approach used for the traffic predictions. Thus, the usage of a simple ANN represents a good trade-off between performance and complexity.

5.1. Input Data

Data provided by a large Italian mobile network operator are used in this study. They report the traffic demand volume, in bit, of 1420 BSs located in the city of Milan (Italy) and in a wide area around it, for two months in 2015, with granularity of 15 minutes. This time periods includes typical weeks, as well as Easter week, when a brief vacation occurs. During typical weeks, people

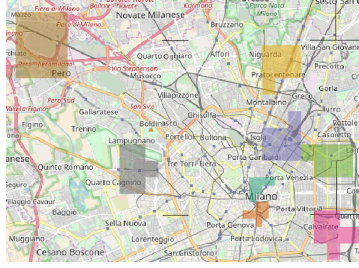


Figure 1: Considered traffic areas: train station (purple), Rho Fiere (brown), Duomo di Milano (orange), Politecnico di Milano (light green), San Siro (grey), a business area (dark green), a residential area (yellow), an industrial area (magenta).

210 follow their usual working and activities routine. As a result, the considered period provides a good representation of the traffic demand volume dynamics. The traffic traces are normalised; hence, the peak of each traffic pattern is equal to the maximum capacity of each BS. Note that this is a pessimistic assumption with respect to energy saving possibilities, since the capacity of the network is usually overdimensioned. For our work, eight portions of the city are selected, which are shown in Fig. 1. These areas were selected as samples of quite different scenarios, and, hence, traffic patterns. All together, the selected areas are representative of the various zones that coexist in a urban environment. The train station area (purple square in Fig. 1) is characterised by intense activity levels, especially at the beginning and at the end of the working hours. The Rho Fiere district (brown in the figure) is an area that hosts big events, fairs and exhibitions, that last for a few days. The Duomo di Milano area (orange square) is a touristic area, with high activity during several hours of the day. The Politecnico di Milano area (light green) hosts a large campus with many students. The San Siro neighbourhood includes a large soccer stadium (grey), and the activity here is quite bursty and variable depending on the scheduled matches and concerts¹. A part of a business neighbourhood (dark green) and some residential streets (yellow) are also considered: the traffic in these areas

220

225

¹In the considered time period, soccer matches were held over weekends.

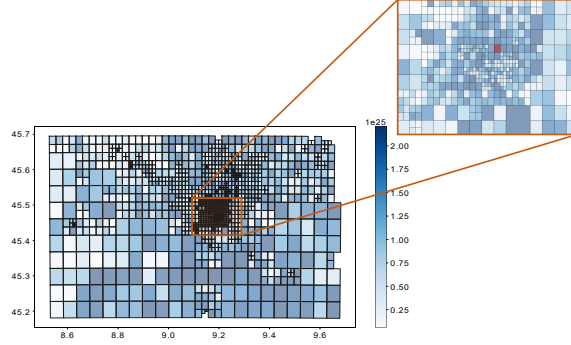


Figure 2: Spatial cross-correlation among cells with lag= 15 minutes.

follows the typical behaviour of people in their daily life. In the business area,
 230 traffic peaks are observed during the central hours of the day, whereas in the
 residential area a traffic rise is observed in the evening. Finally, the industrial
 zone (magenta) is a particular case of a business area. In each of these portions
 of the RAN, we assume that one macro BS and 6 small cell BSs are present, so
 that the service area is covered by one macro cell which overlaps with 6 small
 235 cells. To do this, for each area, we selected 7 traffic patterns recorded in that
 area. The trace which presents the highest traffic demand is chosen as the macro
 cell BS, while the remaining six as micro cell BSs.

5.2. Selection of the ANN input features

In order to predict traffic demand, the ANN must be fed with carefully
 240 selected input features. The investigation of the best choice for the ANN input
 features was made accounting for the temporal and spatial correlations of traffic.
 In particular, we exploit the traffic temporal periodicity (which we observed to
 be present in most traffic patterns) due to the periodicity of human activities,
 and we investigated the possibility of also using the spatial correlation which is
 245 expected to be present among adjacent cells. In Fig. 2, the cross-correlation
 obtained between the traffic at one BS in the city centre, indicated in red, and
 all others is plotted, choosing as time lag one time slot, i.e., 15 minutes. We

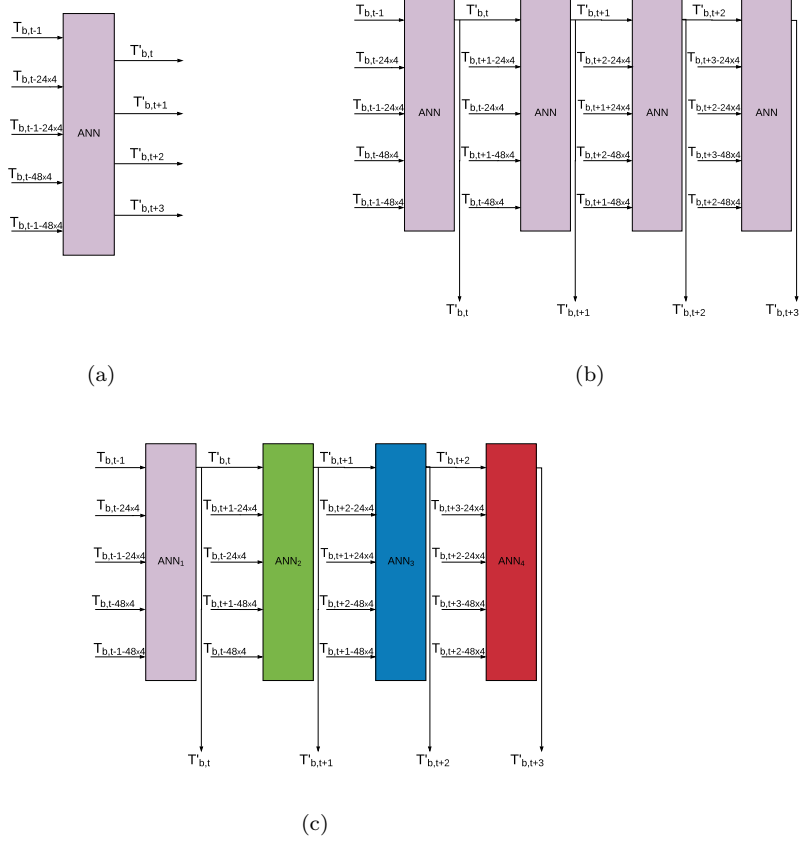


Figure 3: Scheme of the three proposed prediction techniques: (a) 1 ANN-4 outputs, (b) 1 ANN-1 output, (c) 4 ANNs-1 output

can see that correlation only mildly depends on the spatial closeness to the considered BS (darker colours correspond to higher correlation values). Indeed, high correlation values are present even among cells that are very far from each other. For this reason, in this paper we focus only on input features based on the temporal periodicity of traffic patterns.

Let us define by $T_{b,i}$ the traffic demand at BS b and time slot i . For simplicity of notation, in what follows we drop the index b if there is no ambiguity. At the beginning of each time slot t , the traffic demand at time slot t (the time slot that is just beginning), $t+1$ (the following time slot), $t+2$ and $t+3$ are

predicted; predictions are denoted by $T'_t, T'_{t+1}, T'_{t+2}, T'_{t+3}$, respectively.

The prediction tool receives as inputs:

- T_{t-1} : the traffic at the time slot just past, i.e., $t - 1$;
- 260 • $T_{t-(24 \cdot 4)}$: the traffic one day before the current time slot (the factor 4 comes from our time slots being 15 minutes long);
- $T_{t-1-(24 \cdot 4)}$: the traffic one day before the time slot just past;
- $T_{t-(48 \cdot 4)}$: the traffic two days before the current time slot;
- $T_{t-1-(48 \cdot 4)}$: the traffic two days before the time slot just past.

265 5.3. Traffic Forecast Approach

Different ANN-based prediction approaches are tested.

5.3.1. 1 ANN-4 outputs

One ANN for each BS is used. At time t , the ANN outputs the traffic demand samples at time slots $t, t+1, t+2$ and $t+3$ (see Fig. 3a).

270 5.3.2. 1 ANN-1 output

One ANN for each BS is used. The ANN is trained to predict the traffic demand at the current time slot, e.g. at time t , and it is used in cascade to predict also the three future traffic samples, e.g. at time $t + 1, t + 2, t + 3$. This means that the ANN produces the prediction of the traffic demand at time t , namely 275 T'_t , using in input the traffic at previous time slot, T_{t-1} , as well as the traffic of previous days. Once T'_t is computed, for predicting the traffic at time $t + 1$, the same ANN is used but it receives as input the predicted traffic T'_t instead of T_t that is unknown. Similarly, for the prediction of traffic at times $t + 2$ and $t + 3$, predictions are used instead of traffic samples for the unknown values of 280 the input. The logical schema is reported in Fig. 3b.

5.3.3. 4 ANNs-1 output

Four ANNs are used for each BS. Each ANN is dedicated to the prediction of the traffic demand at a given time lag. This means that the 4 future traffic samples are separately predicted, using 4 different ANNs, but the inputs are as
285 in the previous case: predictions are used instead of missing samples whenever needed. The schema is reported in Fig. 3c.

As in [5], each ANN mentioned above is structured in 3 layers: the input layer which has 8 nodes, one hidden layer with 17 nodes, and the output layer with one node, if *1 ANN-1 output* and *4 ANNs-1 output* are employed, or 4 nodes in
290 case of *1 ANN-4 outputs* usage. The number of layers, as well the number of nodes for each layer are among the hyper-parameters which need to be selected. These have been chosen in order to achieve a good trade off between the accuracy and the time needed to train the network. Each ANN is trained minimising the Mean Squared Error (MSE) over the data of the first 47 days of the considered
295 time period.

6. Processing traffic predictions

After the ANN has generated traffic predictions, they must be processed by the MANO to decide about micro cell BS (de)activation, with the objective to save energy, without compromising QoS. In this section, we propose strategies
300 for processing the predictions and deciding microcell activation and deactivation.

6.1. Resource Allocation

Different algorithms can be used to combine traffic predictions in a BS management strategy, based on the approach in [10], which states that a micro cell
305 BS is switched off if its traffic demand is lower than a threshold ρ^* , provided that such amount of traffic can be carried by the macro cell BS. The threshold depends on the energy consumption per carried bit: when the traffic is below ρ^* , the energy needed to carry a unit of traffic in the micro cell is larger than in

the macro, so that it is more convenient to switch off the small cell BS, if this is possible in terms of total capacity. As demonstrated in [10], the optimal value of the threshold is 37% of the maximum load of the BS.

6.1.1. *Max2Max*

In this case, resources can be allocated only at the beginning of each hour: at 00:00, 01:00, etc. At the beginning of each hour, $T'_{b,t}$, $T'_{b,t+1}$, $T'_{b,t+2}$, $T'_{b,t+3}$, the 4 traffic demands corresponding to that hour are predicted for each micro cell BS b , as well as for the macro cell B ; predictions in the macro are denoted by $T'_{B,t}$, $T'_{B,t+1}$, $T'_{B,t+2}$, $T'_{B,t+3}$. Among these 4 samples, the maximum, M'_b , for each micro cell BS b and the maximum, M'_B , for the macro cell are computed:

$$M'_b = \max(T'_{b,t}, T'_{b,t+1}, T'_{b,t+2}, T'_{b,t+3}) \quad (4)$$

$$M'_B = \max(T'_{B,t}, T'_{B,t+1}, T'_{B,t+2}, T'_{B,t+3}) \quad (5)$$

A micro cell BS b is switched off if M'_b is lower than the threshold, and its traffic can be carried by the macro, given that the macro is expected to be carrying an amount of traffic M'_B :

$$\text{if } (M'_b < \rho) \wedge (M'_B + M'_b < C) \rightarrow \text{switch off } b \quad (6)$$

where C is the capacity of the macro cell B . Basically, the decision is taken based on the maximum of the predicted traffic samples. As the decision is taken, M'_B is updated accordingly, to account for the traffic load that will be transferred from the considered BS.

6.1.2. *Max2Max Continuous*

This strategy is very similar to *Max2Max*, but, in this case, it is applied at the beginning of each 15 minute time slot and not only at the beginning of an hour, as in the previous case. The decision to switch off a cell for 4 consecutive time slots (1 hour) can be taken in any time slot.

6.1.3. *I2I*

When this strategy is used, the switch on/off is possible only at the beginning of each hour. Given the four predicted traffic demands belonging to the

considered hour, for each micro cell BS b and for the macro cell BS B , a micro cell BS b is switched off when, for every slot $t + i$ with $i = 0, 1, 2, 3$, the estimated traffic $T'_{b,t+i}$ is lower than the threshold ρ^* and there is enough available capacity on the macro BS:

$$\begin{aligned} \text{if } \forall i = 0, \dots, 3 \quad (T'_{b,t+i} < \rho) \wedge (T'_{b,t+i} + T'_{B,t+i} < C) \\ \rightarrow \text{switch off } b \end{aligned} \quad (7)$$

335 In this case, the decision to switch off is taken if the requested conditions are verified slot by slot.

6.1.4. *I2I Continuous*

When this strategy is used, *I2I* is applied at the beginning of each time slot. As in *Max2Max Continuous*, each micro cell BS remains active or in sleep mode
340 for at least 1 hour (4 consecutive time slots), but a change of state can happen in any time slot.

6.1.5. *I2I Flexible*

This is a further variation of *I2I Continuous*. As before, at the beginning of each time slot, *I2I* is applied. Nevertheless, when a micro cell BS has been
345 put in sleep mode for at least one hour, it remains sleeping if the necessary conditions are verified for one more time slot. This means that when we are at time t , given that the micro cell BS has been deactivated since at least $t-4$, that micro cell BS remains in sleep mode, if $T'_{b,t}$ is lower than ρ^* , provided that $T'_{b,t}$ can be carried by the macro BS during the t -th time interval.

350 Each of the considered micro cell BSs is analysed for its possible deactivation as described above, starting from the least loaded to the most loaded, in the following hour. Given that the load of a micro BS is lower than ρ^* , its energy consumption per bit is larger, if its load is smaller. Thus, giving larger priority to micro BSs in the deactivation procedure leads to minimum network energy
355 consumption [10]. The load during the following hour on each micro BS is given by summing the traffic demand during the 4 time slots belonging to that hour.

6.2. Descending front detection

The presence of noise in traffic patterns may result in incorrect deactivation of the small cell BSs, thus deteriorating QoS. For this reason, the concept of *Descending Front Detection* (DFD) is introduced in the processing of traffic predictions. In particular, the switching from active to sleep mode of a micro cell BS is permitted only if a descending front is detected: if an active micro cell BS is detected to be in a descending phase, the necessary conditions for the micro cell switchoff are checked. Because of the noise inherent in traffic patterns, a negative first derivative is not a sufficiently good indicator of a descending front. Therefore, a moving average filter is used for this purpose. It smooths data by replacing each traffic sample with the average of the neighbouring samples. This operation practically acts as a low-pass filter on traffic patterns. In our case, a triangular smoothing is applied twice. In particular, at time t , the following expression is computed, for $z = t-4, t-5, t-6$:

$$S'_{b,z} = \frac{1}{81} \sum_{j=-2}^2 (3 - |j|) \sum_{i=-2}^2 (3 - |i|) T_{b,z+j+i} \quad (8)$$

where $T_{b,z+j+i}$ is the real traffic demand on BS b at time $z+j+i$. However, notice that for $z = t - 4$ and for $j = 2$ and $i = 2$, $T_{b,z+j+i}$, is $T_{b,t}$, which is not
360 known. Thus, its prediction, $T'_{b,t}$ is used in this case. The maximum z is chosen equal to $t-4$, in order to avoid using other predicted samples.

If $S'_{b,t-4} < S'_{b,t-5} < S'_{b,t-6}$, we conclude that a descending front is detected. If this is the case, the necessary micro cell BS switchoff conditions are checked. If they are verified, as described in section 6.1, the considered micro cell BS can
365 be deactivated.

7. Key Performance Indicators

7.1. Average Relative Error

The ARE (Average Relative Error) measures the average relative error between the real and predicted traffic samples. It is computed as:

$$ARE = \frac{1}{N_{BS}} \sum_{b=1}^{N_{BS}} RE_b \quad (9)$$

where N_{BS} is the number of the considered BSs and RE_b is the RE (Relative Error) on BS b , derived as:

$$RE_b = \frac{1}{H} \sum_{t=1}^H \frac{|T_{b,t} - T'_{b,t}|}{T_{b,t}} \quad (10)$$

where $T_{b,t}$ is the real traffic demand at time t on BS b , $T'_{b,t}$ is the forecast traffic demand at time t on BS b , H is the duration of the testing period, in number of time slots.

7.1.1. Energy Consumption Reduction

When the resource allocation strategies presented in Section 6.1 are used, in each time slot some BSs are active and consume energy, while some others may be in sleep mode and thus consume no (or very little) energy. The energy consumption of each BS is given by (1). In order to measure the effectiveness of these strategies, the energy saving is computed. It is calculated with respect to the *always ON* scenario: this is the case in which all BSs are always active regardless the amount of traffic demand. It is computed as follows:

$$EC_{Red} = 100 \cdot \frac{EC_{on} - EC}{EC_{on}} \quad (11)$$

where EC_{on} is the energy consumption in the *always ON* scenario; EC is the energy consumption with the considered strategy, computed as $\sum_{t=0}^H \sum_{b=1}^{N_{BS}} EC_{t,b}$, where $EC_{t,b}$ is computed as in (1) and N_{BS} is the number of the considered BSs.

7.2. Lost Traffic

In the *always ON* scenario, each BS is always active and able to carry its traffic demand. In case resources are dynamically allocated according to the strategies described above, the situation is different. The *Lost Traffic* is defined as the percentage of the traffic demand that cannot be carried by the network, accounting for the fact that in each time slot some BSs are active and can handle their traffic demand, while some others may be off and thus cannot provide any service. Let us define the traffic that overflows from the micro cell BS b to the macro cell as:

$$O_{b,t} = \begin{cases} T_{b,t} & \text{if } b \text{ is in sleep mode} \\ 0 & \text{if } b \text{ is active} \end{cases} \quad (12)$$

the lost traffic is given by:

$$L = \frac{\sum_{t=1}^H \max(0, T_{B,t} + \sum_{b=1}^{N_{BS}} O_{b,t} - C)}{\sum_{t=1}^H (T_{B,t} + \sum_{b=1}^{N_{BS}} T_{b,t})} \cdot 100 \quad (13)$$

where C is the capacity of a macro BS. The lost traffic is the percentage of traffic that cannot be carried by the macro cell BS when traffic overflows from deactivated small cell BSs.

7.3. Accelerator Factor for each micro cell BS

For each BS b , we measure $AF_{t,b}$ and AF_b which are, respectively, the AF of that BS b , measured at time t and at the end of the considered operating period, i.e. in steady state.

7.4. Accelerator Factor

For each considered portion of network, we measure AF which is the average AF of that area at the end of the considered operating period, i.e. in steady state:

$$AF = \frac{1}{BS} \sum_{b=1}^{BS} AF_b \quad (14)$$

where BS is the number of micro cell BS in the considered portion of RAN and AF_b is the AF, for each BS b in the considered area, computed as in (3).

Table 2: Average relative error, ARE, with the different approaches at different time lags.

ARE	1 ANN- 4 outputs	1 ANN- 1 output	4 ANNs- 1 output	4 ANNs- 1 output (spatial)
ARE_t	0.33	0.33	0.33	0.37
ARE_{t+1}	0.43	0.44	0.43	0.47
ARE_{t+2}	0.52	0.52	0.48	0.52
ARE_{t+3}	0.61	0.57	0.52	0.54

8. Performance evaluation

In this section, we discuss numerical results obtained by experimenting the different prediction, processing and decision algorithms presented in the previous sections on the considered RAN portions. Out of the 61 days for which we have
390 real traffic data, the first 47 are used for the ANN training phase, while the remaining 14 days are used for the run-time phase.

8.1. Choice of the ANN

As a first step, we analyse the effectiveness of the different ANN configura-
395 tions for traffic predictions, using the previously defined ARE (average relative error) as a performance metric. The results provided by the considered ANN configurations, namely *1 ANN-4 outputs*, *1 ANN-1 output* and *4 ANNs-1 output*, for each time lag, are reported in Table 2, averaged over the eight considered geographical areas. Observe that numerical results confirm what is intuitively
400 expected, and was quantitatively shown in [25]: the error increases with the time horizon of the predictions. Moreover, typically, the *1 ANN-4 outputs* provides the largest ARE. This is because, when the other 2 approaches are used, the sample corresponding to the most recent traffic demand, even if only predicted, is provided as an input feature. In Fig. 4, the percentage of the reduction of
405 ARE gained with *1 ANN-1 output* and *4 ANNs-1 output*, with respect to *1*

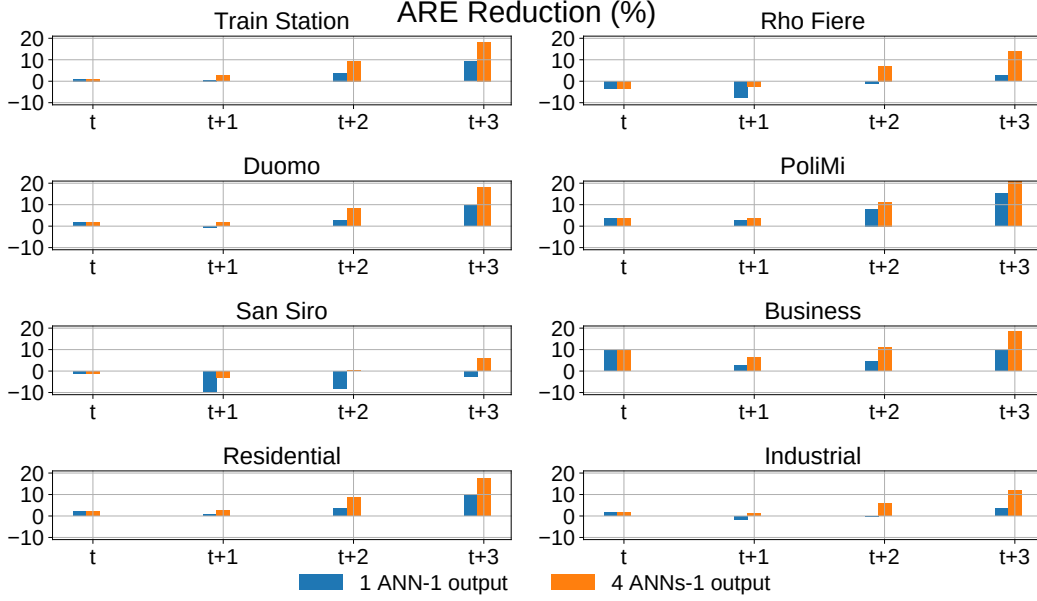


Figure 4: Percentage of the reduction of absolute relative error, ARE, obtained by $1\ ANN-1$ output or $4\ ANNs-1$ output with respect to the $1\ ANN-4$ outputs, for different time lags.

$ANN-4$ outputs, are shown in blue and orange, respectively. The reduction of the estimation error is the largest for $4\ ANNs-1$ output, especially when the time horizon of the predictions is longer. This is because this ANN configuration uses 4 ANNs: during the training phase, each ANN learns how to forecast the desired output, managing the error which affects the input traffic sample derived from a prediction. For these reasons, in the rest of this study we will use $4\ ANNs-1$ output for traffic forecast, unless otherwise specified.

In order to confirm the mild correlation among adjacent cells, we also report the ARE which is obtained when we provide to $1\ ANN-4$ outputs an additional input feature. In order to select this additional input feature, the cross-correlation between the traffic demand of the current BS and each of its adjacent ones, is performed. Then, the argmax function is computed, in order to select the BS bs_{MAX} and the time lag l_{MAX} which provide the largest value of cross-correlation. Thus, when we are predicting the traffic demand at time t , the

420 additional input feature is the traffic demand on bs_{MAX} at time $t-l_{MAX}$. Similarly, for the prediction of the traffic demand at $t+1$, $t+2$ and $t+3$, the traffic demand on bs_{MAX} at time $t+1-l_{MAX}$, $t+2-l_{MAX}$, $t+3-l_{MAX}$, respectively, are given as additional input feature. From Table 2, it is possible to notice that the presence of this feature deteriorates the precision of the forecast.

425 8.2. Dynamic resource allocation performance

We now investigate the performance of the resource allocation strategies presented in Section 6.1. Our solutions are compared against 3 benchmarks: (i) the *TNSM19* approach presented in [5], which allocates the resources of a RAN according to the hourly traffic predictions obtained using an ANN, with
 430 no processing of the ANN outputs; (ii) the *PIMRC18* approach: in this case the traffic is predicted using the LSTM (Long Short Term Memory) network proposed in [25] and resources are allocated based on *I2I*; (iii) the *15 min* approach, similar to the *TNSM19* case, but operated over 15 minutes time slots. With this approach, each small cell BS can be switched to/from sleep
 435 mode as soon as needed, with no constraint on the frequency of switching.

8.2.1. Effect of traffic pattern shape and load distribution

For each strategy and zone, Fig. 5a reports the energy consumption reduction computed with respect to the *always ON* scenario, and Fig. 5b reports the percentage of lost traffic. First, it is possible to confirm that, as expected, the
 440 percentage energy saving directly depends on the shape of the traffic pattern (peak/off-peak ratio, duration of peaks, ...), which is characteristic of the considered area. If the traffic demand is low for many hours, the BS management approach can be very effective: up to 40% of the energy consumed with respect to the always ON approach can be saved, as we see in the San Siro and Rho
 445 Fiere areas. When the traffic demand is larger than ρ^* for longer periods, the small cell BSs can be switched off for shorter periods, and a lower amount of energy is saved. This is the case of the PoliMi and Train Station areas, where the energy saving is lower than 15%. In Fig. 6a, the traffic demand during 5



Figure 5: Comparison of dynamic resources allocation strategy in the various areas: (a) Energy consumption reduction and (b) lost traffic.

days of the simulation of a micro cell BS in the Train Station and San Siro areas is plotted, in dark and light grey, respectively. The former presents a traffic volume usually larger than the threshold, indicated by the black horizontal line, while the latter almost always lower than ρ^* . As a consequence, their sleeping time ratio τ_{sleep} is very different, as can be noticed in Fig. 6b, where τ_{sleep} is plotted, for each dynamic resource allocation approach. In particular, τ_{sleep} is never larger than 0.25 and lower than 0.9 for the micro cell BS located in the Train Station and San Siro areas, respectively. Nevertheless, the frequency switching f_{tr} , reported in Fig. 6c, assumes very close values for both cases, since the micro cell BS located in the Train Station zone is usually ON and is deactivated only during the night, while the one placed in the San Siro area is usually in sleep mode and is active only during public events, when additional capacity is needed in order to satisfy the traffic demand. Fig. 6d reports the traffic demand during 5 days of simulation of two micro cell BSs, BS A and BS B, of the residential district. These two patterns, reported in dark and light grey in the figure, are very similar in shape and volume and are lower than the threshold for most of the time. However, their corresponding τ_{sleep} assume very different values. This is because micro cell BS A presents a lower traffic

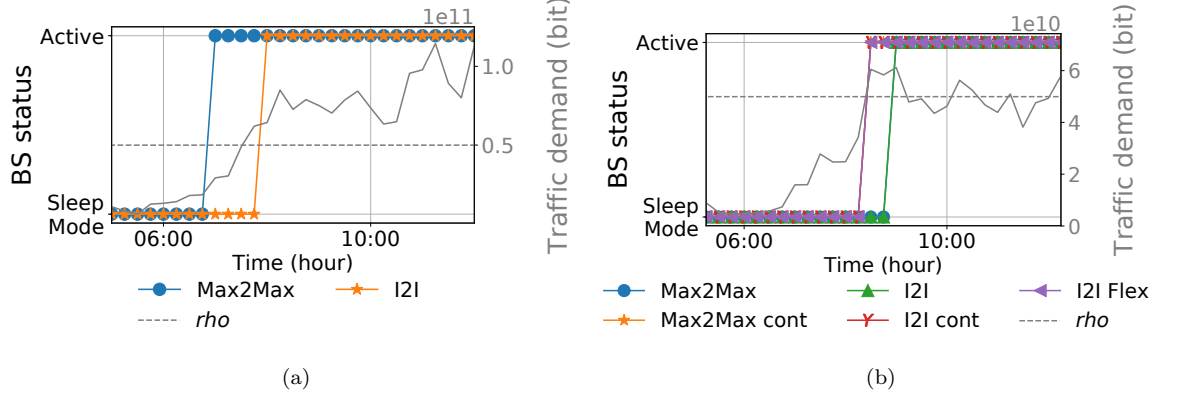


Figure 7: Comparison of dynamic resources allocation strategies in the various areas: (a) *Max2Max* and *I2I* and (b) *Cont* version.

more conservative in energy saving, but better preserve QoS, measured as the percentage of lost traffic. When *TNSM19* is used, resources are allocated under the unrealistic assumption that the traffic demand is uniformly distributed within a whole hour. For this reason, the lost traffic results higher than with the other approaches. With *PIMRC18*, up to 4% of traffic is lost. Even if it provides traffic predictions affected by lower *ARE* (0.29, 0.37, 0.42, 0.47, for the forecast at time t , $t+1$, $t+2$ and $t+3$, respectively), it usually generates more underestimated traffic samples that contribute to QoS deterioration.

The comparison with the *15 min* case is also interesting. The *15 min* case is based only on traffic predictions performed over a time horizon of 15 minutes for which the error is lowest (see table 2). Nonetheless, in this case no processing of the ANN outputs is performed; hence, despite the small error in predictions, the lost traffic is quite large. This is a clear indication of the importance of the processing of ANN outputs.

Let us now focus on the proposed approaches. The lost traffic is lower in the *Max2Max* case than in the *I2I* one, since its switching condition is stricter. Fig. 7a shows the status of a micro cell BS of the Train Station area in orange and blue, when *I2I* and *Max2Max* are used, respectively. The micro cell BS

500 traffic demand is reported in grey and ρ^* with the dashed grey curve. Even if these two approaches use the same prediction samples for resource allocation, *Max2Max* makes the BS active sooner than *I2I*. At 7.00 a.m., predictions of the traffic demand during the following hour are erroneously smaller than the threshold. Nevertheless, *Max2Max* switches BSs if the maximum, among the
505 traffic demand samples belonging to that hour, is smaller than ρ^* and can be carried by the macro cell, supposing that it is managing an amount of traffic which is the maximum traffic demand among the 4 traffic demand samples of that hour. Thus, the micro BS is activated, since its traffic demand cannot be carried by the macro BS because of the capacity constraint.

510 The use of the *cont* variation provides benefits in terms of QoS in both strategies, *Max2Max* and *I2I*. When *cont* is used, the effect of higher errors, which characterises traffic predictions over longer time lags, is further mitigated, so that a more accurate resource allocation can be performed. As a result, the achieved lost traffic is always less than 2%. Specifically, in the areas where the
515 resource allocation is more difficult due to the unpredictability of traffic demand, i.e., Rho Fiere and San Siro, 1.6% and 2% of the traffic is lost. In areas where patterns are more regular, values are always lower than 1.4%. Similar results are given by *I2I Flex*: the lost traffic is lower than the chosen benchmarks because the BS switching can react to the traffic demand every 15 minutes, provided
520 that the last switching has occurred since at least 1 hour. This can be observed in Fig. 7b, where each curve corresponds to the status of a micro cell BS of the Duomo area, obtained with each of the considered allocation approaches. Also its traffic demand (in grey) and ρ^* (grey dashed curve) are reported. The *cont* variation and *I2I Flex* react as soon as the traffic demand increases (at
525 8.30). When *I2I* and *Max2Max* are used, resources are allocated at 8.00 a.m., and the prediction with lag equal to 3 is used for that time slot. Because of the large error which affects this forecast, this sample results lower than the threshold, and the micro cell BS is not activated. Thus, from Figs. 7a and 7b, it is possible to notice that strategies behave similarly if the traffic demand
530 is far from the threshold. Indeed, in this case, the large error, which affects

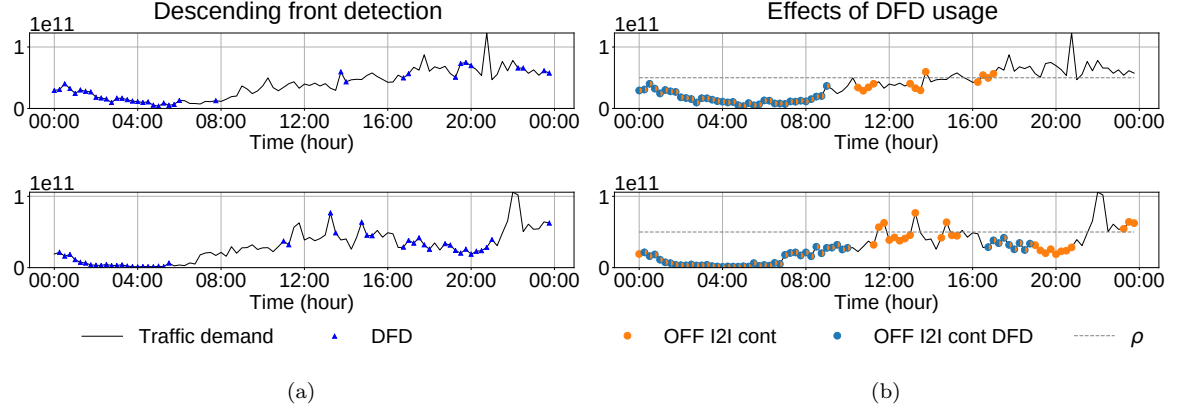


Figure 8: Impact of descending front detection for two areas: (a) detection of fronts and (b) switch off decisions with and without DFD.

typically deeper forecasts used by *I2I* and *Max2Max*, does not impact resource allocation, correctly detecting the value of the traffic demand with respect to the threshold. As soon as the traffic demand moves closer to the threshold, even if based on the same predictions, the resources allocation is different. In case of
535 *max*, conditions for the deactivation are stricter; with *cont* based approaches, more accurate predictions can be used. This results in more likely activation of micro BSs and, consequently, in lower lost traffic.

8.3. Impact of descending front detection

We now investigate the impact of the use of DFD in resource allocation.
540 When DFD is used, the deactivation of a micro cell BS is possible only if a descending front is detected, according to the conditions described in Sec. 6.1. In Fig. 8a, blue triangles mark the detection of a descending front during one day of the run-time phase of 2 micro cell BSs, belonging to the PoliMi and Rho Fiere areas. As can be observed, descending fronts are mostly correctly
545 identified. Since the current predicted traffic demand has lower impact on DFD than past samples, see equation (8), it is possible that DFD is activated after a local minimum.

Fig. 9 reports the energy consumption reduction, in bars, and the lost traffic,

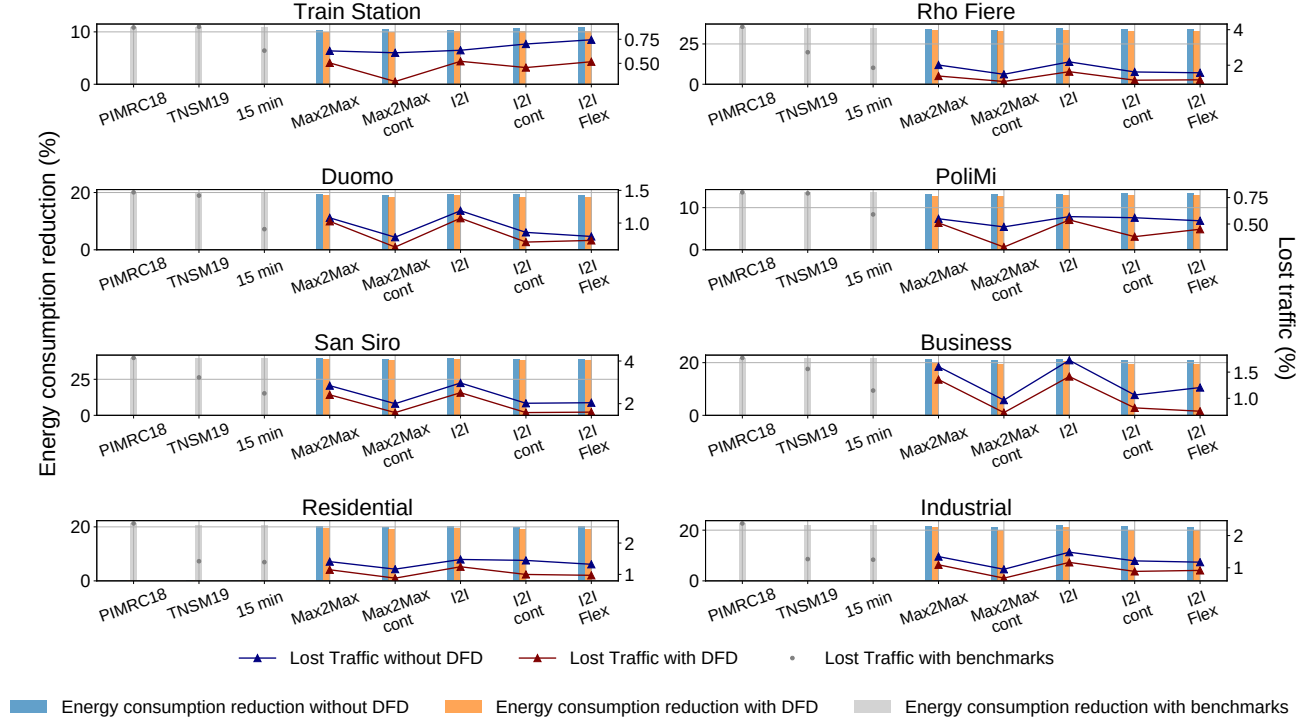


Figure 9: Energy consumption reduction and lost traffic in each area, with each dynamic resource allocation with and without descending front detection, DFD.

indicated by the blue and red lines with circle markers. Blue markers refer to
no DFD, while red markers refer to DFD. The results of the chosen benchmarks
are reported in grey. The figure reveals that the usage of DFD generates a
systematic drop in both energy efficiency, for a small amount, and lost traffic,
for more significant values. The energy consumption reduction remains between
10% and 39%, similar to the case of no DFD, when *TNSM19*, *15 min* and
PIMRC18 are used, but QoS improves significantly: lost traffic is usually below
1%. In the San Siro and Rho Fiere areas, because of the critical characteristics
of traffic patterns, this value is between 1% and 1.5%. The reductions of lost
traffic are due to the stricter conditions to switch off the micro cell BSs. This
can be seen in Fig. 8b, which illustrates an example of the traffic demand, in
black, of the 2 small cells BSs of Fig. 8a. In Fig. 8b, the orange and blue points

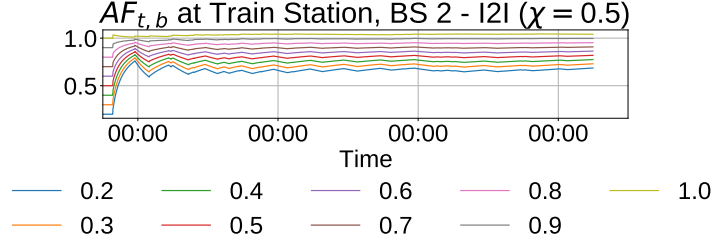


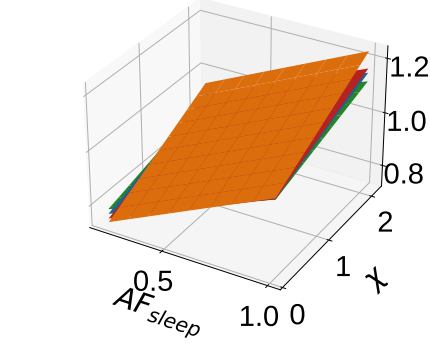
Figure 10: $AF_{t,b}$ for a micro cell BS of the Train Station area, with $\chi=0.5$, varying AF_{sleep}

indicate the time slot during which the considered micro cell BS is in sleep mode when *I2I cont* and *I2I cont with DFD* are used, respectively. During periods of almost constant but noisy traffic demand, if traffic values are close to the threshold ρ^* , incorrect small cell BSs deactivations may occur. Indeed, for those traffic values, even a small error in the traffic predictions can determine a wrong allocation of resources. This is the case reported in the figure: with DFD, incorrect deactivation of the considered small cell BSs is avoided since a descending front is not detected. Without DFD, with the *I2I cont* alone, the estimation error (even if small) makes the predictions lower than ρ^* , and a wrong switch off decision is taken. With DFD, the small cell BS is not switched off because the descending front is not detected. This behaviour explains the slight increase of energy consumption when DFD is applied. However, in spite of a very limited raise in energy consumption, the traffic loss can be reduced by up to 74% with respect to the benchmarks.

8.4. Impact on the BS failure rate

The impact of the proposed dynamic resource allocation schemes on the BS failure rate is now discussed. Each curve in Fig. 10 represents the behaviour of $AF_{t,b}$ versus time for a micro cell BS in the Train Station area, obtained with a different value of AF_{sleep} , when *I2I* is employed, with χ equal to 0.5, that is the value measured in an LTE BS in [16]. At the beginning of each simulation, $AF_{t,b}$ is lower than 1, since the simulation starts at midnight and the micro BS can be put in sleep mode, making $AF_{t,b}$ small. Then, at 7 a.m. it starts

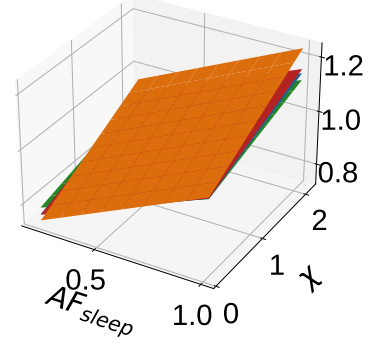
AF_b at Train Station, BS 2



■ I2I ■ I2I DFD
■ I2I cont ■ I2I cont DFD

(a)

AF_b at Train Station, BS 2



■ Max2Max ■ Max2Max DFD
■ Max2Max cont ■ Max2Max cont DFD

(b)

Figure 11: AF_b for a micro cell BS of the Train Station area: (a) with I2I-based approaches and (b) with Max2Max-based approaches.

growing, since the BS is activated due to increasing daily traffic demand. After some fluctuations, due to BS activation and deactivation that follow the daily traffic demand variations, $AF_{t,b}$ stabilises, since τ_{sleep} and f_{tr} stabilise as well. When large values of the parameter AF_{sleep} are considered, $AF_{t,b}$ is large due to more significant BS deterioration in sleep mode.

Fig. 11 reports the value of AF_b , on the z-axis for a BS in the Train Station area; different values of AF_{sleep} in the interval $[0.1, 0.9]$ are considered on the x-axis, and different values of χ in $[0.1, 2.0]$, on the y-axis. Each plotted plane corresponds to a different dynamic resource allocation approach. In particular, the *I2I* and *Max2Max* strategies are considered, with and without the *cont* variant and the use of *DFD*. From these figures, we first notice that the growth of AF_{sleep} and χ implies a growth of AF_b . If AF_{sleep} is large, the time in sleep mode is less beneficial to the BS failure rate; while large values of χ corresponds to the growth of the cost of the BS switching, see (3). When AF_{sleep} and χ are large enough, more conservative dynamic resource allocations provide

lower values of AF_b , than more dynamic ones. Indeed, the *DFD* variant, which uses the strictest switching conditions, provides the lowest AF_b , because of the reduction of the switching frequency f_{tr} , without significant reduction of the sleeping time ratio τ_{sleep} , see Figs. 6b, 6c, 6e and 6f. The largest values of AF_b are obtained when the *cont* variation is used, since, as mentioned in the previous section, it promptly reacts to the low traffic demand. This increases f_{tr} and, as a consequence, AF_b . For small values of AF_{sleep} and χ (bottom left part of the plots in Fig. 11), the situation is different. Indeed, with small values of χ and AF_{sleep} the cost of a BS switching does not significantly impact the BS failure rate, and spending time in sleep mode largely decreases it. Therefore, in this interval of values, the approaches which put the micro cell BS in sleep mode for longer time, provide lower values of AF_b , as in the case of *cont* variants.

Fig. 12 combines energy consumption and AF by representing each dynamic resource allocation algorithm in each area with a marker positioned so that the y coordinate corresponds to the energy consumption reduction and the x coordinate corresponds to the value of AF . Fig. 12a reports the AF values, with AF_{sleep} and χ equal to 0, which corresponds to the ideal case in which in sleep mode the BS failure rate goes to 0, meaning that its lifetime goes to infinity, and the BS switching does not affect it. Results in Fig. 12b are provided for $AF_{sleep} = 0.2$ and $\chi = 0.5$, as measured in [16] for an LTE BS. Finally, in Fig. 12c, the parameters are set pessimistically to 0.9 and 1.9, meaning that the sleep mode only slightly reduces the BS failure rate and the BS switching is highly costly, significantly affecting the BS deterioration. As expected and discussed in section 8, results are clustered according to the geographical area, because of the different achieved energy consumption reduction, which strictly depends on the traffic pattern that is characteristic of each zone. When AF_{sleep} and χ are 0, AF is always lower than 0.7, meaning that the average failure rate of the BS is decreased by 30%. In addition, we notice that AF is directly proportional to the energy consumption reduction, since AF is only affected by the time spent in sleep and active mode, not by the BS switching. As a result, the *15 min* approach always provides the lowest AF , since it rapidly reacts to

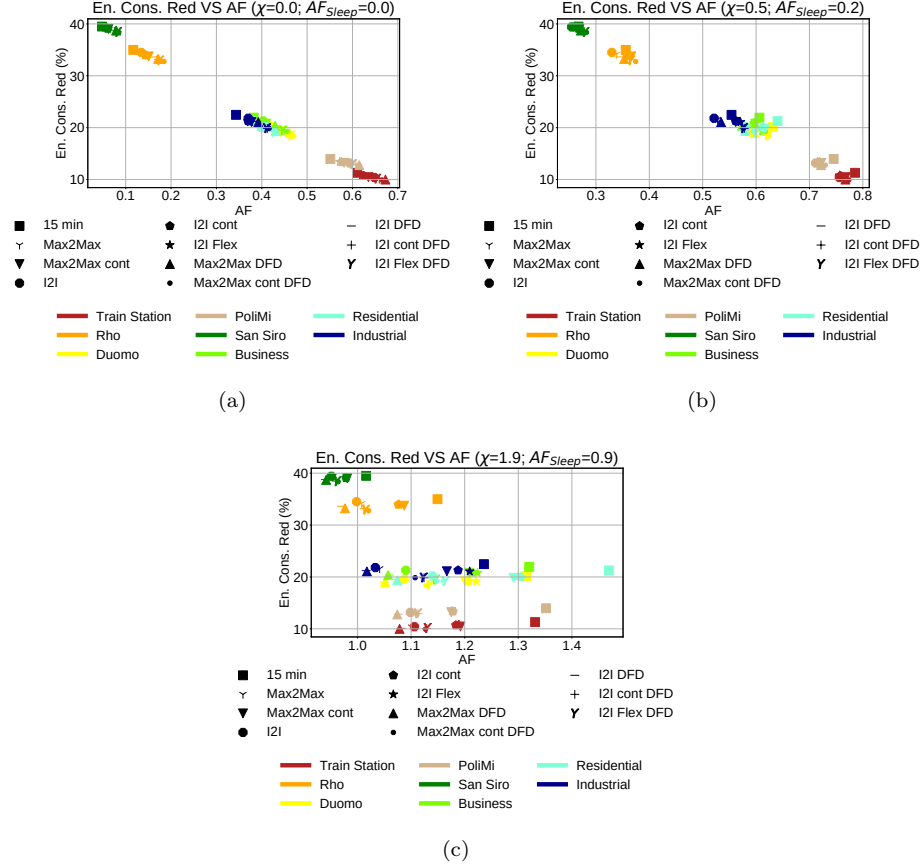


Figure 12: Energy Consumption Reduction and AF , obtained using different dynamic resource allocation, in each area, with χ and AF_{sleep} equal to (a) 0, (b) 0.2, 0.5 and (c) 1.9, 0.9.

the low traffic demand, immediately turning the micro cell BSs into sleep mode.

When AF_{sleep} and χ increase, the dynamism of this resource allocation approach negatively impacts the AF , which results the largest among the ones provided by our strategies, see Figs. 12b and 12c. Indeed, when the parameter χ grows, each BS switching is very costly. Thus, when AF_{sleep} and χ are 0.2 and 0.5, *I2I* and *Max2Max* approaches provide the lowest value of AF , since the most suitable balance between τ_{sleep} and f_{tr} is achieved. This does not occur with the *cont* balance between τ_{sleep} and f_{tr} . This does not occur with the *cont* variant: the large values of τ_{sleep} are not sufficient to compensate for the large values of f_{tr} . Similarly, with the adoption of *DFD*, the small values of τ_{sleep} ,

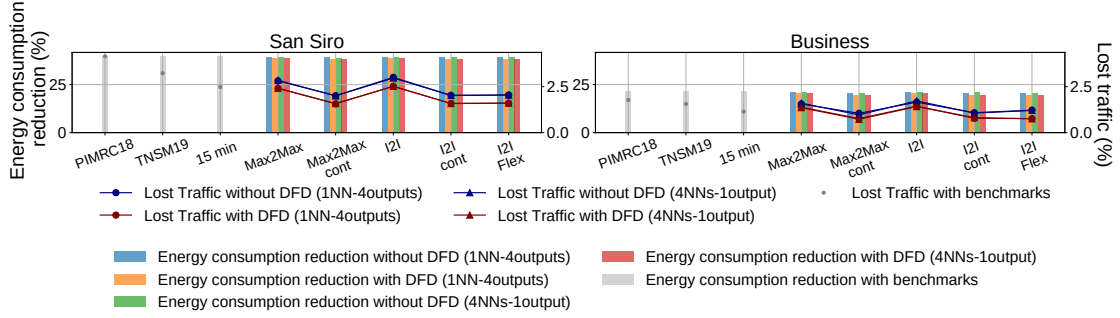


Figure 13: Energy consumption reduction and lost traffic in each area, with each dynamic resource allocation with and without descending front detection, DFD, in San Siro and Business areas, using 1 ANN-4 outputs and 4 ANNs-1 output.

because of the strict conditions for the BS deactivation, generate larger values of AF than with $I2I$ and $Max2Max$. If AF_{sleep} and χ are 0.9 and 1.9, the situation further worsens and AF is usually larger than 1, meaning that dynamic resource allocation increases the BS failure rate, because of the high cost of switching and the low benefit of being in sleep mode. Only for the San Siro area, values of AF lower than 1 are obtained because of the very long time the micro cell BSs spend in sleep mode. In this scenario, $I2I$ DFD and $Max2Max$ DFD are needed to reach the minimum AF values, since their strict deactivation requirements prevent frequent highly costly switching. Furthermore, under high values of AF_{sleep} and χ , the variable traffic patterns observed in the different zones can make even more critical the selection of the proper resource allocation scheme, whose impact on AF may result more significant. Indeed, whereas in the Train Station area the worst performing prediction algorithm increases AF by about 10% with respect to the lowest values obtained under $I2I$ DFD and $Max2Max$ DFD, in other traffic zones, like the Residential area, the worst performing scheme provide AF values that result up to almost 40% higher than the AF value under the best approach, that is anyway larger than 1.

655 8.5. Impact of the traffic prediction technique

Finally, let us consider the impact of the traffic prediction technique. Fig. 13 reports with the blue and the orange bars, the energy consumption reduction achieved with and without DFD, if the traffic demand is forecast with *4 ANNs-1 output*, in San Siro and Business areas. The resulting lost traffic is 660 shown with the red and blue lines with triangle markers, if the DFD is used or not, respectively. Similarly, the green and red bars in Fig. 13 indicate the energy consumption reduction obtained with and without DFD, when the traffic demand is forecast with *1 ANN-4 outputs*, which is the ANN that we identified as the one performing worst in predicting traffic. The obtained lost traffic is 665 reported, respectively, with the red and blue lines with circle markers. In spite of the larger estimation error with respect to *4 ANNs-1 output* (see Table 2), performance is very similar: the values of lost traffic and energy consumption are almost equal to the previous case. Indeed, lost traffic drops up to 1%, while energy consumption is reduced between 9% and 40%. Similar results are 670 achieved in the other areas. This means that the choice of an effective processing algorithm can have more impact on performance than the choice of the ANN. Only with a careful processing, the ANN prediction errors are mitigated, and a good trade-off between energy consumption reduction and QoS is achieved.

9. Lesson Learnt

675 In this section we discuss the main aspects which have emerged in our work. First, allocation of heterogeneous hierarchical RAN resources according to the traffic demand is promising, but the provided energy saving of each BS is strictly related to its traffic demand pattern, as well as to the traffic patterns over the whole considered area. Dynamic resource allocation requires the knowledge of 680 the actual traffic demand and, hence, machine learning approaches are needed to accurately predict it so as to enable network management mechanisms that adapt to traffic variability. This is interesting in perspective, for the promising possibilities offered toward the deployment of new networks that are easily and

automatically reconfigurable. However, machine learning approaches become
685 particularly effective only if their outputs are integrated into decision processes
that are driven by a deep domain knowledge, which cannot be eliminated if
the desired objectives are to be achieved. If the traffic predictions are carefully
processed, QoS deterioration is avoided, while significant energy saving can be
achieved. Prediction processing requires both the understanding of traffic pat-
690 terns over long time scales, so as to detect the overall trend of increasing or
decreasing traffic, as well as strategies to combine predictions at different time
lags.

Finally, prediction processing and the consequent dynamic resource allocation
affect the BS failure rate in different ways. Switching a BS is harmful to its
695 failure rate while the time spent in sleep mode prevents its deterioration. The
actual impact of the combination of these two phenomena depends on the HW
components of the BS, as well as on the RAN management strategy. In case
the switching of a BS is not costly, less strict switching conditions can be ap-
plied: the BS failure rate is not affected while larger energy saving is achieved.
700 Conversely, when the BS is sensitive to switching, more conservative resource
allocations should be employed. For existing networks, not designed for highly
dynamic resource allocation, conservative approaches better prevent BSs from
HW failure; however, in perspective, with the deployment of new devices suited
for strongly dynamic networks, less conservative approaches, which frequently
705 activate and deactivate BSs, can be used, and higher energy saving is expected.

10. Conclusions

In this paper, the traffic demand of BSs of a portion of a RAN is forecast
with the objective of enabling BS management strategies that aim at reducing
the RAN energy consumption. Results show that, in order to achieve good per-
710 formance trade-offs, measured in energy saving, QoS and impact on BS failure
rate, the traffic predictions need to be carefully processed, understanding the
traffic patterns over long time scales, detecting the overall trend of increasing

or decreasing traffic, as well as combining predictions at different time lags. As
 next steps of our work, we will design dynamic RAN management strategies
 715 that optimise both the energy consumption and the BSs lifetime.

References

- [1] K. Johansson, A. Furuskar, P. Karlsson, J. Zander, Relation between base
 station characteristics and cost structure in cellular systems, in: 2004 IEEE
 15th International Symposium on Personal, Indoor and Mobile Radio Com-
 720 munications (IEEE Cat. No. 04TH8754), Vol. 4, IEEE, 2004, pp. 2627–
 2631.
- [2] D. Renga, M. Meo, Dimensioning renewable energy systems to power mo-
 bile networks, IEEE Transactions on Green Communications and Network-
 ing 3 (2) (2019) 366–380. doi:10.1109/TGCN.2019.2892200.
- 725 [3] D. Pompili, A. Hajisami, T. X. Tran, Elastic resource utilization framework
 for high capacity and energy efficiency in cloud ran, IEEE Communications
 Magazine 54 (1) (2016) 26–32.
- [4] D. Sabella, A. De Domenico, E. Katranaras, M. A. Imran, M. Di Girolamo,
 U. Salim, M. Lalam, K. Samdanis, A. Maeder, Energy efficiency benefits of
 730 ran-as-a-service concept for a cloud-based 5g mobile network infrastructure,
 IEEE Access 2 (2014) 1586–1597.
- [5] G. Vallero, D. Renga, M. Meo, M. A. Marsan, Greener ran operation
 through machine learning, IEEE Transactions on Network and Service
 Management 16 (3) (2019) 896–908.
- 735 [6] G. Vallero, D. Renga, M. Meo, M. Ajmone Marsan, Processing ann traffic
 predictions for ran energy efficiency, in: Proceedings of the 23rd Interna-
 tional ACM Conference on Modeling, Analysis and Simulation of Wireless
 and Mobile Systems, 2020, pp. 235–244.

- [7] L. Budzisz, F. Ganji, G. Rizzo, M. A. Marsan, M. Meo, Y. Zhang, G. Koutitas, L. Tassiulas, S. Lambert, B. Lannoo, et al., Dynamic resource provisioning for energy efficiency in wireless access networks: A survey and an outlook, *IEEE Communications Surveys & Tutorials* 16 (4) (2014) 2259–2285.
- [8] T. Shankar, et al., A survey on techniques related to base station sleeping in green communication and comp analysis, in: 2016 IEEE International Conference on Engineering and Technology (ICETECH), IEEE, 2016, pp. 1059–1067.
- [9] S. Buzzi, I. Chih-Lin, T. E. Klein, H. V. Poor, C. Yang, A. Zappone, A survey of energy-efficient techniques for 5g networks and challenges ahead, *IEEE Journal on Selected Areas in Communications* 34 (4) (2016) 697–709.
- [10] M. Dalmaso, M. Meo, D. Renga, Radio resource management for improving energy self-sufficiency of green mobile networks, *ACM SIGMETRICS Performance Evaluation Review* 44 (2) (2016) 82–87.
- [11] H. Ghazzai, M. J. Farooq, A. Alsharoa, E. Yaacoub, A. Kadri, M.-S. Alouini, Green networking in cellular hetnets: A unified radio resource management framework with base station on/off switching, *IEEE Transactions on Vehicular Technology* 66 (7) (2017) 5879–5893.
- [12] N. B. Rached, H. Ghazzai, A. Kadri, M.-S. Alouini, A time-varied probabilistic on/off switching algorithm for cellular networks, *IEEE Communications Letters* 22 (3) (2018) 634–637.
- [13] D. Renga, H. A. H. Hassan, M. Meo, L. Nuaymi, Energy management and base station on/off switching in green mobile networks for offering ancillary services, *IEEE Transactions on Green Communications and Networking* 2 (3) (2018) 868–880.
- [14] M. Ali, M. Meo, D. Renga, Cost saving and ancillary service provisioning

in green mobile networks, in: *The Internet of Things for Smart Urban Ecosystems*, Springer, 2019, pp. 201–224.

[15] L. Chiaraviglio, M. Listanti, E. Manzia, Life is short: The impact of power states on base station lifetime, *Energies* 8 (12) (2015) 14407–14426.

770 [16] L. Chiaraviglio, F. Cuomo, M. Listanti, E. Manzia, M. Santucci, Fatigue-aware management of cellular networks infrastructure with sleep modes, *IEEE Transactions on Mobile Computing* 16 (11) (2017) 3028–3041.

[17] L. Chiaraviglio, F. Cuomo, M. Listanti, E. Manzia, M. Santucci, Sleep to stay healthy: Managing the lifetime of energy-efficient cellular networks, 775 in: *2015 IEEE Global Communications Conference (GLOBECOM)*, IEEE, 2015, pp. 1–7.

[18] C. Natalino, L. Chiaraviglio, F. Idzikowski, C. R. Francês, L. Wosinska, P. Monti, Optimal lifetime-aware operation of green optical backbone networks, *IEEE Journal on Selected Areas in Communications* 34 (12) (2016) 780 3915–3926.

[19] P. Wiatr, J. Chen, P. Monti, L. Wosinska, Energy efficiency versus reliability performance in optical backbone networks, *Journal of Optical Communications and Networking* 7 (3) (2015) A482–A491.

[20] D. Bega, M. Gramaglia, M. Fiore, A. Banchs, X. Costa-Perez, Deepcog: 785 Cognitive network management in sliced 5g networks with deep learning, in: *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, IEEE, 2019, pp. 280–288.

[21] J. Guo, Y. Peng, X. Peng, Q. Chen, J. Yu, Y. Dai, Traffic forecasting for mobile networks with multiplicative seasonal arima models, in: 790 *2009 9th International Conference on Electronic Measurement & Instruments*, IEEE, 2009, pp. 3–377.

- [22] P. Cortez, M. Rio, M. Rocha, P. Sousa, Multi-scale internet traffic forecasting using neural networks and time series methods, *Expert Systems* 29 (2) (2012) 143–155.
- 795 [23] M. Z. Shafiq, L. Ji, A. X. Liu, J. Wang, Characterizing and modeling internet traffic dynamics of cellular devices, *ACM SIGMETRICS Performance Evaluation Review* 39 (1) (2011) 265–276.
- [24] P. Cortez, M. Rio, M. Rocha, P. Sousa, Internet traffic forecasting using neural networks, in: *The 2006 IEEE International Joint Conference on Neural Network Proceedings*, IEEE, 2006, pp. 2635–2642.
- 800 [25] H. D. Trinh, L. Giupponi, P. Dini, Mobile traffic prediction from raw data using lstm networks, in: *2018 IEEE 29th Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, IEEE, 2018, pp. 1827–1832.
- 805 [26] S. Troia, R. Alvizu, Y. Zhou, G. Maier, A. Pattavina, Deep learning-based traffic prediction for network optimization, in: *2018 20th International Conference on Transparent Optical Networks (ICTON)*, IEEE, 2018, pp. 1–4.
- [27] L. Yao, T.-S. Tsai, Novel hybrid scheme of solar energy forecasting for home energy management system, in: *2016 IEEE International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)*, IEEE, 2016, pp. 150–155.
- 810 [28] S. Wang, J. Guo, Q. Liu, X. Peng, On-line traffic forecasting of mobile communication system, in: *2010 First International Conference on Pervasive Computing, Signal Processing and Applications*, IEEE, 2010, pp. 97–100.
- 815 [29] H. Pan, J. Liu, S. Zhou, Z. Niu, A block regression model for short-term mobile traffic forecasting, in: *2015 IEEE/CIC International Conference on Communications in China (ICCC)*, IEEE, 2015, pp. 1–5.

- 820 [30] G. Auer, O. Blume, V. Giannini, I. Godor, M. Imran, Y. Jading, E. Katranaras, M. Olsson, D. Sabella, P. Skillermark, et al., D2. 3: Energy efficiency analysis of the reference systems, areas of improvements and target breakdown, *Earth* 20 (10).
- [31] L. Chiaraviglio, P. Wiatr, P. Monti, J. Chen, J. Lorincz, F. Idzikowski,
825 M. Listanti, L. Wosinska, Is green networking beneficial in terms of device lifetime?, *IEEE Communications Magazine* 53 (5) (2015) 232–240.
- [32] S. Arrhenius, About the reaction speed during the inversion of cane sugar by acidic acids, *magazine for physical chemistry* 4 (1) (1889) 226–248.
- [33] Y.-L. Lee, J. Pan, R. Hathaway, M. Barkey, *Fatigue testing and analysis: theory and practice*, Vol. 13, Butterworth-Heinemann, 2005.
830