

Digital Twin Assisted Risk-Aware Sleep Mode Management Using Deep Q-Networks

Meysam Masoudi, Ebrahim Soroush, Jens Zander, and Cicek Cavdar

Abstract—Base stations (BSs) are the most energy-consuming segment of mobile networks. To reduce BS energy consumption, different components of BSs can sleep when BS is not active. According to the activation/deactivation time of the BS components, multiple sleep modes (SMs) are defined in the literature. In this study, we model the problem of BS energy saving utilizing multiple sleep modes as a sequential Markov decision process (MDP) and propose an online *traffic-aware* deep reinforcement learning approach to maximize the long-term energy saving. However, there is a risk that BS is not sleeping at the right time and incurs large delays to the users. To tackle this issue, we propose to use a digital twin model to encapsulate the dynamics underlying the investigated system and estimate the risk of decision-making (RDM) in advance. We define a novel metric to quantify RDM and predict the performance degradation. The RDM calculated by DT is compared with a tolerable threshold set by the mobile operator. Based on this comparison, BS can decide to deactivate the SMs, re-train when needed to avoid taking high risks, or activate the SMs to benefit from energy savings. For deep reinforcement learning, we use long-short term memory (LSTM), to take into account the long and short-term dependencies in input traffic, and approximate the Q-function. We train the LSTM network using the experience replay method over a real traffic data set obtained from an operator's BS in Stockholm. The data set contains data rate information with very coarse-grained time granularity. Thus, we propose a scheme to generate a new data set using the real network data set which 1) has finer-grained time granularity and 2) considers the bursty behavior of traffic data. Simulation results show that using proposed methods, considerable energy saving is obtained, compared to the baselines at cost of negligible number of delayed users. Moreover, the proposed digital twin model can predict the performance of the DQN proactively in terms of RDM hence preventing the performance degradation in the network in anomalous situations.

Index Terms—5G, Base station, Digital twin, sleep modes, deep learning, energy saving.

I. INTRODUCTION

Mobile network operators anticipate a surge in their network energy consumption due to exponentially increasing network traffic and emerging new services. Therefore, in order to maintain the network sustainability, *energy efficiency* becomes a key factor in the future deployment designs [1]. Studies show that base stations (BSs) consume around 80% of the overall network energy consumption. Thus, reducing BSs'

energy consumption can boost the energy efficiency of the network [2]–[5]. In particular, base station sleeping is one of the promising ways of reducing energy consumption in the mobile networks. During a day, there are multitudes but short durations in which no user is connected to the BS while it is active and consumes energy. In such idle durations, BS can put some of its components into sleep and hence save energy, however, it may lead to longer serving delay for a number of new arrivals. In this study, serving delay is defined as the time that a user is waiting to be served after its arrival. We look into this energy saving opportunity while keeping under control incurred serving delay to the users.

In the literature, the concept of BS sleeping for the purpose of energy saving has been thoroughly investigated [6]–[13]. The authors in [6], propose a traffic-aware BS sleeping strategy and derive a sleeping policy to maximize the energy saving. In [7], the authors propose a scheme to put BS into sleep mode (SM) until a time when a certain number of users are queued in the BS. In a similar setup, the study in [8] investigates the impact of bursty arrivals on the incurred delay and energy saving. The studies in [9]–[12] investigate the tradeoffs between total energy consumption and overall incurred delay to the users. The authors in [9], derive a closed-form energy–delay tradeoff which can be used in designing BS sleeping policies. In [10], the authors propose a framework for BS energy saving considering user association. In [11], the authors solve the problem of joint clustering and BS sleeping. In [12], the authors derive an optimal sleeping policy where the policy takes into account the maximum tolerable delay for the user. They propose a hysteresis-based sleeping strategy and show that there is a linear relation between amount of incurred delay and energy saving when the delay constraint is lower than a threshold. All the aforementioned studies consider only a single level SM.

Recently, multi level SMs have been introduced in the literature for potentially better energy saving performance at the BSs [14]. According to the (de)activation transition time of BSs' transceiver chain components, the study in [4] defines 4 levels of SMs known as advance sleep modes (ASMs) each with a different (i) time duration, and (ii) energy consumption:

- 1) SM1: This SM is the fastest, and the most shallow, whose duration is comparable to the symbol time T_s . Whenever a base station (BS) does not have any traffic over the entire band of sub-carriers during symbol time per antenna, the PA can be switched off to save power. This mode is available and known in current technologies as micro-scale discontinuous transmission (μ D TX).
- 2) SM2: A slightly deeper sleep level can be reached by

Copyright (c) 2015 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

M. Masoudi, J. Zander, and C. Cavdar are with the School of Electrical and Computer Science, KTH Royal Institute of Technology, Sweden, e-mail: {masoudi,cavdar,jenz}@kth.se. Ebrahim Soroush is with Zi-Tel Company, e-mail: ebrahim.soroush@gmail.com. This work is supported by European Celtic-Next AI4Green project.

switching off one more set of components with slower (de)activation transition time than that of SM1, if this longer delay can be afforded by the system. In this mode, the transition is on the transmission time interval (TTI) scale, i.e., 1 msec (a duration of a *sub-frame* constituting 2 resource blocks (RBs) or 14 resource elements (REs) assuming frequency division duplexing (FDD) Frame Structure 1 [15])

- 3) SM3: This mode has a deeper sleeping level but with a minimum duration of a 10 msec frame (10 sub-frames).
- 4) SM4: This is the deepest SM with a minimum duration of 100 frames (1 sec).

Implementing such SMs can save more energy at the cost of incurring more delay to the users that have to wait until the BS is back in its active state. Therefore, novel yet efficient schemes are required to utilize ASMs without any adverse impact on the offered quality-of-service (QoS).

The problem of deciding on the sleeping as well as the level of SMs are complex. Machine learning (ML) techniques have shown promising capabilities to efficiently handle the complexity and tune the parameters in systems with ASMs [16]–[23]. The authors in [16] evaluate the energy consumption of 5G and beyond BS with multi level SMs. The authors in [17]–[19] investigate 4-level SMs with pre-defined order of (de)activation of SMs. They decide on the number of repetitions for each SM. In [17] a heuristic algorithm is proposed to implement the ASMs. They investigate the impact of the synchronization signaling periodicity on the energy saving performance of ASMs. In [18], the authors propose a dynamic algorithm based on Q-learning for SM selection. In [19], the authors propose a traffic-aware strategy to build a codebook to map the traffic load to the possible actions using Q-learning algorithm. In these studies, the authors assume a pre-defined order of the decisions which may not result in optimal energy saving. The authors in [20] relax the pre-defined order of SMs and propose a traffic-aware reinforcement learning algorithm to select SMs. However, they train their network on a fixed traffic pattern. The authors in [21] use reinforcement learning to adjust the BSs' configurations, e.g., bandwidth and MIMO parameters, to increase the sleeping time of BS. Since control and signaling can limit the energy saving of ASMs, the study in [22] proposes a control/data plane separation in 5G BSs which allows implementing SMs with longer durations. This separation can also be leveraged in other network settings such as co-coverage scenarios when basic coverage cells may carry all periodic control signaling, leaving capacity-enhancing cells with higher ability to sleep and save energy. The study in [23], performs BS sleeping for the capacity cell using tabular reinforcement learning method which is trained over real network data. However, the impact of abnormal or bursty arrivals of users on the BS sleeping energy performance are not investigated.

ML-based algorithms utilize data extracted from the network to learn the network behavior and to make proper decisions to improve the network performance. However, ML-based techniques performance may deteriorate due to 1) lack of good training 2) abnormal behavior of data that has not been observed before 3) inaccurate design of the algorithm,

etc. When the network is experiencing a state that is not observed before, the performance of ML-based techniques is questionable. Furthermore, the traffic pattern may change during time and re-training may be required. It is necessary to tackle these open issues and answer questions such as “Under what circumstances ML-algorithm performance is not reliable in the cellular network?”, “when is it needed to re-train?” and “if an ML-based algorithm is trained in a BS, can it be used in other BSs?” These questions are stemmed from a fact that there is always a risk associated to the use of ML-based algorithms under different circumstances. Consequently, there is a crucial need to quantify and measure the potential risk of using these algorithms in the network which is a critical challenge and an open problem, not tackled by existing solutions.

Current systems utilize conservative approaches to deal with the risk of activating any energy saving feature in the network. For instance, most energy saving features are never activated, and few are used only at night. Currently, there is no intelligent risk management approach for energy saving, deployed in the network. In order to determine the proper time of activation/deactivation, threshold-based approaches are used in which when a certain KPI goes below a certain threshold, energy saving feature is turned off. This approach is reactive in which it only considers instantaneous behavior, without involving any prediction, learning, or proactive approaches. However, if the risks can be predicted in advance, these features can be activated or deactivated dynamically and hence the network can benefit more by using these features.

One way of predicting the risk of using ML algorithms is to utilize the paradigm of digital twin (DT) [24]. A DT is a virtual representation of a physical entity which can mimic the behavior of the physical entity to simulate, predict, forecast behaviour, and possibly control the entity. This representation can be a data driven or analytical model of the system [25]. A typical DT-assisted system has three main parts: 1) physical entities 2) virtual entities and 3) the interaction interface between them [25]. DT can be integrated to the next generation of mobile networks, e.g., 6G, to support the realization of the smart control and optimization of the network [24], [26]. DT's applications are quite wide. It can focus on the end-to-end view of the network providing a higher-level analysis of services or network infrastructure, or can focus on a specific domain like the RAN, e.g., ASMs. In this study, DT is used to realize risk-aware intelligent ASMs decision making at the BSs. In particular, DT is composed of a data driven and analytical representations of a BS using ASMs. Using DT, we can predict the performance of the ASM management algorithm in advance, given the required parameters. Based on the performance evaluation, the system can raise an early alarm to deactivate the ASMs and check for possible algorithm re-training, provided that a risky situation is identified. Risky situation occurs when data used to train the algorithm becomes invalid or network performance is predicted to be below a tolerable threshold determined by the mobile operator. Hence, assisted by DT, we can identify the times in which using sleep modes will cause unacceptable performance degradation, i.e., imposing delay to large number of users.

In this paper, using a real network data set obtained from a Swedish operator, we propose a deep Q-learning (DQN) algorithm and a DT representation of the system to save energy at BS and avoid unexpected performance degradation of the DQN algorithm at the BS. Preliminary version of this study, is published in [23]. The main contributions of this study are summarized as follows:

- We obtain a data set from an operator in Sweden. We propose a framework to generate fine time granularity data from the real network data considering the bursty behavior of cellular network traffic.
- We model the problem of sleep mode management in 5G BSs as a Markov Decision Process. We propose a DQN algorithm to solve the problem of minimizing the energy consumption of a BS using ASMs while maintaining the users' QoS. Given the traffic input, the proposed algorithm can dynamically decide on the level and duration of ASM.
- We propose a DT model-based on the Markov model that can characterize and predict the performance of the proposed sleep mode management algorithm.
- We propose a novel metric to quantify the risk of decision making (RDM) and propose a framework to monitor the network and prevent the BS from taking risky decisions.
- We show that combining the RDM metric, and the proposed framework for risk monitoring, DT can detect the abnormal behavior of traffic and can retrain the algorithm based on new observed data, if required.
- Through analytical and simulation results, we show that when the RDM is high or there is a miss match between DT and physical network observation, the sleeping features can be deactivated temporarily and the algorithm can be re-trained with new data to avoid the risk of incurring large serving delay to the users.

The rest of the paper is organized as follows. In Section II, we explain the system model, adopted power and traffic generation models. In Section IV, we explain the proposed deep Q-learning algorithm. In Section III-C, we present the Markov model of the sleep mode management algorithms. We evaluate the performance of the proposed algorithm in Section V. We conclude the paper in Section VI.

II. SYSTEM MODEL

A. Network Scenario

In this study, we consider a non-stand alone 5G deployment in which there are two BSs, namely capacity BS (CapBS) and coverage BS (CovBS). The CovBS is an LTE BS and provides ubiquitous coverage and cell specific signaling. The CapBS is responsible for providing high data transmission rate, on-demand services, and user-based signaling. Since CovBS is responsible for providing coverage and signaling, it cannot take advantage of deep sleep modes. However, CapBS can go to sleep and wake up when it is needed. In this study, we focus on CapBS which from now on we refer to it as BS. We also consider that BS can go to SM1-SM3 when there is no user in the network. Table I presents the power consumption and minimum duration of each SM for macro-base stations with

TABLE I: Power and Minimum Durations from [4], [27]

Mode	Active		SM1	SM2	SM3
	100%	0%			
Pow. (W)	702.6	114.5	76.5	8.6	6.0
Dur. (sec)	T_s : fast modes (FMs)		1 msec	10 msec	

max power of 49 dBm over 3 sectors with a bandwidth of 20 MHz [4].

In mobile networks, BSs are required to communicate with users with periodic signaling that are transmitted in the first and fifth LTE sub-frames. Hence it is not possible to leverage from long and deep SMs. For this reason, 5G New Radio (5G NR) standards allow operators to adjust their periodicity of synchronization signals to 5, 10, 20, 40, 80 and 160 ms, which makes it possible to leverage from deeper/longer SMs. 5G deployments support standalone and non-standalone deployment. In standalone deployment, operators deploy a separate 5G core and NR for 5G BSs. In non-standalone deployment, the network uses the existing LTE radio access for signaling and existing evolved packet core for backhauling. This approach allows the operators to introduce 5G new services while reusing the existing LTE networks. Moreover, in such networks it is possible to separate control and data plane and utilize 5G BSs for data provisioning and use LTE BSs for signaling.

B. Power Model

The adopted BS power consumption model follows the one implemented by IMEC in [27] and reported in [4]. Equation (1) formulates the operational power consumption of a BS and Table I summarizes the power consumption values with the configuration parameters summarized in Table II. When BS is active, it consumes $P_{idle} + P_t$ where the first term is the idle power consumption when no SM is activated, and there is no active user serving, and P_t is the dynamic power consumption which is proportional to the BS's load. When BS is in sleep mode, the operational power consumption of BS operating in SM i is P_{SM_i} . In this study, the BS can operate in three modes, namely fast mode (FM), SM2, and SM3. FM merges three cases: 1) active with 0 to 100% load, 2) idle with 0% load, and 3) SM1. In this mode, BS makes a decision and hops instantaneously into any operational mode depending on the arriving traffic without any need for optimization. The remaining two sleep modes incur a transition delay from the time of decision until completion of the action, which causes an additional queuing delay of traffic arriving in the transition period.

$$P_c = \begin{cases} P_{idle} + P_t & \text{If BS is in active mode} \\ P_{SM_i} & \text{If BS is in SM}_i, \quad i \in \{1, 2, 3\} \end{cases} \quad (1)$$

The BS consumes transient power for switching, between operational modes [28]. Denote the total power consumption of switching in duration of T by P_{sw}^T , power consumption per switching by P_{sw} , and the number of switching in duration of T by F_m , then we have

$$P_{sw}^T = F_m P_{sw} \quad (2)$$

C. Traffic Model

1) *Data Set Description*: Cellular traffic varies hourly depending on the daily activity of users in the area. Understanding the traffic patterns of cellular networks is extremely valuable for better network resource management. In this study, we use a data set from a Swedish mobile network operator. The data set is an anonymized averaged data rate for both uplink and downlink, collected by the operator from multiple BSs in Stockholm for two months in 2018, with one sample taken every 5 minutes. This coarse-grained data set guarantees the credibility of our daily traffic pattern analysis and modeling. This data set can provide us with long term dependencies of network traffic, however, this information is too coarse to provide information about very short time instances comparable to sleep mode duration which is in order of milliseconds. On the other hand, it is not possible to obtain data in such fine-grained granularity due to the limitations in the deployed equipment at BSs measuring the KPIs. Therefore, we should adopt a model to generate traffic data with short term dependencies while maintaining the daily pattern and long term dependencies of traffic data.

2) *Distribution of traffic*: The data collected by the operator needs to be preprocessed due to the existence of the incomplete information in the provided data rates at some time instances. Moreover, the information on the average provided rate cannot be directly used in our study, since we need the demand per user. The preprocessing includes the following steps. First, we eliminate the empty and incomplete information from the data set. Then, we focus on the aggregated downlink traffic of one of such BSs in one location in Stockholm, only considering the FDD downlink band for simplicity, and without loss of generality. Very similar studies can be performed to simultaneously take into account all different BS bands, considering both uplink and downlink. This coarse grained original data set is further used as follows to obtain the required information to generate short term traffic data.

For data preparation, we first collate the average data rates of same time index over all days into one set. Denote the set of data at time index i by \mathcal{O}_i :

$$\mathcal{O}_i = \{o_1, o_2, \dots, o_\kappa\}, \quad \text{for } i \in \mathcal{T} \quad (3)$$

where κ is 60 in our study since we have data for 60 days and \mathcal{T} denotes the set of the time index of data during the day. The cardinality of set \mathcal{T} is $24 * 60/t_s$ where t_s is sampling time which is 5 minutes in our study. Then, for each set, i.e., each 5 minutes, we calculate the mean, i.e., $E(\mathcal{O}_i)$, and variance, i.e., $Var(\mathcal{O}_i)$, of the original data. Finally, based on derived statistics, we can generate demanded data rate for each time step and each user, using the steps explained in the next subsection. It is worth mentioning that this coarse grained time granularity data provides no information about the time of arrivals or the time duration between two arrivals which is a key parameter to determine the idle time of the BS. Therefore, an accurate and yet tractable model is required to generate the random arrivals of the users.

3) *Bursty Behavior of Traffic Data*: In order to model the user arrivals, various distributions such as the Poisson distri-

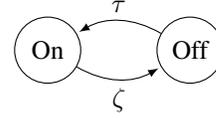


Fig. 1: Data traffic model as interrupted Poisson process

bution [12], log normal distribution [19], hyper exponential distribution [18], and interrupted Poisson process (IPP) [29] have been used in the literature. Among these models, IPP and hyper exponential distribution can model the bursty behavior of incoming traffic while they are mathematically tractable [29]. It can be shown that the IPP is stochastically equivalent to a hyper exponential H2 renewal process [30].

In IPP, as illustrated in Fig. 1, arrivals have two states, i.e., ON and OFF states. In ON state, arrivals follow the Poisson model with parameter λ and in OFF state there is no arrival. Assume that the transition rate from ON to OFF (OFF to ON) is denoted by ζ (τ), the steady state probabilities are,

$$P_{ON} = \frac{\tau}{\tau + \zeta}, \quad P_{OFF} = \frac{\zeta}{\tau + \zeta} \quad (4)$$

In order to complete the steps of data generation, during ON state users are generated according to the Poisson distribution with parameter λ . In OFF state, there will be no arrivals i.e., $\lambda = 0$. Denoting U as the number of arrivals per given time T , the mean and variance of U is given by [30]:

$$E(U) = \frac{\lambda\tau}{\tau + \zeta}T \quad (5)$$

$$Var(U) = \frac{\lambda\tau T}{\tau + \zeta} + \frac{2\lambda^2\tau\zeta T}{(\tau + \zeta)^3} \left[1 - \frac{(1 - e^{-(\tau+\zeta)T})}{(\tau + \zeta)T} \right] \quad (6)$$

For $(\tau + \zeta)T \gg 1$, the variance can be approximated as,

$$Var(U) \approx \frac{\lambda\tau T}{\tau + \zeta} \left[1 + \frac{2\lambda\zeta}{(\tau + \zeta)^2} \right]. \quad (7)$$

4) *Data generation*: In previous two subsections, we derived the statistics of the number of arrivals and the data sets. In this section we explain how to map the parameters of the IPP model so that it matches with the statistics of the original data set. Denote the data rate request per arrival by ψ_j and the aggregated data rate request per given time T by Ψ .

$$\Psi = \sum_{j=1}^U \psi_j \quad (8)$$

$$E(\Psi) = E(U)E(\psi) \quad (9)$$

$$Var(\Psi) = E(U)Var(\psi) + Var(U)E(\psi)^2. \quad (10)$$

In Equation (8), ψ -s are independent random variables and for a sufficiently large value of U , e.g., more than 30, Ψ is normally distributed, regardless of distribution of ψ -s, with mean and variance of $E(\Psi)$ and $Var(\Psi)$, respectively.

In order to generate the user arrivals and their data rates with the original data set statistics, the model parameters, i.e., τ , ζ , λ , $E(\psi)$, and $Var(\psi)$, should be set in a way that Equations (9) and (10) hold.

Let us assume that ψ is exponentially distributed with mean and variance of $\bar{\psi}$ and $\bar{\psi}^2$, respectively. At each time index i , $E(\Psi)$ and $Var(\Psi)$ are set to $E(\mathcal{O}_i)$, and $Var(\mathcal{O}_i)$, respectively. Now, we have 4 unknown parameters and two equations. Since the information of these parameters is not available in the data set, we have to assume two parameters and derive the other two parameters. Let's assume that τ and ζ are given ¹, then by inserting (5) and (7) into (9) and (10), λ and $E(\psi)$ are derived as follows:

$$\lambda = \frac{\tau + \zeta}{\tau T \left[\frac{Var(\Psi)}{E(\Psi)^2} - \frac{2\zeta}{\tau T(\tau + \zeta)} \right]} \quad (11)$$

$$E(\psi) = \frac{(\tau + \zeta)E(\Psi)}{\tau \lambda T}. \quad (12)$$

With (11) and (12) we have all the parameters to generate the user arrivals with their data rate demand.

D. BS Model

In this study, we assume that the BS has two operating states namely, active mode and sleep modes. If BS is active, users are served with rate μ . At the same time, BS can put new service requests with arrival rate of λ . The BS goes to sleep mode only if the last user in the BS is served. If BS is in sleep mode, all new users must wait to be served until the BS turns into active mode again.

III. RISK-AWARE SLEEP MODE MANAGEMENT

In this section, we propose a framework for risk-aware sleep mode management. AI or Machine learning algorithms are prone to anomalous and unknown data and their performance can be inadequate or sub-optimal. A mechanism is needed to define the risk, monitor the performance of AI with respect to the input traffic to find out whether re-training is needed. As shown in Fig.2, we propose a DT that receives the network data as an input, communicates with and mimics the behaviour of the real system composed of the physical network and AI module. DT continuously assesses the risk by analyzing and predicting the performance pro-actively. Based on this risk assessment and prediction, the management framework decides whether to use the AI, retrain the network, or temporarily deactivate the AI module corresponding to deactivating SMs. In the following we explain this approach in detail. First, we introduce the concept of risk associated to the BS sleeping algorithms. We explain the DT model for BS ASM management to estimate the risk. Then, we combine the intelligent sleep mode management algorithm (SMA) and DT, and explain the framework for risk-aware BS sleep mode management, as is depicted in Fig.3.

A. Risk-Aware Sleep Mode Management

In this study, risk refers to a situation in which SM management algorithm takes sleep decisions resulting in delays

¹In the operators' deep packet inspection level data set, the time of service requests and the service durations are available. Hence, the statistics of the τ and ζ can be derived from the data set.

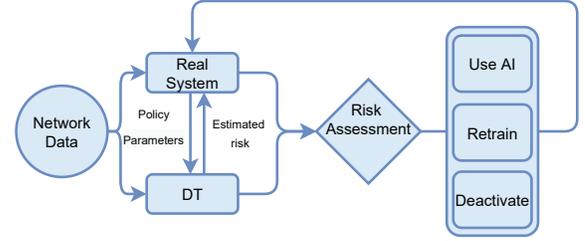


Fig. 2: Digital twin assisted decision making

in connection setup time of future user arrivals. The more users experiencing delay, the higher value of risk should be associated to the situation. Therefore, if the risk could be measured in advance, BS can avoid delaying large number of users. For this purpose, we utilize the DT concept as a virtual representation of the BS sleeping process.

In Fig.3, we illustrate the proposed DT assisted risk-aware SM management framework. In this scheme, the physical network provides the required information, e.g., duration of sleeping, ON/OFF state duration, sleeping time, and arrival rate, to the virtual network to update the model parameters to be used in analytical modeling in DT. The virtual network is composed of two modules. The first one updates the parameters according to the physical network (only for the initial parameters it uses the training data). The second one is a Markov model with the updated parameters. The virtual network together with the RDM prediction module construct the DT, details of which will be discussed later in Section III-C and III-D. These updated parameters and the performance metrics are used to calculate the RDM in the network using Equation (16). The calculated RDM in the DT, i.e., RDM_{dt} , is the expected or predicted value of risk in the next time window, which is defined as the duration in which the risk is evaluated. Therefore, the DT can predict the risk of sleeping by utilizing the Markov model and its updated parameters. Then, the estimated risk is compared with the threshold set by the operator. It will also be compared with the next the actual risk, i.e., RDM_a , which will be available at the end of the next time window².

If the predicted risk is higher than the predefined threshold, regardless of RDM_a , the SMs should be deactivated. By this mechanism, operator intent is taken into account. Otherwise, if RDM_{dt} is below the predefined threshold, BS can activate the SMs and benefit from BS sleep mode management algorithms. If the actual risk is higher than the predicted risk, the algorithm might need retraining. If the actual risk is below the predicted risk, the network decides based on the value of RDM_{dt} .

This algorithm is designed and explained in Section IV. Using the risk monitoring procedure in Fig.3, the operators can make sure that their sleep mode management algorithm does not take sleep decisions when there are unexpected user arrivals or is not missing the opportunity of energy saving when it is safe to save energy.

²The time window should be long enough to collect enough samples from the network to calculate the risk and also short enough to avoid delaying large number of users, e.g., for instance in order of seconds.

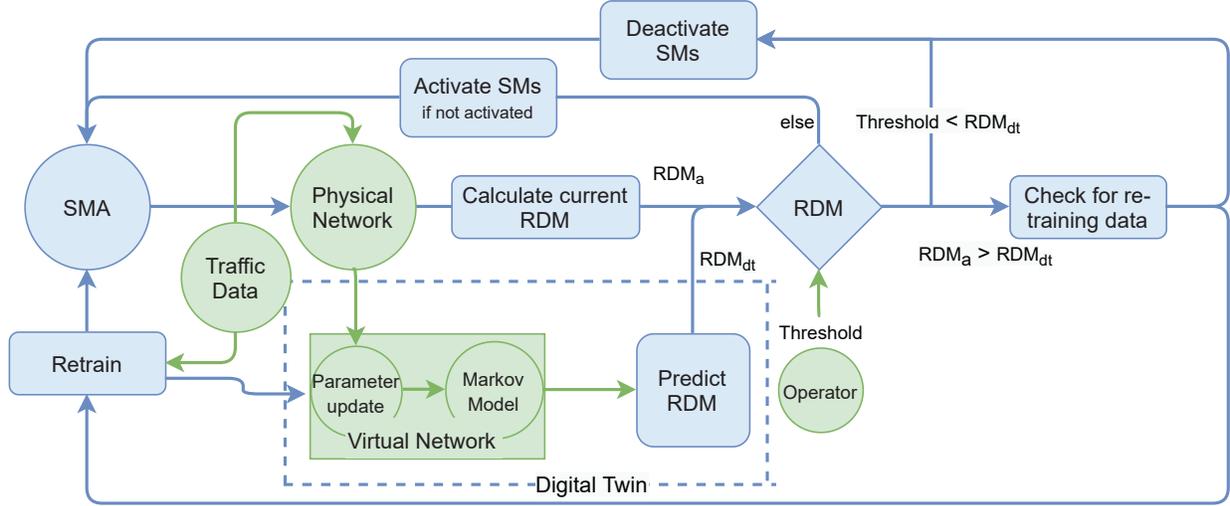


Fig. 3: Digital twin-based risk-aware sleep mode management

B. Digital twin model

In this study, as is depicted in Fig.3, the digital twin has three main parts, 1) parameter update, 2) network model, and 3) prediction. The former two constructs the virtual network, representative of the physical network. The model is continuously updated from real-time data, and uses machine learning and reasoning to help decision-making. The network is modeled as hidden Markov process and in the prediction phase we use the updated virtual network to estimate and predict the future performance metric of the BS sleeping, e.g., sleeping duration, probability of delaying users, and the risk for each state of the network. In the following we explain the Markov model for BS sleep mode management algorithm.

C. Virtual Network Model: A Hidden Markov Process

In the DT, we require a realistic analytical model that can estimate the actual behavior of ML-based BS sleep mode management as a virtual network. In particular, we model the system as a hidden Markov model³ depicted in Fig. 4 where the states, input traffic, and model parameters are obtained by the interaction with the physical network environment. The environment is defined as the load in the system and level of sleeping. The comprehensive description of the environment is provided in Section IV-A1. In the following, we present the hidden Markov model, required information, performance metrics, and a framework for risk monitoring in the system.

The Markov model should encompass all operating states of the BS to be able to predict the performance of the BS sleep modes. The states contain information on 1) sleep mode level that the BS is in, 2) number of users in active mode, and 3) whether there are arrivals or not (ON state or OFF state as depicted in Fig. 1). We model the state transition diagram of the BS as a Markov model illustrated in Fig. 4. For instance, state $S(i, j)$, $i \in \{1, 2, 3\}$ and $j \in \{1, 2\}$, denotes the state

that the BS is in SM_i and the traffic arrivals are in the state j where $j = 1$ means ON state and $j = 2$ means OFF state. If the BS is in state $A_{m,j}$, $m \in \{1, 2, \dots, M\}$ and $j \in \{1, 2\}$, then it is active and m number of users are in the BS. We assume that maximum M number of users can be served and the BS can change its state according to the state transition diagram. When all users in the BS are served, BS goes to one of the SM_i -s and stays there until a new user arrives. Before arrivals of the users, the state changes from $S(i, 2)$, OFF state to $S(i, 1)$, ON state. When a new user arrives, BS goes from state $S(i, 1)$ to state $A(1, 1)$ switching from SM_i to FM; however, it must wait for the wake-up time of SM_i , the duration to activate the sleeping components. When the BS is in any of the sleep modes there is a transition possibility from any of the sleep mode state, i.e., $S(i, j)$ $i \in \{1, 2, 3\}$ to the other SM_i if $j=2$, there is no traffic arrival.

D. Performance Metrics

The proposed Markov model can enable us to derive the steady state probabilities and later the calculation of the risk value. In this section, according to Fig. 4, given the model parameters, and the power model, we derive the sleeping probability, and we define a new parameter quantifying the risk of taking wrong action. Let's define the following probabilities:

- $\nu_{i,j}$ be the probability that BS is in state $S(i, j)$.
- $u_{m,j}$ be the probability that BS is in state $A(m, j)$.

where $S(i, j)$ is the state in which BS is in sleep mode i and the arrival state is j and $A(m, j)$ is the state in which BS is active and there are m number of users being served at the BS and the arrival state is j , and $j = 1$ means ON state and $j = 2$ means OFF state.

1) *Probability of Sleeping*: The probability that the BS is in sleep mode is summation of the probabilities that BS is in SM_i which is given by,

$$V_s = \sum_{i=1}^2 \sum_{j=1}^2 (\nu_{i,j}) \quad (13)$$

³An IPP model plotted in Fig. 1 is a Poisson-emission hidden Markov model with two hidden states [29] in which one of the states has zero emission rate. Therefore, the state diagram depicted in Fig. 4 is a hidden Markov model.

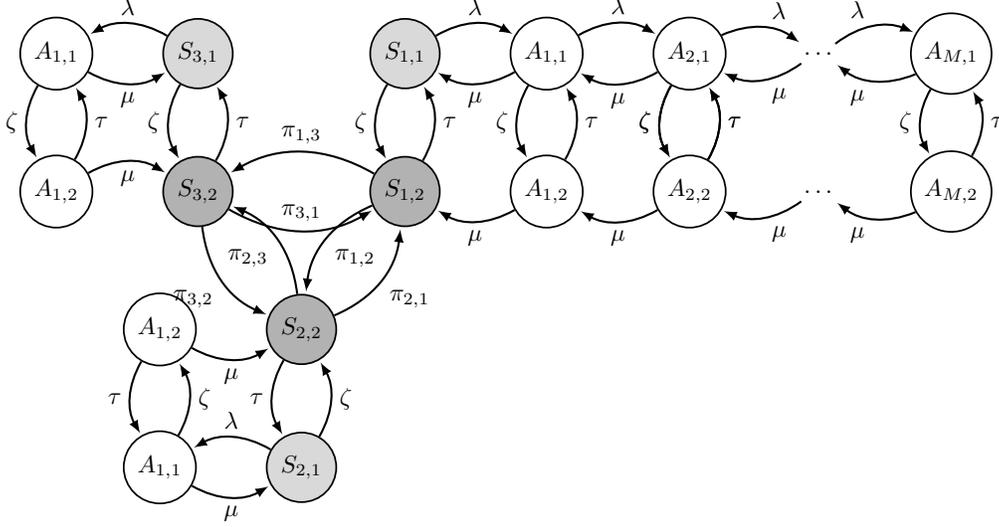


Fig. 4: Markov model for advanced sleep mode management.

In the Appendix A, we have written the balance equations for Markov model in Fig.4. The balance equations are used to calculate the steady state probabilities, i.e., $\nu_{i,j}$ and $u_{m,j}$.

2) *Number of Switching*: Each time BS switches from one sleeping state to another, it adds extra energy consumption or cost to the BS. Therefore, it is crucial to keep track of the number of switchings which is denoted by F_m and is defined as the number of times that BS switches between SMs or activate/deactivates its components. F_m is derived as,

$$F_m = 2\mu \sum_{i=1}^2 u_{1,i} + \sum_{j=1}^3 \nu_{j,2} \left(\sum_{k=1/j}^3 \pi_{k,2} \right). \quad (14)$$

Each time we pass from active to sleep mode in Fig. 4, i.e., *go to SM_i*, later BS has to take a *go to FM* action. The first term in Equation (14) counts the number of switching to SMs and wake-ups. When BS is sleeping and arrivals are in the OFF state, it can switch between SMs. The second term in Equation (14) counts the number of switchings between SMs.

E. Risk of Decision Making

When BS is in SM and new users arrive, the new arrivals must wait until the BS wakes up and goes to the active mode again. Denoted by U_s , the average number of users in the BS when BS is any of SMs is calculated as,

$$U_s = \lambda \nu_{1,1} + \lambda \nu_{2,1} + \lambda \nu_{3,1} \quad (15)$$

where λ is user arrival rate and $\nu_{i,1}$, $i \in \{1, 2, 3\}$ are the probabilities of being in SM_{*i*} while arrivals are in ON state. From the operators' perspective, the smallest possible value for U_s is more preferable since minimum number of users experience delay until they are being served. Although this metric can reflect the favor of the operators, it is incapable of measuring the performance of the BS sleeping algorithms. For instance, when the network is busy, i.e. λ is high, the probability of sleeping will be low. However, one inappropriate sleeping decision may result in a large number of users experiencing delay. On the other hand, in off peak hours,

when λ is low and V_s is high, we may have very few users experiencing delay. Therefore, the parameter U_s cannot solely reflect the higher risk of incurring high delay at peak hours. To tackle this issue, a novel metric called risk of decision making (RDM) is introduced to measure the risk when we take a non-optimal action by taking into account the U_s and the probability of sleeping (in both ON and OFF state). Risk of wrong decision is formulated as,

$$\begin{aligned} RDM &= E(\text{Number of users per time unit} | \text{BS be in SMs}) \\ &= \frac{U_s}{V_s} = \frac{\lambda(\nu_{1,1} + \nu_{2,1} + \nu_{3,1})}{V_s}. \end{aligned} \quad (16)$$

When the network is in peak hours, U_s might be high and V_s is low, hence there is not much room to save energy. Since the network is crowded, any wrong sleeping decision may result in a large number of users experiencing delay, and hence RDM is higher as indicated in Equation (16). In the off peak hours, U_s might be low and V_s is high. Since less number of users are in the network, RDM is lower which is indicated in Equation (16). RDM value can also be used to determine the right time/moments to deactivate the BS sleeping algorithm, in order to avoid delaying large number of users. There are two scenarios in which RDM will be high.

High RDM in Busy Hours: At peak hours, when the network is crowded, it is not beneficial to put BS into sleep modes due to the risk of incurring delay. Therefore, the BS sleeping algorithm can be disabled temporarily until the sleeping becomes beneficial again. It is true that the learning algorithm may learn these moments and avoids sleeping in such events, but it is possible that the algorithm choose the wrong action due to many reasons such as abnormal behavior of traffic, lack of proper training, and change of traffic pattern.

High RDM due to Abnormal Behavior of Traffic: When the behavior of users or the traffic pattern changes, e.g., arrival rate changes, and the algorithm is not previously trained for it, the sleeping decision may yield to unnecessary wake-ups or inappropriate sleepings. It takes some time for the BS to

realize the abnormal behavior of the input traffic which results in unnecessary energy consumption and/or unacceptable incurred delay. When any non-optimal action is taken, it takes some time for the algorithm to converge and update its policy accordingly. Comparing the RDM with a threshold, we can prevent incurring such delays. The threshold value is set by the operator depending on the hours of the day, assured QoS, etc. Using DT, we can estimate the value of RDM. If the experienced value of RDM is above the estimated one, it means the current traffic pattern is different than the expected traffic. Therefore, the BS should disable the sleeping features, i.e., SMs, until either the algorithm is trained over new data set or the traffic behavior becomes normal again.

F. Interaction with DT and Physical Network

After defining the DT, underlying Markov model, and the performance metrics, we now explain how the parameters of the hidden Markov model within the DT can be obtained and updated. The parameters of hidden Markov model is function of user arrival patterns, rates, ON state duration, and OFF state duration. Therefore, these parameters can change during the day. In order to obtain the parameters of the model, i.e., λ , τ , ζ , from the training data set, we can use a well-known backward-forward algorithm, i.e., Baum–Welch algorithm [31]. The Baum–Welch algorithm finds the maximum likelihood estimate of the parameters of a hidden Markov model given a set of observed input sequences. This algorithm computes the statistics of the input traffic sequence $O = \{o_1, o_2, \dots, o_T\}$, and then updates the maximum-likelihood estimate of the model parameters, i.e., ON/OFF state duration and arrival rate. The procedure of Baum–Welch algorithm to extract the IPP parameters is defined in [29].

When BS is in the active mode, BS's transmit power P_t is adapted to match the traffic load. Assume that BS is capable of serving x bit/sec and resources are equally shared between all served users, using a fair scheduler [32]. The user departure rate is $\mu = x/l$ where l is the user's requested file size and $x = B \log(1 + hP_t)$, and $h = (g/N_0B)$ where g represents the channel gain, B is the bandwidth, and N_0 denotes the noise density [8].

Another set of parameters, i.e., $\pi_{i,j}$ -s $i \in \{1, 2, 3\}$ and $j \in \{1, 2\}$, are the transition rates between SMs while BS is in OFF state. Meanwhile, this transition rates represent the learned SM management policy. The initial transition rates can be calculated during the training phase. However, during the test phase these rates can be updated by keeping the information about the average number of times BS switches from one SM to another SM in a given time duration. The policy can be fed back to hidden Markov model to update the predicted BS sleeping performance metrics, e.g., risk of decision making.

IV. PROPOSED SLEEP MODE MANAGEMENT ALGORITHM

A. Reinforcement Learning Elements

1) *Action set and Environment*: The state of the system at time index i , is denoted as s_i . This state is composed of current traffic and encoded information about previous traffic which is fed back by a recurrent neural network module, explained

in Section IV-C. The current traffic is the number of utilized PRBs at time index i . Each action has a transition period, and the system is said to be in the last state until the new state is in action. The action set denoted by \mathcal{A} is defined as $\mathcal{A} = \{\text{FM}, \text{SM2}, \text{SM3}\}$. Here, FM and SM k , with $k \in \{2, 3\}$, denote the go-to-fast-mode and go-to-SM k mode, respectively.

The system environment information at time index i is fully captured via the following 5-tuple:

$$e_i = \left(p_i^{(b)}, l_i^{(b)}, p_i^{(c)}, l_i^{(c)}, l_i \right) \quad (17)$$

where $p_i^{(b)}$ and $p_i^{(c)}$ are the basic and capacity BS power consumption, $l_i^{(b)}$ and $l_i^{(c)}$ represent the basic and capacity BS accommodated load in RBs, while l_i is the total traffic load that needs to be served at the i^{th} time index, respectively. That is $l_i = l_i^{(b)} + l_i^{(c)} + d_i$ where d_i represents the number of delayed RBs at the i^{th} time index, which is the difference between the arriving traffic and total served traffic. We also define the total power consumption at the i^{th} time index as $p_i = p_i^{(b)} + p_i^{(c)}$. The users arrive to the network as is explained in Section II-C3. At each time index, the amount of load in the environment is determined by the number of users that are in network in the same time index. The value of risk in the environment can be calculated based on the current load, service time and other estimated model parameters.

2) *Power Saving Reward*: When the CapBS is operating in fast mode (FM), no serving delay occurs. However, the system may miss the opportunity of saving more energy if deeper SM is possible. We define a normalized energy saving metric as

$$\tilde{r}_{i,p} = \frac{(p_{SM1} - p_i)^+}{p_{SM1} - p_{SM3}}, \quad (18)$$

where $(x)^+$ is a operator that takes the value of x if it is positive and is zero otherwise, $p_i, i \in \mathcal{A}$, is the power consumption of SMs and $r_{i,p} \in [0, 1]$.

3) *Delay Penalty/Reward*: When the system receives a request while it is in one of the deep SMs, i.e., $d_i \neq 0$, or in the case when $d_i = l_i^{(c)} = 0$, the normalized delaying penalty can be calculated as

$$r_{i,d} = -\frac{d_i}{l_{\max}^{(c)}}, \quad (19)$$

which takes 0 when $d_i = 0$, and 1 when $d_i = l_{\max}^{(c)}$, respectively. When the system is in FM with $l_i^{(c)} \neq 0$, although no extra power saving is attained, there should exist a reward for the avoided delay, i.e., $(d_i = 0) \& (l_i^{(c)} \neq 0)$, the reward is

$$r_{i,d} = \frac{l_i^{(c)}}{l_{\max}^{(c)}}. \quad (20)$$

4) *Total Reward*: The incurred reward function at the end of the i^{th} time index due to the action taken at the end of time index $i - 1$ can now be plausibly defined as

$$r_i = (1 - \alpha)r_{i,p} + \alpha r_{i,d}, \quad (21)$$

where α is a weight parameter to prioritize power saving or serving delay. When an action is taken, the system will be frozen for the minimum duration of that action. Therefore,

the incurred reward is calculated at the end of such duration. Let w_{a_i} be the transition duration of a_{i-1} , as defined in Table I, the total reward at time index i is calculated as

$$\bar{r}_i = \frac{\sum_{j=i-w_k-1}^i ((1-\alpha)r_{j,p} + \alpha r_{j,d})}{w_k}. \quad (22)$$

We assume that when the BS is activated, it stays in FM mode for 14 symbol times to avoid ping-pong effect [33]. Since every switching from one state to another consumes energy [34], we add a negative reward to the total reward if change of action happens, i.e.,

$$\bar{r}_i^* = \begin{cases} \bar{r}_i & \text{if } a_{i-1} = a_i \\ \bar{r}_i - \frac{1}{w_k} & \text{else} \end{cases} \quad (23)$$

B. From Q-learning to Deep Q-learning

Reinforcement learning is a class of machine learning solutions which learns from interactions to achieve a certain goal. The learner, which is called the agent, at each time step (for discrete time) interacts with the environment by taking an action and observes the reward of taken action. Based on the received reward, the agent updates a table, known as Q table, to keep track of values of each action for all states. Q-learning is an RL algorithm that calculate this value using,

$$Q^{New}(s_i, a_i) \leftarrow Q^{Old}(s_i, a_i) + \eta [R_{i+1} + \gamma \max_a Q(s_{i+1}, a_{i+1}) - Q(s_i, a_i)] \quad (24)$$

where $R_{i+1} = \bar{r}_{i+1}$ is the instantaneous reward at time index i , (s_i, a_i) is current state-action pair, $\eta \in (0, 1]$ is learning rate, and $\gamma \in (0, 1]$ is discount factor.

The computational requirements of Q-learning grow exponentially with the number of states and actions. The problem with large state spaces, which is called the curse of dimensionality, stems from 1) the memory needed for large tables, 2) the time and data needed to fill them accurately 3) many encountered states will be new, will have never been seen before. Therefore, it is not feasible to find the optimal and exact value of Q-function. Alternatively, it is possible to approximate the Q-function with limited computational resources. The function approximation is an instance of supervised learning which can be combined with RL methods to deal with large state space sizes. It helps generalizing learnings from a limited set of past states with a successful approximation over a larger state space. In this study, instead of using Equation (24), we approximate the Q-function using a special type of recurrent neural networks i.e. Long Short Term Memory (LSTM) explained in the following section.

C. Deep Q-learning with LSTM

When the problem has a large state set, it is not possible to visit all states. In this partially observed environment, the RL agent needs to encode the information about the current input and state-action trajectory. The most common way of doing this is to use recurrent neural networks (RNNs) which are powerful models for processing sequential data such as time series data. In RNNs, there is a feedback loop which

makes the output as part of the input for next time step. With this approach, the information about states can propagate over time. However, conventional RNNs fundamentally cannot learn long-term dependencies between data sequences. In recent years, LSTM architectures have been gaining popularity due to their ability to model long-term dependencies. An LSTM architecture includes input gate, forget gate, cell state, and output gate. Input gate updates the cell state which can save information of the current state. Forget gate decides whether this information should be kept or discarded, and the output gate decides the next state. With these gates, networks are able to dynamically learn the structure of longer input sequences and effectively associate memories to adapt which parts of sequence to remember and which parts to forget for the task at hand. Interested readers can refer to [35] for in-detailed explanation of LSTM architecture.

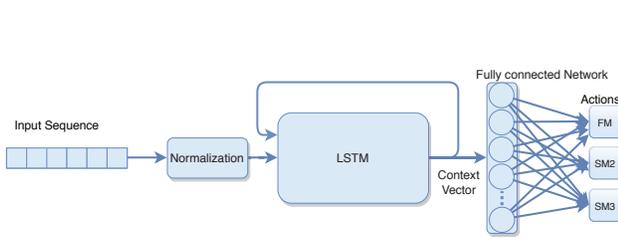
Fig. 5 illustrates the structure of the proposed Deep learning algorithm and the training procedure. According to Fig. 5a, the RL agent gets the normalized input sequence and the context vector, i.e., the mobile traffic data and a hidden vector (generated in previous time index), as inputs to the first layer of the LSTM module. This module, encodes the current state into a context vector. This context vector describes the state of the problem in a high dimensional space. The dimension is a hyper-parameter model and the value is chosen depending on the complexity of the data and problem. This vector is fed to a fully connected network which connects the context vector to the actions and estimates the value of Q-function for each action.

To train our DQN model, we use experience replay method [36] with two incentives. The first reason is the intrinsic characteristics of sleep mode management problems, such as diverse input traffic patterns. Hence, the proposed method should be online and learn directly from the environment. Secondly, this method removes the existing correlation in the traffic data sequence, therefore, it makes the training phases more stable [36]. As depicted in Fig. 5b, in this method, at each time step the agent interacts with the Environment and performs its current actions, then stores the tuple, $(s_i, a_i, R_{i+1}, s_{i+1})$, in the replay memory of length M . We randomly select N batches from the replay memory where each batch contains B consecutive samples. Therefore, we have N uncorrelated batches and a total number of $N \times B$ samples. We use this batches to train and update the weights of the DQN in the opposite direction of its gradient. By random sampling, we make sure that the training set is composed of enough span of samples and our model can generalize well with inexperienced states.

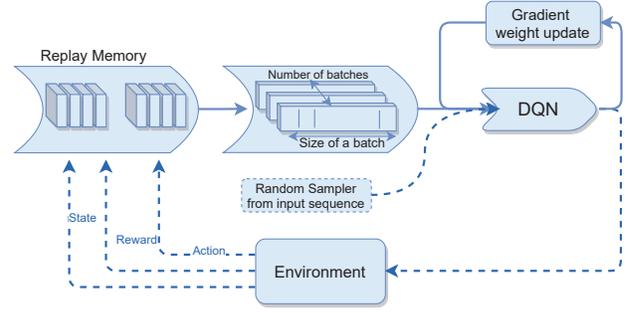
Since the rewards are normalized, we can use a multi-level binary cross entropy loss function defined as,

$$\mathcal{L}(q, R) = \frac{1}{B} \sum_{b=1}^B \ell_b, \quad (25)$$

where q is estimated Q-function, R is the reward, B is the



(a) Deep Q-learning non-linear architecture



(b) Experience replay training method

Fig. 5: Structure of proposed DQN Sleep Mode Management Algorithm.

batch size, and ℓ_b is,

$$\ell_b = \frac{1}{|\mathcal{A}|} \sum_{a=1}^{|\mathcal{A}|} \zeta_a \ell_{b,a}, \quad (26)$$

where ζ_a is a rescaling weight given to the loss of each action for loss balancing⁴. We set the ζ_a to the ratio of non-zero values of the data set to the size of the data set. $\ell_{b,a}$ is given by,

$$\ell_{b,a} = -\nu_b [R_a \cdot \log \sigma(q_{b,a}) + (1 - R_a) \cdot \log(1 - \sigma(q_{b,a}))] \quad (27)$$

where ν_b is a rescaling weight given to the loss of each batch element which is constant in our study. $\sigma(\cdot)$ is a Sigmoid function to map its input to a value between 0 and 1. R_a and $q_{a,b}$ are rewards due to action a and estimated Q-function for given action a and batch element b . In this study, we use Adam optimizer [37] to update the weight parameters to minimize the loss function. Having estimated Q-function in hand, the optimal policy to take an action is

$$\pi^*(s) = \arg \max_a Q(s_i, a_i | \mathbf{v}), \quad (28)$$

where \mathbf{v} is the context vector calculated by the LSTM layer.

We summarized the procedures of training and testing phases of SMA in the Algorithm 1. In the Algorithm 1, we explain the training and test phases of the SMA. We assume that the DQN model is given. For the training phase, we randomly take B values from the data set and feed the data to the DQN model. We get the action as output of the DQN and compute the loss for this action. Then, by using Adam optimizer we minimize the loss and update the DQN model.

For the test phase of the Algorithm 1, first, the stream of input traffic is fed to the DQN model and DQN estimates the value of the Q-function for each action. Based on these values, the agent selects the best action and passes the policy to the BS.

It is worth mentioning that the sleep mode management algorithm is used only when RDM is below the predefined threshold. In other words, sleep mode management algorithm

⁴When the number of zeros in data is much higher than the number of other values, the loss due to the more frequent values contributes much more to the total loss. Therefore, the network learns to minimize the loss corresponding to the more frequent values to minimize the loss over the whole data set. However, in our setup, the loss due to non-zero values is more important to be minimized. Therefore, rescaling the loss is needed to avoid this issue.

Algorithm 1: DQN sleep mode management algorithm

Input: Traffic data \mathbf{X} , batch size B , training data length L , prediction length P_L , number of epochs N_P
Output: For retraining: Training data set, trained DQN model. For SM management: Predicted best action.
Initialize: DQN Model: LSTM layers, Fully connected network;
 Normalize traffic data \mathbf{X} ;
Train:
 for $Epoch = \leftarrow 1 : N_P \times N$ do
 • To get batch:
 – Randomly select B values from $[0, |\mathbf{X}| - (L + P_L)]$, to form b_i -s
 – Construct batch data $\mathbf{d}_i = \mathbf{X}[b_i : b_i + L - 1]$ and $\mathbf{f}_i = \mathbf{X}[b_i + L : b_i + L + P_L - 1]$ to form $x_i = (\mathbf{d}_i, \mathbf{f}_i)$
 • Forward batch x_i -s to DQN model M
 • Get a_i from DQN
 • Compute loss based on Eq. (25)-(27).
 • Update DQN model using Adam optimizer
Test:
 Get SM active from Algorithm 2;
 while SM active do
 • Get data stream from live network
 • Estimate Q-function using DQN
 • Choose best action a using policy:
 $\pi^*(s) = \arg \max_a Q(s_{i+1}, a_{i+1} | \mathbf{v})$
 • Perform action a
 • pass policy and reward to BS;

will be in use according to the Algorithm 2. This algorithm explains the procedure of RDM monitoring in the network. According to this algorithm, the network information such as, average arrival rates and transition probabilities, is given to the BS. The BS uses this information and updates the model parameters and predicts the RDM. The value of RDM is used to check if the risk is above or below the threshold. If the value of RDM is above threshold, the SMs are deactivated. If the value of RDM is below the threshold the SMs can be activated. It is important to keep track of the RDM value for a while, let's say T_w , to make sure that the decisions are not made based on instantaneous changes in the network. To avoid spontaneous decisions, moving average operator can be used to cancel out the instantaneous changes in data traffic and RDM. Therefore, the sleep mode management algorithm comes back to loop only if the average value of RDM is low and it is safe

Algorithm 2: Risk management algorithm in the network

Input: Traffic data, transition rates between SMs
Output: Decide to run which part of Algorithm 1 or not run Algorithm 1.
Condition=True;
while Condition **do**
 Network status update:
 • Update the parameters of hidden Markov model.
 • Calculate the estimated RDM, i.e., RDM_{dt} .
 • Calculate the actual RDM in the network, i.e., RDM_a ;
 if RDM_a or RDM_{dt} is higher than threshold **then**
 | Deactivate the SMs;
 else if $RDM_a > RDM_{dt}$ **then**
 | • Deactivate the SMs;
 | • Run Algorithm 1: Train;
 else
 | • Monitor RDM_a for duration T_w ;
 | • if $RDM_a <$ threshold: Activate the SMs.
 if SM active **then**
 | Run Algorithm 1: Test.

TABLE II: Studied System Configuration Parameters

Parameter	Value
Modulation	16-QAM (4 bits per RE), $M = 16$
Antennas	1 per cell sector, one sector per cell
Bandwidth	20 MHz for CovBS, CapBS
RBs per band	100 per 20 MHz (LTE-A Compatible)
REs per TTI	14 (Based on FDD Frame Structure 1)
Symbol Time	$T_s = \frac{1}{14 \times 1000}$ sec $\approx 72 \mu\text{sec}$
Hidden-dimension/ Batch size	256, 200
Episode per epoch	2048
τ, ζ	0.1, 0.5
Risk Threshold	1.2

to save energy.

V. SIMULATION RESULTS

An event-based numerical simulation is implemented in Python, running the proposed adaptive algorithm and its baseline counterparts. The simulation parameters are defined in Table II.

1) *Baselines:* We compare the performance of the SMA algorithm with optimal BS sleeping (OBS) where future traffic arrivals are assumed to be known as an upper bound on energy saving gain when no delay is incurred to the users. We follow the procedure explained in Algorithm 3 [23].

In OBS, the BS is fully aware of all future incoming users and their activity times. Scanning all inactivity periods, BS can fill the inactive time with the proper SM. After the initialization step, OBS fills the inactive time with the deepest possible SM, i.e., SM3. Then, the remaining inactive time will be filled by the next deep SM, i.e., SM2. The remaining inactive time will be filled by SM1. The OBS algorithm, keeps track of the number of actions for each SMs and the energy performance of the BS.

This procedure guarantees the highest energy saving gain with no incurred delay to users, yet utilizing the future

Algorithm 3: Optimal BS Sleeping

Input:
Cell load $l^{(c)} = \max\{0, l - l_{\max}^{(b)}\}$ for every symbol time over the entire time horizon. SM counter vector. Duration of SMs in terms of symbol time, i.e.,
 $w = (w_1, w_2, w_3) = (1, 14, 140)$.
Output: Energy performance, total SM counter.
Initialization:
 • Set $w = (w_1, w_2, w_3) = (1, 14, 140)$.
 • Set total SM counter, TSMC= $[0, 0, 0]$.
 • z : A vector containing the run lengths of zero loads (in symbol times).
 • e : A vector with the same length as z containing the energy consumption of the each duration in z .
 • $j = 0$, an index that runs over entries of z .
SMs Fitting:
Set: $n[j] = z[j]$
for $j = 1 : 1 : \text{length}(z)$ **do**
 SMC = 0
 for $k = 3 : -1 : 1$ **do**
 while $n[j] \geq w_k$ **do**
 SMC $[k] + = 1$,
 TSMC $[k] + = 1$,
 $n[j] - = w_k$
 Energy Evaluation :
 Update the energy consumption, e using Equations (1) and (2).

knowledge⁵. We compare the performance of SMA with fixed sleep mode where only SM1, i.e., the shallowest SM, is used. SM1 is already implemented in the BSs and can be activated without requiring any intelligence. In this scheme, BS can go to SM1 when no user is served. In these baselines all users are served instantaneously with no incurred delay. We also compare the results with Q-learning algorithm in [20].

2) *LSTM Hyper-parameters:* Choosing the right hyper-parameters is crucial to design a network, including the deep learning model with the best performance. Through simulations, we swipe over different parameters and compare the performance of the whole network, i.e., loss and accuracy, which is summarized in Table III. The loss function is defined in (25) and to calculate the accuracy we compare the decision made by the optimizer with the correct decision to make. We pick the parameter values with the best performance for the rest of simulations. The best performance is achieved with a two-layer LSTM and hidden vector of size 50 and batch size of 200.

3) *Data Generation:* We use mobile traffic data provided by a Swedish operator. As explained in Section II-C, data sets needs to be processed to be used in our simulation tools. In Fig. 6a, we plot the normalized average load per 5 minutes for original data set and the generated load for the same 5 minutes. The comparison shows that the generated data preserves the average daily behavior of the original data. The fluctuations are due to the randomness of arrivals and fluctuations in the data set. We can see that the generated load in 5 minutes, which is the finest available time granularity, is following the

⁵Since the output of OBS depends on the future inputs and not just past and current inputs, the method is non-causal.

TABLE III: Hyperparameter values of the network

LSTM Performance		LSTM Parameters		
accuracy	loss	layers	hidden size	seq length
0.96940	0.08858	1	50	100
0.96972	0.08706	1	100	100
0.96924	0.08882	1	150	100
0.97964	0.07855	2	50	100
0.96000	0.09223	2	100	100
0.96068	0.09214	2	150	100
0.96024	0.09438	3	50	100
0.95708	0.09212	3	100	100
0.96712	0.08960	3	150	100

trend in the data set. In Fig. 6b we show the arrivals in time scale of 5 seconds.

According to Section II-C4, the parameters of IPP model, i.e., τ , ζ , λ , as well as the statistical parameters of distributions of request per arrivals must be set in a way that mean and variance of the generated load matches with the mean and variance of the original data set, for the same duration. However, according to Equation (11), for a given τ and ζ , there might be no solution for λ and $E(\psi)$. If the ratio of $\frac{\text{var}(\Psi)}{E(\Psi)^2}$, i.e., dispersion ratio of data set, is lower, the feasible region for λ is wider. In Fig. 6c, we illustrated the feasible region of data generation parameters for the mean value of 100 Kbps and variance of 5 Kbps. As can be seen from this figure, when τ is larger (or ζ is smaller), i.e., duration of ON state is higher, the feasible region for λ is larger. This figure also illustrates that if the combination of τ , ζ , and λ , lies in the infeasible region, the data generation is not reliable. Moreover, during the test phase if similar situation happens, one cannot rely on the prediction of the parameters τ , ζ , λ , and the associated risk.

4) *Learning performance*: In Fig.7, we present the learning curve of the proposed DQN algorithm. We compared the DQN selected action with the pre-calculated best action based on the chosen alpha to decide whether the action was correct within an episode. In this figure, each episode contains a sequence of states, actions, and rewards with a sequence length of 100. The solid blue line shows the average ratio of correctly chosen actions for 100 simulation runs of each episode. The dark-shaded and light-shaded areas show the one and three standard deviations of the accuracy for the given episode. A high standard deviation or larger shaded area indicates the uncertainty in the model to select the actions and hence may take different actions at a specific episode. A low standard deviation indicates that the algorithm is more determined and tends to choose similar actions at a particular episode. In the early episodes, the model learns the best actions by exploration and choosing various actions and, hence, high standard deviation. After 100 episodes, the DQN model is adequately trained in this setup.

5) *Energy saving vs incurred delay*: As per Equation (21), the parameter α is a tuning parameter adjusting the weight given to power saving vs minimization of delay, in the reward function. In Fig. 8, we present the performance of the SMA with regards to the α . The higher α is, the more weight is put to minimize delay and less weight is assigned to energy saving. Hence, with higher α less energy is saved. We compare

the performance of SMA with optimal BS sleeping (OBS) algorithm defined in Algorithm 3, as an upper bound on energy saving. OBS is always optimal and non-causal due to the knowledge of future information. For $\alpha = 0$, where the total reward includes only energy saving reward, SMA can achieve the optimal energy saving but at cost of incurring delay to the users. When latency is more prioritized, the BS becomes more conservative to choose deeper SMs and prefers shallower SMs. Therefore, less energy is saved in favor of less incurred delay. SMA outperforms the fixed SM where only SM1 is in use and BS cannot benefit from deeper SMs. Moreover, SMA performs better than Q-learning algorithm [20] because SMA can find better long-short term dependencies in traffic data and hence leverage this information to find a better energy saving policy.

6) *Daily energy saving*: Fig. 9 shows the energy saving percentage of the SMA, only SM1 and OBS over a day. The values are calculated with regards to the energy consumption of the BS without using any SM and using the energy saving algorithms, i.e., OBS, Only SM1, and SMA. Fig. 9 illustrates that between 2:00-6:00 AM are the best hours of a day with the highest energy saving potentials. The least opportunity for energy saving happens in the evening at about 18:00. It is worth mentioning the data set under consideration is for an area which is dominated by the industrial buildings. The traffic pattern and hence the energy saving patterns depend on the type of area. For instance, in residential area, there might be a peak demand in the evening when people are at home while during the working hours network may experience low load within this area. During 2:00-6:00 AM, where the most energy saving is attained, SMA achieves considerable energy saving very close to the OBS. Slightly higher energy saving is achieved compared to OBS due to the delay tolerance allowing user arrivals when a BS is in deep sleep. OBS wakes up and does not allow any arrivals when BS is in sleep. During the peak hour, since there is a risk of delaying more number of users, SMA avoids going to deep sleeps and hence less energy saving is achieved compared to the OBS. When only SM1 is activated, at very low load, the opportunity of saving energy is missed because we only use the shallowest SM. At very high load, there is less opportunity to activate longer and deeper SMs leading to similar performance between SMA and only-SM1 to avoid causing performance degradation. In this figure, we plot the ratio of delayed users compared to the total number of users in each hour. Between 2:00-6:00 AM, when the traffic is too low, the ratio of delayed users is higher because of two reasons, 1) the traffic load is low; therefore, even delaying one user increases the ratio significantly. 2) SMA achieves very high rewards due to high energy saving and choosing deeper SMs; hence, it accepts a small number of delaying users. Similar reasoning is applicable to other hours. It is worth noting that Fig. 9 is very dependent on the value of α . For instance, a small value of α scales up the ratio of delayed users in favor of saving more energy (please see Fig. 8).

7) *Digital twin performance evaluation*: In Section III-C, we present the DT model that can characterize the behavior of the sleep mode management algorithm. In order to calculate the probability of sleeping, we need to have information

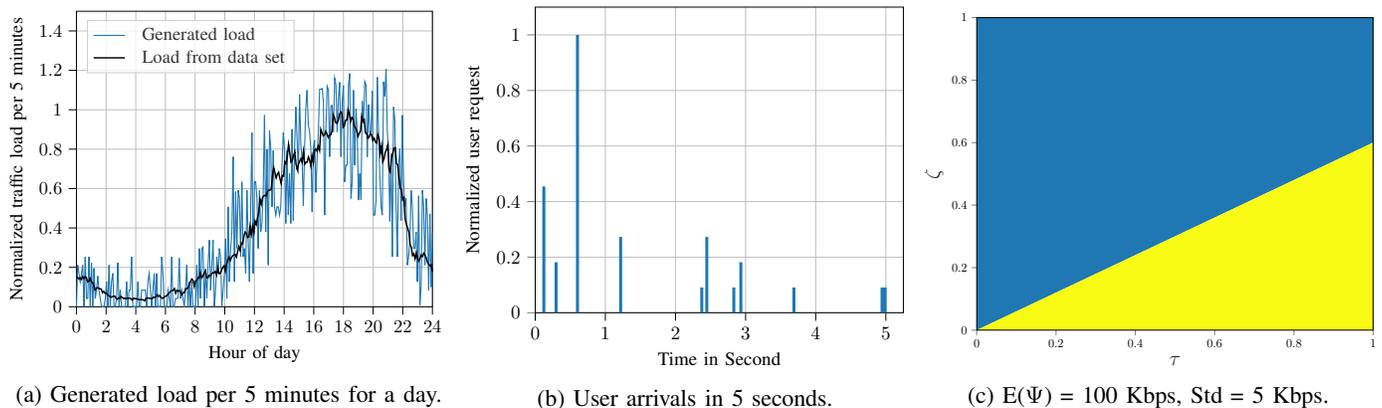


Fig. 6: Network load generation as input to the network for $\tau = 0.1$ and $\zeta = 0.5$. Feasible region for λ for $E(\Psi) = 100$ Kbps. Yellow shows the feasible and blue represents the infeasible region.

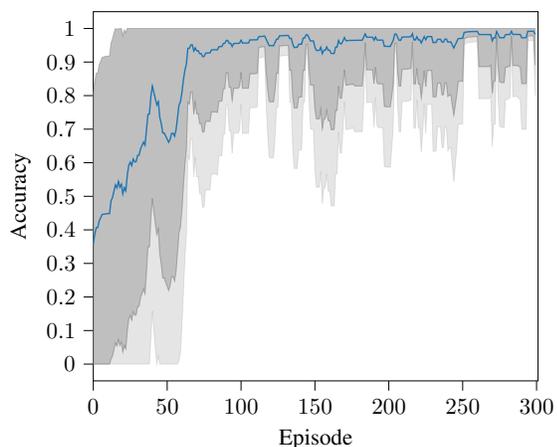


Fig. 7: Learning curve the DQN model for $\alpha = 0.7$. Dark shaded region depicts the one standard deviation and the light shaded are depicts the 3 standard deviations variations.

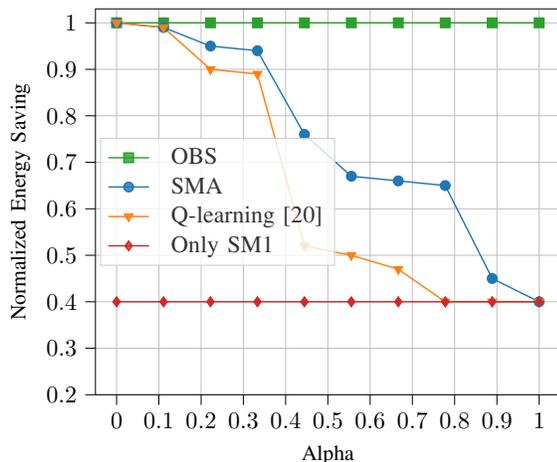


Fig. 8: Energy saving as a function of tuning parameter α in the reward function.

regarding the transition rates. This information can be obtained from the trained data, e.g., transition rates between sleep modes, and be updated in the physical network, e.g., arrival

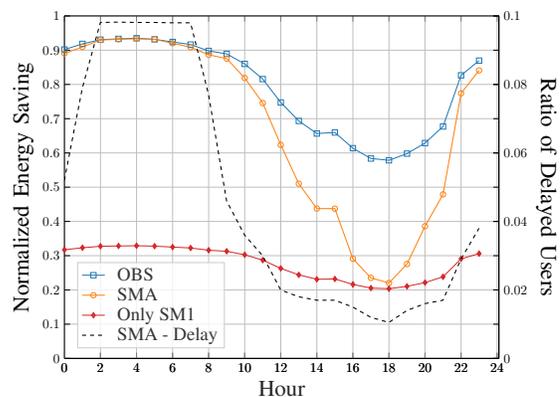


Fig. 9: Daily profile of energy saving and ratio of delayed users.

rate. In Fig. 10, we compare the probability of sleeping derived in Equation (13) with the statistics collected from the implementation of the SMA. When running the SMA, for a given arrival rate, we measure the activity and inactivity time of the BSs. Assuming that during the inactivity time of the BS one of the SMs will be used⁶, we calculate the probability of sleeping by dividing the inactivity duration by the total duration, i.e., activity plus inactivity time. It can be seen from this figure that the DT can properly model the behavior of the SMA for different arrival rates.

In Fig. 11, we illustrate the breakdown of probability of being in each sleep mode. In Fig. 11a, we compare the probability of sleeping and being active from the analysis, i.e., Markov model, and simulation, i.e., SMA. In Fig. 11b, utilizing the hidden Markov model, we illustrated the breakdown of the probability of being in each of the possible operating state of BSs, i.e., three sleep modes and being active. Assuming that the user arrival random process is wide sense stationary, the probability of being in each state corresponds to the fraction of the time BS is in each of the operating states. This information can be used to determine the most visited state by the BS for further energy saving improvement.

⁶We also assume that the random process of the user arrival is wide sense stationary.

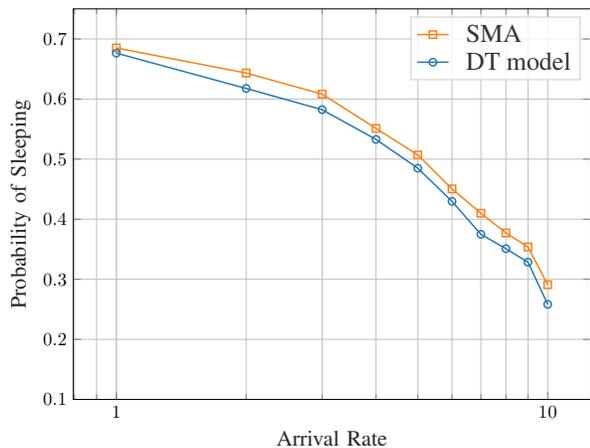


Fig. 10: Probability of sleeping in SMA vs estimated probability of sleeping via hidden Markov model.

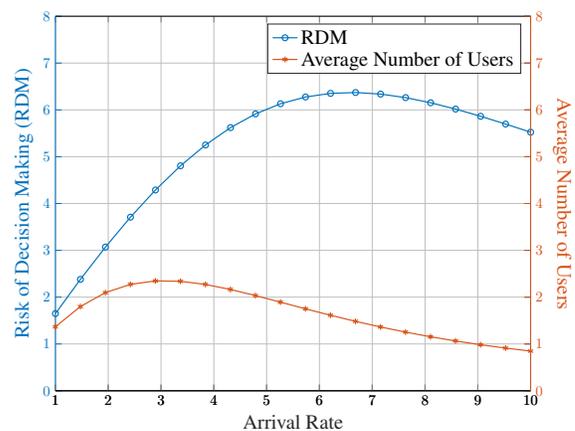
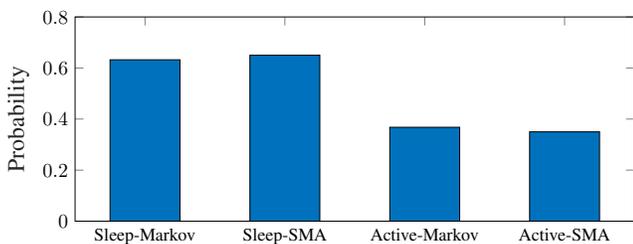
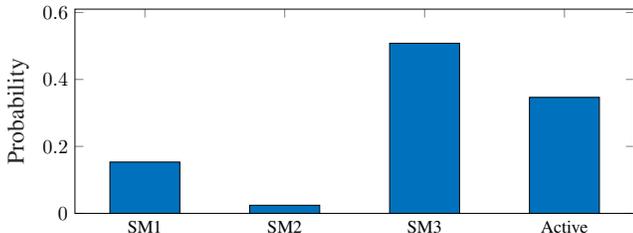


Fig. 12: Comparison of RDM and the average number of users experiencing the serving delay.



(a) Comparison of hidden Markov model estimation with SMA in terms of sleeping probability.



(b) Breakdown of probability of sleeping for each sleep modes using DT.

Fig. 11: Comparison of hidden Markov model estimation and SMA performance.

8) *RDM performance*: According to the discussion in Section III-E, the expected number of users cannot represent the risk of decision making. In Fig. 12, we show that the RDM metric defined in Equation (16) is capable of reflecting the risk of decision making in BS sleeping. When the arrival rate is high, i.e., BS is under the high load, BS is active for most of the times and less number of users should wait for BS to wake up and serve them. On the other hand, if the arrival rate is low, fewer number of users are active in the network and hence less number of users might wait until they are served. In both cases, the number of delayed users due to BS sleeping are small, however, the risk of BS sleeping is higher in the latter case. According to Fig. 12, unlike the average number of users metric, the RDM metric can perfectly differentiate between these two cases and show higher risk at the higher arrival rates.

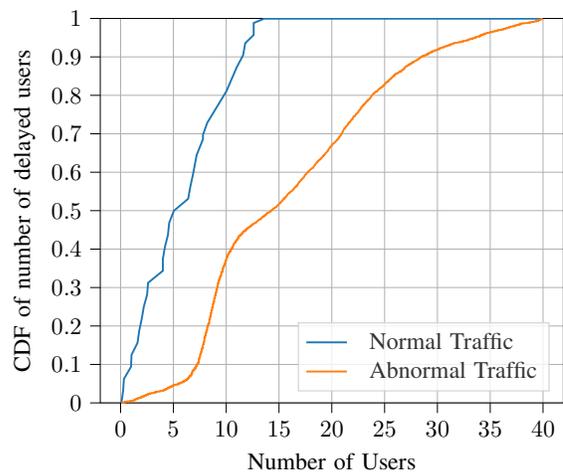


Fig. 13: CDF of delayed users

9) *Abnormal Traffic Behavior*: Despite the appealing performance of ML-based algorithms, their performance under abnormal situations, such as unseen abrupt changes, must be investigated. In Fig. 13, we illustrate the cumulative distribution function (CDF) of the number of users who have been delayed due to BS sleeping in duration of one hour for two types of traffic, 1) normal traffic, and 2) abnormal traffic. Normal traffic is a traffic that the algorithm is trained over same data set. The abnormal traffic is an input traffic which is generated from a different data set. According to Fig. 13, with 80 % probability, less than 10 requests are delayed in the observation period while in the abnormal traffic with probability of 50 % more than 15 users are delayed during the same observation period. According to this figure, with very high probability, large number of users are being delayed due to BS sleeping under abnormal traffic circumstances. Therefore, there is a need to detect the abnormal situations and avoid the risk of delaying large number of users.

In Fig. 14, we plot the performance of Algorithm 2. Using the procedure in Algorithm 2, BS calculates the risk value and when the risk is high it temporarily deactivates SMs. When SMs are deactivated, BS calculates the potential risks and

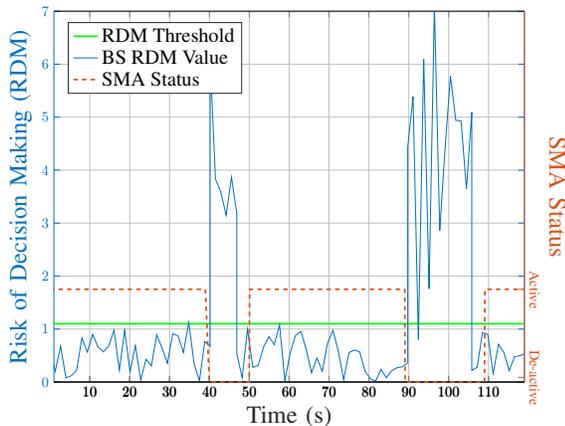


Fig. 14: Performance of risk management algorithm

when the risk is small again for a while, i.e., for duration of T , it activates the SMs again. The performance of the risk management algorithm is dependent on the RDM threshold which triggers the (de/re)activation of the SMA. When RDM threshold is high, BS accepts more risk and hence more number of users can be delayed. Hence SMA will be active for a longer duration. On the other hand, very low value of RDM threshold makes the BS very sensitive to delaying the users and hence make less use of SMA to save energy. This threshold can be set by the operator depending on the time and location of the BS.

VI. CONCLUSION

In this paper, we proposed a deep Q-learning algorithm to reduce the energy consumption of 5G base stations using multi-level sleep modes (SMs). In particular, we considered 3 SMs each with distinct (de)activation time and energy consumption. Since there is always a risk associated to machine learning algorithms, we proposed a framework to predict and manage the associated risk. First, we defined a novel metric to quantify the risk of decision making (RDM). We demonstrated that the proposed metric can represent the risk of taking non-optimal decisions. Then, we proposed a digital twin model (DT) that can predict the performance of the proposed DQN algorithm by estimating the expected probability of sleeping and the expected RDM. Finally, thanks to the proposed DT, BS can detect the abnormal behavior of traffic data and deactivate the SMs to avoid any extra performance degradation. We evaluated the performance of proposed algorithm using real network data obtained from a BS in Stockholm and compared it with the baselines. Simulation results confirmed that with the proposed DT, we can estimate the performance of the proposed SM management algorithm. With the help of this model, we can avoid incurring large delays to the users due to abnormal behavior of input traffic. Moreover, the simulation results showed that considerable energy saving can be achieved with a good compromise with the serving delay, considering that number of users that will be delayed can be controlled thanks to the DT assisted learning mechanism.

APPENDIX

In this section, we write the balance equation for the Markov process depicted in Fig. 4.

$$\nu_{1,2}(\tau + \pi_{1,3} + \pi_{1,2}) = \nu_{2,2}\pi_{2,1} + \nu_{3,2}\pi_{3,1} + \zeta\nu_{1,1} + u_{1,2}\mu_{1,2} \quad (29)$$

$$\nu_{2,2}(\tau + \pi_{2,3} + \pi_{2,1}) = \nu_{1,2}\pi_{1,2} + \nu_{3,2}\pi_{3,2} + \zeta\nu_{2,1} + u_{1,2}\mu_{2,2} \quad (30)$$

$$\nu_{3,2}(\tau + \pi_{3,1} + \pi_{3,2}) = \nu_{1,2}\pi_{1,3} + \nu_{2,2}\pi_{2,3} + \zeta\nu_{3,1} + u_{1,2}\mu_{3,2} \quad (31)$$

$$\nu_{1,1}(\lambda_1 + \zeta) = u_{1,1}\mu_{1,1} + \tau\nu_{1,2} \quad (32)$$

$$\nu_{2,1}(\lambda_2 + \zeta) = u_{1,1}\mu_{2,1} + \tau\nu_{2,2} \quad (33)$$

$$\nu_{3,1}(\lambda_3 + \zeta) = u_{1,1}\mu_{3,1} + \tau\nu_{3,2} \quad (34)$$

$$u_{1,1}(\zeta + \mu_{1,1} + \mu_{2,1} + \mu_{3,1} + \lambda) = \tau u_{1,2} + \lambda(\nu_{1,1} + \nu_{2,1} + \nu_{3,1}) + \mu u_{2,1} \quad (35)$$

$$u_{1,2}(\tau + \mu_{1,2} + \mu_{2,2} + \mu_{3,2}) = \zeta u_{1,1} + \mu u_{2,2} \quad (36)$$

$$(\zeta + \mu + \lambda)u_{m,1} = \lambda u_{m-1,1} + \mu u_{m+1,1} + \tau u_{m,2}, \quad m \in [2, M-1] \quad (37)$$

$$(\tau + \mu)u_{m,2} = \mu u_{m+1,2} + \zeta u_{m,1}, \quad m \in [2, M-1] \quad (38)$$

$$(\mu + \zeta)u_{M,1} = \lambda u_{M-1,1} + \tau u_{M,2} \quad (39)$$

$$(\tau + \mu)u_{M,2} = \zeta u_{M,1} \quad (40)$$

$$\sum_{j=1}^2 \left(\sum_{i=1}^3 \nu_{i,j} + \sum_{m=1}^M u_{m,j} \right) = 1 \quad (41)$$

In Equations (29)-(41), we have $2M + 6$ unknowns and $2M + 6$ independent equations. Since all equations are linear, well-known techniques can solve this system of equations. Assume matrix A be the coefficient matrix, X be the vector of unknowns, and B be the vector of constants. One can easily solve the linear equation sets, i.e., $AX = B$, using existing solvers, such as Matlab or IBM CPLEX.

REFERENCES

- [1] M. Masoudi, M. G. Khafagy, A. Conte, A. El-Amine *et al.*, "Green mobile networks for 5G and beyond," *IEEE Access*, vol. 7, pp. 107 270–107 299, 2019.
- [2] ETSI, "Energy consumption and CO2 footprint of wireless networks," Report RRS05-024, 2011.
- [3] M. Masoudi and C. Cavdar, "Device vs edge computing for mobile services: Delay-aware decision making to minimize power consumption," *IEEE Transactions on Mobile Computing*, vol. 20, no. 12, pp. 3324–3337, 2020.
- [4] B. Debaille, C. Desset, and F. Louagie, "A flexible and future-proof power model for cellular base stations," in *2015 IEEE 81st Vehicular Technology Conference (VTC Spring)*, May 2015, pp. 1–7.
- [5] M. Masoudi, O. T. Demir, J. Zander, and C. Cavdar, "Energy-optimal end-to-end network slicing in cloud-based architecture," *IEEE Open Journal of the Communications Society*, 2022.
- [6] P. Chang and G. Miao, "Joint optimization of base station deep-sleep and DTX micro-sleep," in *2016 IEEE Globecom Workshops (GC Wkshps)*. IEEE, 2016, pp. 1–6.
- [7] J. Wu, S. Zhou, and Z. Niu, "Traffic-aware base station sleeping control and power matching for energy-delay tradeoffs in green cellular networks," *IEEE Transactions on Wireless Communications*, vol. 12, no. 8, pp. 4196–4209, 2013.
- [8] J. Wu, Y. Bao, G. Miao, S. Zhou, and Z. Niu, "Base-station sleeping control and power matching for energy-delay tradeoffs with bursty traffic," *IEEE Transactions on Vehicular Technology*, vol. 65, no. 5, pp. 3657–3675, 2015.
- [9] Z. Niu, X. Guo, S. Zhou, and P. R. Kumar, "Characterizing energy-delay tradeoff in hyper-cellular networks with base station sleeping control," *IEEE Journal on Selected Areas in Communications*, vol. 33, no. 4, pp. 641–650, 2015.
- [10] K. Son, H. Kim, Y. Yi, and B. Krishnamachari, "Base station operation and user association mechanisms for energy-delay tradeoffs in green cellular networks," *IEEE journal on selected areas in communications*, vol. 29, no. 8, pp. 1525–1536, 2011.

- [11] J. Kim, H.-W. Lee, and S. Chong, "Traffic-aware energy-saving base station sleeping and clustering in cooperative networks," *IEEE Transactions on Wireless Communications*, vol. 17, no. 2, pp. 1173–1186, 2017.
- [12] X. Guo, Z. Niu, S. Zhou, and P. Kumar, "Delay-constrained energy-optimal base station sleeping control," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 5, pp. 1073–1085, 2016.
- [13] P. Chang and G. Miao, "Optimal operation of base stations with deep sleep and discontinuous transmission," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 11, pp. 11 113–11 126, 2018.
- [14] G. Andersson, A. Västberg, A. Devlic, and C. Cavdar, "Energy efficient heterogeneous network deployment with cell DTX," in *2016 IEEE International Conference on Communications (ICC)*. IEEE, 2016, pp. 1–6.
- [15] C. Cox, *An Introduction to LTE: LTE, LTE-Advanced, SAE, VoLTE and 4G Mobile Communications*, 2nd ed. Wiley Publishing, 2014.
- [16] S. K. G. Peesapati, M. Olsson, S. Andersson, M. Masoudi, and C. Cavdar, "An analytical energy performance evaluation methodology for 5G base stations," in *2021 17th International Conference on Wireless and Mobile Computing (WiMob)*, Bologna, Italy, Oct. 2021.
- [17] F. E. Salem, A. Gati, Z. Altman, and T. Chahed, "Advanced sleep modes and their impact on flow-level performance of 5G networks," in *2017 IEEE 86th Vehicular Technology Conference (VTC-Fall)*, Sep. 2017.
- [18] F. E. Salem, Z. Altman, A. Gati, T. Chahed, and E. Altman, "Reinforcement learning approach for advanced sleep modes management in 5G networks," in *2018 IEEE 88th Vehicular Technology Conference (VTC-Fall)*, Chicago, USA, Aug. 2017.
- [19] F. E. Salem, T. Chahed, Z. Altman, and A. Gati, "Traffic-aware advanced sleep modes management in 5G networks," in *2019 IEEE Wireless Communications and Networking Conference (WCNC)*, Marrakech, Morocco, Apr. 2019.
- [20] A. El-Amine, M. Iturralde, H. A. H. Hassan, and L. Nuaymi, "A distributed Q-Learning approach for adaptive sleep modes in 5G networks," in *2019 IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE, 2019, pp. 1–6.
- [21] S. K. G. Peesapati, M. Olsson, M. Masoudi, S. Andersson, and C. Cavdar, "Q-Learning based radio resource adaptation for improved energy performance of 5G base stations," in *2021 IEEE 32nd Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, Helsinki, Finland, Sep. 2021.
- [22] H. Pervaiz, O. Onireti, A. Mohamed, M. A. Imran, R. Tafazolli, and Q. Ni, "Energy-efficient and load-proportional eNodeB for 5G user-centric networks: A multilevel sleep strategy mechanism," *IEEE Vehicular Technology Magazine*, vol. 13, no. 4, pp. 51–59, 2018.
- [23] M. Masoudi, M. G. Khafagy, E. Soroush, D. Giacomelli, S. Morosi, and C. Cavdar, "Reinforcement learning for Traffic-adaptive sleep mode management in 5G networks," in *2020 IEEE 31st Annual International Symposium on Personal, Indoor and Mobile Radio Communications*, London, United Kingdom, Aug. 2020.
- [24] F. Tao, H. Zhang, A. Liu, and A. Y. Nee, "Digital twin in industry: State-of-the-art," *IEEE Transactions on Industrial Informatics*, vol. 15, no. 4, pp. 2405–2415, 2018.
- [25] Y. Ma, H. Zhou, H. He, G. Jiao, and S. Wei, "A digital twin-based approach for quality control and optimization of complex product assembly," in *2019 International Conference on Artificial Intelligence and Advanced Manufacturing (AIAM)*. IEEE, 2019, pp. 762–767.
- [26] Y. Lu, X. Huang, K. Zhang, S. Maharjan, and Y. Zhang, "Communication-efficient federated learning for digital twin edge networks in industrial iot," *IEEE Transactions on Industrial Informatics*, 2020.
- [27] imec, GreenTouch Project, "Power model for wireless base stations." [Online]. Available: <https://www.imec-int.com/en/powermodel>
- [28] J. Yang, W. Wang, and X. Zhang, "Hysteretic base station sleeping control for energy saving in 5G cellular network," in *2017 IEEE 85th Vehicular Technology Conference (VTC Spring)*. IEEE, 2017, pp. 1–5.
- [29] J. Liu, B. Krishnamachari, S. Zhou, and Z. Niu, "Deepnap: Data-driven base station sleeping operations through deep reinforcement learning," *IEEE Internet of Things Journal*, vol. 5, no. 6, pp. 4273–4282, 2018.
- [30] J. E. Beyer and B. F. Nielsen, "Predator foraging in patchy environments: the interrupted poisson process (IPP) model unit," *DANA-CHARLOTTENLUND-*, vol. 11, 1996.
- [31] A. McCallum, "Hidden markov models baum welch algorithm," *Introduction to Natural Language Processing CS585; University of Massachusetts Amherst: Massachusetts, USA*, 2004.
- [32] I. Sousa, M. P. Queluz, and A. Rodrigues, "A survey on QoE-oriented wireless resources scheduling," *Journal of Network and Computer Applications*, vol. 158, p. 102594, 2020.
- [33] S.-E. Elayoubi, L. Saker, and T. Chahed, "Optimal control for base station sleep mode in energy efficient radio access networks," in *2011 Proceedings IEEE INFOCOM*. IEEE, 2011, pp. 106–110.
- [34] M. Dolfi, C. Cavdar, S. Morosi, P. Piuanti, J. Zander, and E. Del Re, "On the trade-off between energy saving and number of switchings in green cellular networks," *Transactions on Emerging Telecommunications Technologies*, vol. 28, no. 11, p. e3193, 2017.
- [35] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [36] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [37] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.