# Distributed Q-Learning for Energy Harvesting Heterogeneous Networks

Marco Miozzo[‡], Lorenza Giupponi[‡], Michele Rossi[†], Paolo Dini[‡]

[‡]*CTTC, Av. Carl Friedrich Gauss, 7, 08860, Castelldefels, Barcelona, Spain*
[†]*DEI, University of Padova, Via G. Gradenigo, 6/B, 35131, Padova, Italy*
{mmiozzo, lgiupponi, pdini}@cttc.es, rossi@dei.unipd.it

*Abstract*—We consider a two-tier urban Heterogeneous Network where small cells powered with renewable energy are deployed in order to provide capacity extension and to offload macro base stations. We use reinforcement learning techniques to concoct an algorithm that autonomously learns energy inflow and traffic demand patterns. This algorithm is based on a decentralized multi-agent Q-learning technique that, by interacting with the environment, obtains optimal policies aimed at improving the system performance in terms of drop rate, throughput and energy efficiency. Simulation results show that our solution effectively adapts to changing environmental conditions and meets most of our performance objectives. At the end of the paper we identify areas for improvement.

*Index Terms*—Mobile Networks, HetNet, Sustainability, Renewable Energy, Energy Efficiency, Q-Learning.

## I. Introduction

In the near future, mobile network operators will have to handle a much higher capacity demand [1], especially within urban areas. In fact, it is expected that fifth generation (5G) mobile networks will support $1,000$ times more capacity per unit area than 4G. On the downside, this trend will affect the greenhouse gases emissions of ICT ecosystems, which already consume about $1500$ TWh of energy annually, approaching $10\%$ of the world's electricity generation and accounting for $2-4\%$ of the carbon footprint due to human activities. This calls for a radical change in the system design philosophy, shifting from a coverage and capacity oriented optimization (typical of 3G and 4G) to an energy oriented one. Besides, operators are also looking at the massive deployment of small scale factor base stations (BSs), which are referred to as *small cells* (SCs). SCs will provide capacity extension to macro BSs, giving rise to a multi-tier heterogeneous network (HetNet). In addition, the reduced energy consumption of small cells could be exploited to self-power these SCs through the use of e.g., small size solar panels, as we advocate in this paper, through the concept of *energy harvesting HetNets*. This is expected to reduce the cost associated with the purchase of energy from the power grid, and so the carbon footprint of ICT [2].

A proper management of the SCs calls for a lightweight and scalable architecture, including new management procedures. In this respect, Self Organized Networking (SON) will be key to bring intelligence and autonomous adaptability to network elements, by diminishing human involvement, while maximizing system performance, reducing operational costs, meeting QoS requirements and improving the overall energy efficiency (EE) [3]. This paradigm is of paramount importance, especially considering that SCs will be operated in an uncoordinated fashion. Previous research work demonstrates that sleep strategies (or switch ON-OFF) are a valuable means to reach these goals [4]. However, in the case of energy harvesting (EH) HetNets, we also need to consider the erratic and intermittent nature of renewable energy sources, which entails some additional complexity.

In this work, we consider solar energy as a reference renewable energy source (RES), due to its widespread availability, the good efficiency of photovoltaic (PV) technology and its competitive cost. On this matter, we observe that one may obtain some prediction of the amount of energy that is harvested on a daily basis and, taking into account bad weather conditions, may over dimension the PV panel and the associated energy storage (usually made of Lithium ion rechargeable cells), to meet a certain service availability criterion. This practice is commonly adopted for off-grid installations (see, e.g., rural areas), where network elements cannot be easily connected to the power grid and have to solely rely on harvested or diesel-generated energy. In this case, the resulting panel size for urban scenarios is not viable for macro BS and still too large even for SCs [2], which are supposed to be installed in street furnitures (i.e., traffic lamps, street lights, transportation hubs, etc.).

To overcome this limitation, we target hybrid two-tier deployments where macro BSs reside in the first tier and are powered by the power grid, whereas SCs operate within the second one supplied by solar panels. In this scenario, several optimizations are possible, such as offloading some of the data traffic from the macro BS to the SCs or switching ON or OFF the SCs, based on the traffic demand and the energy offer. This makes it possible to reduce the requirements (i.e., solar panel size and battery capacity) for the SCs, and the energy cost of the macro BS, at the cost of some additional complexity in terms of hardware (the installation of SCs) and optimization algorithms (for traffic offloading). The design of these algorithms, along with the quantitative assessment of their effectiveness, is the main objective of the present paper.

The increasing interest in EH cellular networks is testified by the rich literature on this topic [4]. In [5], the authors present a design based on stochastic geometry for the management of $k$-tier HetNets powered by RESs. Their model controls the fraction of time that each tier can be kept on, according

to its energy reserve. Similarly, in [6], the authors propose an algorithm to control the BS power consumption as a function of the energy reserve and the expected amount of renewable energy that will be stored. However, neither of these two works considers the temporal variations in traffic and in harvested energy processes, which is fundamental for a realistic model of the scenario. In [7], the authors focus on off-grid mesh networks of EH BSs. First, the problem of dimensioning the renewable energy "add-on" (solar panel and battery) is solved by considering typical daily traffic and harvested energy profiles for different cities. Then, an optimization approach, for two-tier networks, is proposed, based on SCs sleep modes. However, the proposed optimization approach is based on historical data, and is consequently unable to adapt to the dynamic system conditions, in terms of harvested energy or traffic demand, as it would be desirable in a realistic setting.

In this paper, we overcome the above mentioned problems, and we propose to model the SC network by means of a multi-agent system where each agent makes autonomous decisions, according to the Decentralized SON (D-SON) paradigm. In this context, we propose a distributed on-line solution based on a multi-agent Reinforcement Learning (RL) algorithm, known as *distributed Q-learning*. Through RL, each agent (SC) independently learns a proper radio resource management (RRM) policy, so as to jointly maximize the system performance in terms of throughput, drop rate and energy consumption, while adapting to the dynamic conditions of the environment, in terms of energy inflow and traffic demand. The performance of this algorithm is then assessed for two-tier networks, considering realistic models for the data traffic and for the energy harvested. While still preliminary, these results are encouraging and show that our approach is viable as the designed algorithm meets most of our design goals.

The remainder of the paper is organized as follows. In Section II we present the system model, whereas the RL algorithm is presented in Section III. In Section IV we discuss some performance results. In Section V we draw our conclusions and discuss future research directions.

## II. SYSTEM MODEL

In this section we describe the network and energy management models. We consider a two-tier HetNet composed of heterogeneous LTE BSs, which includes one macro BS and $N$ SCs. The macro BS is connected to the power grid and provides baseline coverage to the whole cell. The SCs are deployed in a hotspot manner to increase the capacity where needed (e.g., shopping hall, city center, etc.). Also, these SCs are solely powered through solar-harvested energy and are controlled in a distributed fashion by means of Q-learning agents, as we detail in the next section.

At the physical layer, LTE is based on OFDMA. The total transmission bandwidth $B$ is divided into $R$ resource blocks (RBs) of 1 msec each (referred to as TTI). Each SC $i$ has a set $\mathcal{U}_i$ of associated users, which depend on its geographical location and on the distribution of the users (see Section IV-A for further details). For the BS power consumption, we adopt

TABLE I
POWER MODEL PARAMETERS FOR VARIOUS TYPES OF BS.

| BS Type | $P_0$ [W] | $\beta$ [W] |
|---------|-----------|-------------|
| Macro   | 750.0     | 600         |
| Small   | 105.6     | 39          |

the model presented in [8]. In particular, the energy consumption of a LTE BS can be approximated by the linear function $P = P_0 + \beta \rho$, where $\rho \in [0, 1]$ is the traffic load of the BS, normalized with respect to its maximum capacity, and $P_0$ is the baseline power consumption. Typical values of $\beta$ and $P_0$ are reported in Table I for macro and small BSs. Regarding the type of SC, we consider medium scale factor "metro cells", as the Alcatel-Lucent 9764 Metro Cell Outdoor, featuring a maximum transmission power of 38 dBm. The (time-varying) BS capacity (in terms of number of resource blocks allocated to the users) is defined based on [9]. This includes the simulation of the wireless channels and the selection of the modulation and coding scheme (MCS) for each user, based on the particular channel conditions and on the (dynamically computed) system interference. For the SC management, we assume a slotted time model with a slot duration of 1 hour. This time granularity is deemed appropriate to track variations in the system load and in the EH process.

## III. ALGORITHM

In this section, we present a decentralized multi-agent radio resource management algorithm for the SCs in the second tier.

### A. Q-learning-based Radio Resource Management (RRM)

We consider a network setup of $N$ distributed agents (the SCs), which can be modeled by means of a multi-agent system, as it fulfils the following conditions: (1) the intelligent RRM decisions are made by multiple intelligent and uncoordinated agents; (2) the agents partially observe the overall scenario; and (3) their inputs to the intelligent decision process differ from agent to agent, since they come from spatially distributed sources of information. In particular, the inputs to the RRM algorithm depend on the SC's particular location and on the geographical distribution of the users (i.e., the load). The objective of the algorithm is for each agent to learn, through real-time interactions with the environment, an energy management policy by means of a Q-learning approach. The decision making process of each agent is defined by a Markov Decision Process with state vector $\vec{x}_t = \{x_t^1, x_t^2, \ldots, x_t^N\}$, where $x_t^i$ is the state associated with SC $i$ (described in the next Section III-B), at time $t$. Based on $x_t^i$, each agent $i$ *independently* chooses an action $a_t^i$ from an action set $\mathcal{A}$. As a result of the execution of this action, the environment returns an *agent dependent* reward $r_t^i$, which allows the local update of a Q-value, $Q(x_t^i, a_t^i)$, indicating the appropriateness of selecting action $a_t^i$ in state $x_t^i$. The Q-value is computed according to the rule:

$$Q(x_t^i, a_t^i) \leftarrow Q(x_t^i, a_t^i) + \alpha[r_t^i + \gamma \min_a Q(x_{t+1}^i, a') - Q(x_t^i, a)] \tag{1}$$

where $\alpha$ is the learning rate, $\gamma$ is the discount factor, $x^i_{t+1}$ is the next state for agent $i$ and $a'$ is the associated optimal action. For more details on RL and Q-learning the reader is referred to, e.g., [10], [11].

### B. States, actions and rewards

In this section we provide details on the Q-learning algorithm, by defining state, action set and reward function, for the $N$ agents.

**State:** The local state $x^i_t$ is defined by:

$$x^i_t = \{S^i_t, B^i_t, L^i_t\}, \tag{2}$$

where $S^i_t$ is the state of the renewable energy source based on the incoming amount harvested energy (e.g., day and night), $B^i_t$ is the normalized battery energy level, $L^i_t$ is the normalized load for SC $i$ in slot $t$, which depends on the number of users served by this SC. We uniformly quantize $S^i_t$, $B^i_t$ and $L^i_t$ into 2, 5 and 3 levels, respectively.

**Action set:** The set of possible actions $\mathcal{A}$ consists of the two actions of switching ON and OFF the SC. We have not considered the option of modulating the load $\rho$ between 0 and 1, due to the energy profile of SCs. In fact, the $\beta$ parameter in Table I for the SCs is usually small, and therefore the parameter $\rho$ has a marginal impact on their energy consumption. When a SC is switched OFF, the associated users have to connect to the macro BS. However, in case the macro BS is not able to provide them with service, they will be dropped, until the next time slot, when a variation of system state may lead to different RRM decisions.

**Reward function:**

$$r^i_t = \begin{cases} 0 & B^i_t < B_{\mathrm{th}} \text{ or } D_t < D_{\mathrm{th}} \\ \kappa T^i_t & B^i_t \geqslant B_{\mathrm{th}} \text{ and } D_t \geqslant D_{\mathrm{th}} \text{ and SC } i \text{ is ON} \\ 1/B^i_t & B^i_t \geqslant B_{\mathrm{th}} \text{ and } D_t \geqslant D_{\mathrm{th}} \text{ and SC } i \text{ is OFF} \end{cases} \tag{3}$$

where $T^i_t$ is the normalized throughput of SC $i$ in slot $t$, $D_t$ is the instantaneous *system* drop rate, defined as the ratio between the total amount of traffic dropped and the traffic demand in the entire network (accounting for macro and small BSs). $D_{\mathrm{th}}$ is the maximum tolerable drop rate. Finally, $B_{\mathrm{th}}$ is a threshold on the battery level. The rationale behind (3) is the following. The condition in the first line implies a zero reward when the battery level falls below $B_{\mathrm{th}}$ ($B^i_t < B_{\mathrm{th}}$) or the system drop rate is below $D_{\mathrm{th}}$ ($D_t < D_{\mathrm{th}}$). This incentivizes the SC to turn itself OFF to save energy, as this implies a higher reward. When $B^i_t < B_{\mathrm{th}}$, this is necessary to promote the energetic self-sustainability of the SC, whereas when $D_t > D_{\mathrm{th}}$, the system performance is deemed sufficient. Thus, the SC can be switched OFF and offload the macro BS at a later time. In the second and third line of (3), the reward is proportional to the throughput when the SC is turned ON and is instead proportional to the inverse of the energy buffer level when the SC is OFF. Note that the SC, after a learning phase, will choose to remain ON (and offload the macro BS) when the reward in the second line is higher, i.e., when $\kappa T^i_t > 1/B^i_t$. Note that $1/B^i_t$ may dominate over $\kappa T^i_t$ in case battery level

and throughput are both low. In this case, the SC switches OFF to save energy. The constant $\kappa$ is used to balance the impact of the two terms (throughput *vs* energy saving).

## IV. PERFORMANCE EVALUATION

### A. Simulation Scenario

We consider $N$ SCs operating within a square macro cell area with a side of 1 km ($N$ is varied as a free parameter in Section IV). The macro BS is placed in the center of this area, whereas SCs are randomly positioned with the constraint that their cells do not overlap. Specifically, we pick a transmission power of 38 dBm for SCs, which translate into a coverage radius of 50 m. 120 users (UEs) are uniformly placed within the coverage area of each SC. The number of UEs has been selected so that the SCs are congested during the traffic peaks. The traffic of these users follows a urban profile [7] (i.e., traffic peaks are concentrated around working hours). For what concerns the distribution of traffic among users, we adopt the model in [8], configuring $20\%$ of the UEs as heavy users (their data volume is 900 MB/h), while the remaining UEs are ordinary users (112.5 MB/h). For the renewable energy sources we consider the Panasonic N235B solar modules, which have single cell efficiencies of about $21\%$, delivering about $186 \mathrm{W/m}^2$. For SCs, an array of $16 \times 16$ (4.48 m$^2$) solar cells has been chosen. The battery size of the small cell is 2 kWh (panel and battery sizes have been chosen so that SC batteries can be replenished in a full winter day). Harvested energy traces have been obtained using the SolarStat tool [12], considering the city of Los Angeles as the deployment location. These traces have been translated into a Markov process with 12 energy states that, as shown in [12], provide an excellent approximation of the harvested energy process, and so are used for this purpose in our simulations. Fig. 1 shows typical profiles for the traffic demand and the harvested energy across two subsequent days. Interestingly, we see that the maxima in the energy inflow and in the traffic demand are not aligned. This means that some optimization actions that could be taken are e.g., saving energy resources and use them when the next traffic peak occurs.

The decentralized Q-learning algorithm of Section III is independently implemented by each SC. The learning rate is set to $\alpha = 0.5$ and the discount factor to $\gamma = 0.9$ for all SCs, according to our simulation analysis. The constant $\kappa$ (see (3)) is set to 10 as this provides a good tradeoff for the considered system parameters. The Q-learning algorithm also implements exploration features [10], i.e., random states are visited by the learning agents with probability $\varepsilon = 0.1$. In the following plots, we refer to "QL" as our Q-learning solution. We compare QL against a greedy scheme ("greedy" in the figures) where SCs are put into a sleep mode (OFF) when their battery level $B_t$ drops below $B_{\mathrm{th}}$, and they are switched ON when $B_t \geqslant B_{\mathrm{th}}$. The battery threshold $B_{\mathrm{th}}$ is set to $30\%$ of the battery capacity. The threshold on the instantaneous traffic drop rate is set to $D_{\mathrm{th}} = 0.05$. Simulations are run for 420 consecutive days, where 60 of them are used for the training phase, while the results from the remaining 360
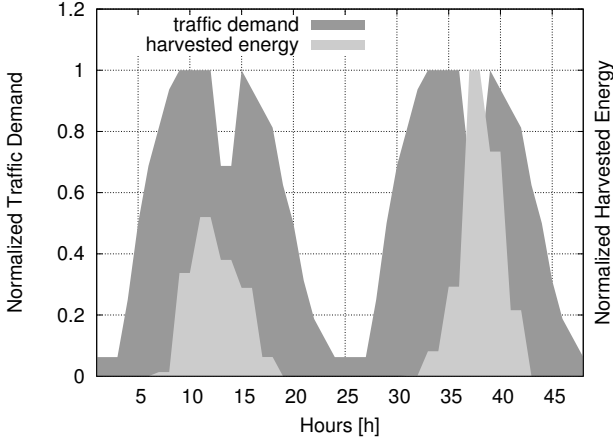
Fig. 1. Examples of total traffic demand and amount of energy harvested.
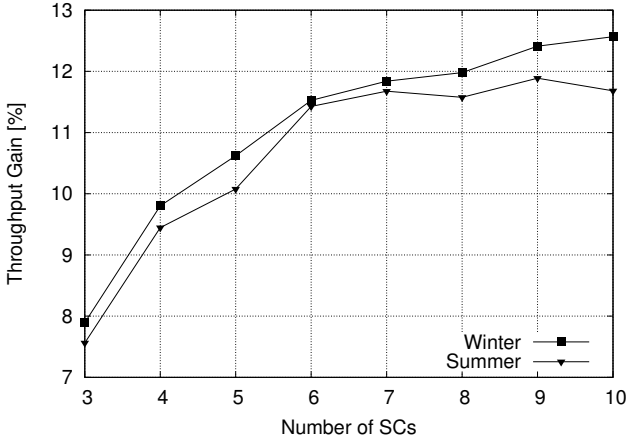


Fig. 2. Throughput gain [%] of QL with respect to the greedy scheme.

days are used to evaluate the behaviour of QL and the greedy approach. In the following plots, we treat separately the *winter* and the *summer* months, as the energy harvesting statistics are very different in these two cases. Specifically, we consider as *winter* the months of January, February, October, November and December, while the remaining months are classified as *summer*.

### B. Numerical Results

In Fig. 2, we show the system throughput gain provided by QL with respect to the greedy scheme. It can be observed that the QL approach offers improvements of up to 14%, during the winter months.

In Fig. 3, we plot an example of the temporal system behavior for a HetNet including 3 SCs and a macro BS for the last week of December. Here, from top to bottom we show temporal traces concerning traffic demand and instantaneous harvested energy (in the same plot), battery level, policy adopted at the SCs (y-label "Action") and normalized load at the macro BS (y-label "Macro Load"). From these results various observations can be made. First, the policy adopted by

QL tends to save energy during the night, and this makes it possible to offload more the macro BS during the day, as it can be seen in the bottom plot of Fig. 3 in correspondence of the points marked with "(a)". Also, the impact of our reward function (see (3)) can be appreciated in correspondence of label "(b)". Here, the QL keeps the SCs ON, as the traffic demand is high, and in this case sleeping would cause congestion at the macro BS. We remark that QL is capable of doing this as it proactively saves some of the harvested energy when the energy inflow is abundant. In contrast, the greedy scheme shows a more aggressive behavior and, as a result, it has no residual energy to compensate for an upsurge in the traffic load.

We observe that the energy harvesting traces are the same for all SCs. We implement this choice since it is expected that the level of solar irradiation will not change much within a macro cell area. In addition, this sort of synchronization with respect to the experienced energy inflow from RESs is enforced by the traffic demand processes, as different SCs will as well undergo similar traffic profiles. This implies that, in the considered setup, SCs are often switched ON/OFF simultaneously. This can be appreciated from Fig. 4, where the average load is plotted as a function of the hour of the day for a network with 3 SCs. The greedy scheme usually leads to a higher load for the macro BS during the morning peak hours, where the batteries are likely to be drained, and therefore most of the SCs must be turned OFF. On the contrary, QL loads the macro BS slightly more during most of the day in order to put some of the SCs to sleep (saving energy at these SCs) and serve more traffic during the morning peak.

The traffic drop rate as a function of the number of SCs is shown in Fig. 5, where it can be observed that the QL algorithm considerably reduces the drop rate compared to the greedy scheme. The throughput improvements directly translate into improvements in terms of user QoS, as depicted in Fig. 6, where the Jain's fairness index (JFI) is plotted as a function of the number of SCs. If $T^i$ is the throughput experienced by UE $i$, $T_{\mathrm{req}}^i$ is the capacity requested by this user and $N_u$ is the number of UEs, the JFI is defined as $\mathrm{JFI} = [\sum_{i=1}^{N_u}(T^i/T_{\mathrm{req}}^i)]^2/[N_u\sum_{i=1}^{N_u}(T^i/T_{\mathrm{req}}^i)^2]$. In terms of JFI, QL provides an improvement higher than 5% with respect to the greedy scheme, as fewer users are dropped.

We define by *battery outage* the amount of time a SC spends with a battery level $B$ below the threshold $B_{\mathrm{th}}$. In this case, the SC has to be momentarily put into sleep, independently of the adopted policy. Based on our model, the battery outage of the greedy scheme is always higher than 4 hours per day, reaching a maximum of 8 hours in the winter. On the other hand, QL offers an average battery outage below 1 hour, except in winter, when it gets close to 2 hours. This is achieved thanks to the intelligent behavior of QL, which proactively reacts to the reward function, defined to optimize the battery outage.

In Fig. 7, we show the average cell load for the macro BS. It can be observed that in general, when implementing QL, the macro cell ends up serving more load. The reason is that with QL, during off-peak hours, the SCs tend to save some of
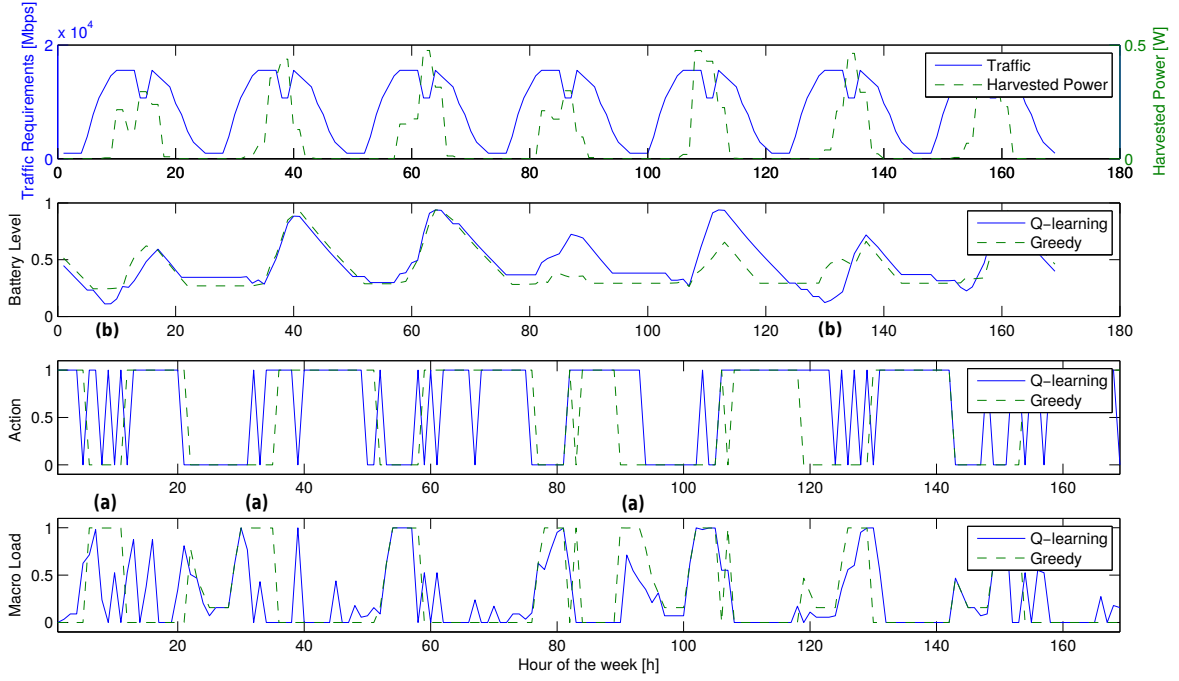
Fig. 3. Example temporal behavior for a HetNet with 3 SCs and one macro BS. Temporal traces show the status of the SCs.
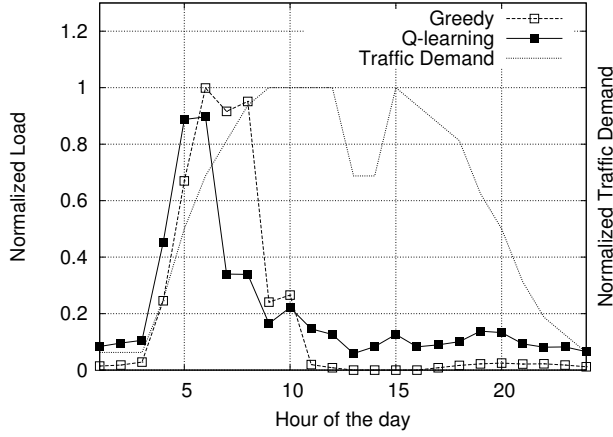


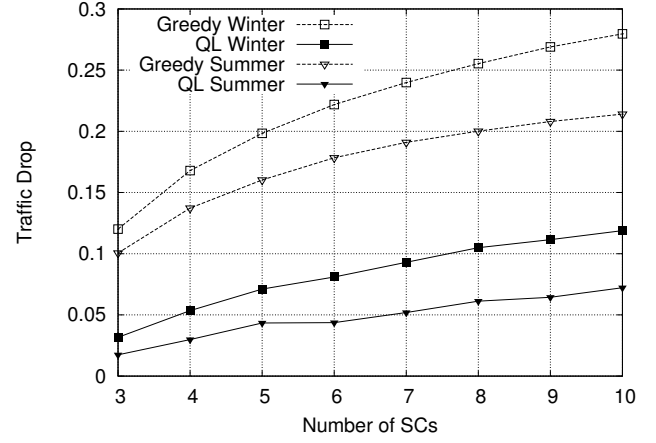Fig. 4. Average hourly load for the macro BS in a network with 3 SCs.



Fig. 5. Traffic drop rate for QL and greedy.

the harvested energy by turning themselves OFF for a longer period of time than with the greedy solution. As a result of this, the macro BS may result to be more loaded. The saved energy is then used by QL to compensate for the traffic peaks, where more SCs are turned ON. Overall, this reduces the amount of traffic dropped, increases the average throughput and loads more the macro BS when the traffic volume is moderate.

In Fig. 8, we look at the EE, which is defined as EE $= E_S/T_S$, where $E_S$ is the total energy drained by the macro BS from the power grid and $T_S$ is the system throughput. As we can see, QL offers a higher EE than the greedy scheme. However, the EE diminishes for an increasing number of SCs

because the macro BS has to serve a higher number of UEs when the SCs are switched OFF. Finally, when we look at the total amount of energy spent by the system, it is proportional to the served traffic (which is higher for the QL option), so that it approximatively amounts to $7.5$ MWh in a year for a greedy solution, while it varies from $7.5$ (with 3 SCs) to $8.3$ MWh (with 10 SCs) when QL is adopted.

As a final remark, it is worth mentioning that the same system implemented without energy harvesting capabilities (i.e., where the SCs are grid-connected) would consume from $9.6$ (3 SCs) to 17 MWh (10 SCs) in a year, which implies an increment in terms of used energy of more than $50\%$.
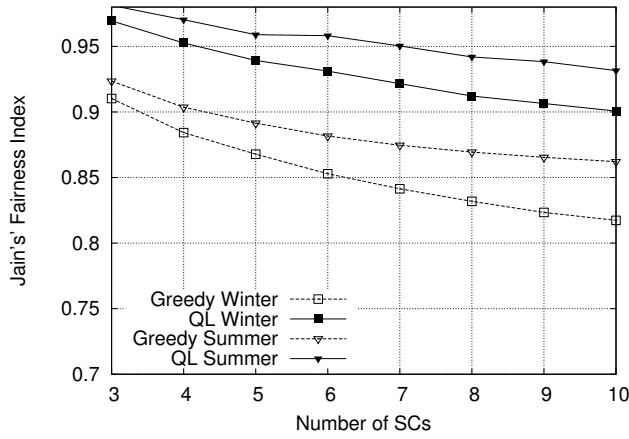
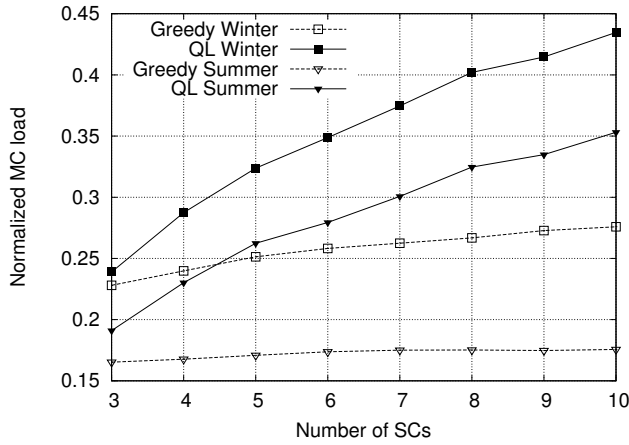Fig. 6. Jain's fairness index for QL and greedy *vs* number of SCs.



Fig. 7. Macro BS load for QL and greedy *vs* number of SCs.



Fig. 8. Energy efficiency improvement [%] of QL with respect to greedy *vs* number of SCs.

## V. CONCLUSIONS AND WAY FORWARD

In this paper, we have presented a distributed Q-learning algorithm for the management of energy harvesting SCs in heterogeneous networks. Our scheme is designed to increase the system throughput, offload the macro BSs and decrease the drop rate at the macro BS. Our simulation results are encouraging and show that the proposed approach is viable, as the algorithm meets most of our design goals and also improves the energy efficiency of the system.

Nevertheless, there are various aspects that need to be further investigated. First, we would like to enhance the decisions made by the distributed small cells so that they will cooperatively compute optimal policies accounting for common (and global) performance objectives. Note that in the current algorithm this cooperation is only marginally achieved through the use of the global drop rate $D_t$ in the reward functions that are locally computed by the small cells (see (3)). Finally, we need to explore further reward functions so as to still obtain performance gains even when the number of small cells is large. In particular, we also plan to embed the energy efficiency metrics into the learning algorithm.
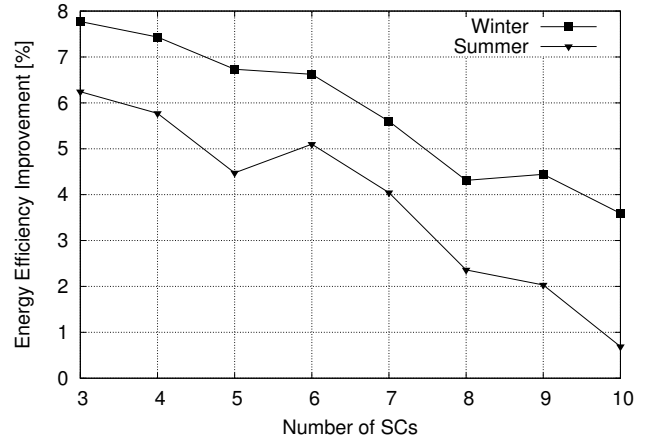
## REFERENCES

[1] Cisco Systems Inc., "Cisco visual networking index global mobile data traffic forecast update 2013-2018," White Paper, http://www.cisco.com/, Feb. 2013.
[2] G. Piro, M. Miozzo, G. Forte, N. Baldo, L. A. Grieco, G. Boggia, and P. Dini, "HetNets Powered by Renewable Energy Sources," *IEEE Internet Computing*, vol. 17, no. 1, pp. 32–39, 2013.
[3] 3GPP, "TS 36.927 v12.0.0; LTE; E-UTRAN; potential solutions for energy saving for E-UTRAN," 2014.
[4] H. Al Haj Hassan, L. Nuaymi, and A. Pelov, "Renewable energy in cellular networks: A survey," in *IEEE Online Conference on Green Communications (GreenCom)*, Oct. 2013.
[5] H. Dhillon, Y. Li, P. Nuggehalli, Z. Pi, and J. Andrews, "Fundamentals of Heterogeneous Cellular Networks with Energy Harvesting," *IEEE Transactions on Wireless Communications*, vol. 13, no. 5, pp. 2782–2797, May 2014.
[6] D. Valerdi, Q. Zhu, K. Exadaktylos, S. Xia, M. Arranz, R. Liu, and D. Xu, "Intelligent energy managed service for green base stations," in *IEEE Global Communications Conference (GLOBECOM)*, Miami, FL, US, Dec. 2010.
[7] M. Marsan, G. Bucalo, A. Di Caro, M. Meo, and Y. Zhang, "Towards zero grid electricity networking: Powering BSs with renewable energy sources," in *IEEE International Conference on Communications (ICC)*, Budapest, Hungary, Jun. 2013.
[8] EU EARTH: Energy Aware Radio and neTwork tecHnologies, "D2.3: Energy efficiency analysis of the reference systems, areas of improvements and target breakdown," Deliverable D2.3, www.ict-earth.eu, 2010.
[9] M. Mezzavilla, M. Miozzo, M. Rossi, N. Baldo, and M. Zorzi, "A Lightweight and Accurate Link Abstraction Model for System-Level Simulation of LTE Networks in ns-3," in *ACM MSWIM*, Paphos, Cyprus Island, Oct. 2012.
[10] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press, 1998.
[11] C. J. Watkins and P. Dayan, "Technical note: Q-learning," *Machine Learning*, vol. 8, pp. 279–292, 1992.
[12] M. Miozzo, D. Zordan, P. Dini, and M. Rossi, "SolarStat: Modeling Photovoltaic Sources through Stochastic Markov Processes," in *IEEE Energy Conference (ENERGYCON)*, Dubrovnik, Croatia, May 2014.