

# Distributed Q-Learning Based Dynamic Spectrum Access in High Capacity Density Cognitive Cellular Systems Using Secondary LTE Spectrum Sharing

Nils Morozs, David Grace and Tim Clarke  
Department of Electronics, University of York  
Heslington, York YO10 5DD, United Kingdom  
E-mail: {nm553, david.grace, tim.clarke}@york.ac.uk

**Abstract**—In this paper a distributed Q-learning based dynamic spectrum access (DSA) algorithm is applied to a cognitive cellular system designed for providing ultra high capacity density with only secondary access to an LTE channel. Large scale simulations of a stadium temporary event scenario show that the distributed Q-learning based DSA scheme provides robust quality of service (QoS) and extremely high system throughput densities to the users of the stadium network, whilst successfully coexisting with a primary network of macro eNodeBs on the same LTE channel. It is also shown that incorporating spectrum awareness or spectrum sensing based admission control into the DSA algorithm in this scenario does not improve its performance. Therefore, distributed Q-learning based DSA is a viable and easily implementable solution for facilitating secondary LTE spectrum sharing in high capacity density cognitive cellular systems.

**Keywords**—Cognitive Cellular Systems, Small Cells, Reinforcement Learning, Spectrum Sharing, Dynamic Spectrum Access

## I. INTRODUCTION

One of the fundamental tasks of a cellular system is spectrum management, concerned with dividing the available spectrum into a set of resource blocks or subchannels and assigning them to voice calls and data transmissions in a way which would provide a good quality of service (QoS) to the users. Flexible dynamic spectrum access (DSA) techniques play a key role in utilising the given spectrum efficiently. This has given rise to novel wireless communication systems such as cognitive radio networks [1] and cognitive cellular systems [2]. Such networks employ intelligent opportunistic DSA techniques that allow them to access licensed spectrum underutilized by the incumbent users.

The classical and most common application of spectrum sharing in cognitive radio networks to date is use of the TV white spaces. Such networks aim to reuse the spectrum allocated to TV broadcasters for other wireless communications, whilst eliminating harmful interference to the incumbent TV receivers, e.g. [3][4]. A more recent problem investigated by researchers, mobile network operators (MNOs) and regulators is LTE and LTE-Advanced spectrum sharing [5]. In many cases LTE spectrum sharing is required by two or more co-primary MNOs. This can be facilitated by an emerging framework known as licensed shared access (LSA) [5]. Here, LSA licenses for the use of LTE spectrum are issued upon agreement for a specific geographical area and time duration required. Another type of LTE spectrum sharing actively investigated

within the LTE research community, is resource allocation in heterogeneous networks (HetNets) consisting of a number of LTE femto-cells overlaid by a high power macro-cell whilst sharing the same LTE channel, e.g. [6][7]. In these scenarios, the problem is often tackled by using game theory or machine learning principles.

One of the scenarios currently considered in the EU FP7 ABSOLUTE project is a temporary cognitive cellular infrastructure that is deployed in and around a stadium to provide extra capacity and coverage to the mobile subscribers and event organizers involved in a temporary event, e.g. a football match or a concert [8]. This scenario is depicted in Figure 1, where a small cell network is deployed inside the stadium to provide ultra high capacity density to the event attendees, and an eNodeB on an aerial platform can be deployed to provide wide area coverage, if required. In this particular study, we investigate how the stadium small cell network can share LTE spectrum with the local macro eNodeBs (eNBs) in the area, as a secondary system using cognitive DSA mechanisms. There is currently no evidence in the literature of investigating the feasibility of providing high capacity density using a cognitive cellular system with only secondary access to LTE spectrum.

An emerging state-of-the-art technique for intelligent DSA is reinforcement learning (RL); a machine learning technique aimed at building up solutions to decision problems only through trial-and-error [9]. The most widely used RL algorithm

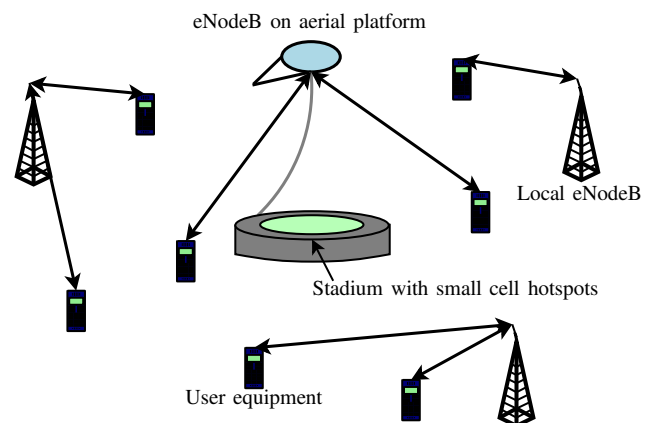


Figure 1. Example of a temporary event network which supplements the local cellular infrastructure

in both artificial intelligence and wireless communications domains is Q-learning. Therefore, most of the literature on RL based DSA focuses on Q-learning and its variations, e.g. [10][11]. This paper uses a distributed Q-learning based DSA algorithm, proposed in our previous work [12][13]. It has been shown to work effectively in scenarios where a cognitive cellular system has to prioritize among resources in its own dedicated spectrum. However, it has not been applied to problems where the spectrum is shared with incumbent systems.

The purpose of this paper is to assess the performance of distributed Q-learning based DSA in the cognitive cellular system with secondary LTE spectrum sharing scenario described above. We also aim to investigate the benefits of incorporating spectrum awareness and spectrum sensing based admission control into the DSA scheme employed by the cognitive network to aid coexistence between the primary and the secondary cellular systems.

The rest of the paper is organized as follows: Section II describes the stadium temporary event network scenario. Section III describes the distributed Q-learning based DSA algorithm, previously used only with a dedicated spectrum band. In Section IV we discuss several ways in which this DSA algorithm can be extended to facilitate coexistence between the primary and the secondary LTE systems. Simulation results are discussed in Section V. Finally, the conclusions are given in Section VI.

## II. TEMPORARY EVENT NETWORK SCENARIO

The cognitive cellular system investigated in this paper is designed for a stadium event scenario, where a small cell LTE network architecture is installed in a large stadium to provide high mobile data capacity to the users attending the event.

The network architecture is depicted in Figure 2, where the users are located in a circular spectator area 53.7 - 113.7 m from the centre of the stadium. The spectator area is covered by 78 eNBs arranged in three rings at 1 m height, e.g. with antennas attached to the backs of the seats or to the railings between the different row levels. Seat width is assumed to be 0.5 m, and the space between rows - 1.5 m, which yields the total capacity of 43,103 seats.

This cognitive small cell network has access to a 20 MHz LTE channel, also used by a network of 3 macro eNBs whose coordinates, with respect to the centre point of the stadium, are  $(-600, -750)$ ,  $(100, 750)$  and  $(750, -800)$  m. The users outside of the stadium are randomly distributed across a circular area with 1.5 km radius, i.e. the area covered by the macro eNBs. The small cell network is assumed to have secondary access to the spectrum, aiming to form a high capacity density cellular system by reusing the LTE spectrum of the primary macro eNB network around it.

## III. DISTRIBUTED Q-LEARNING BASED DYNAMIC SPECTRUM ACCESS

In pure distributed reinforcement learning based DSA the task of every eNB is to learn to prioritise among the available subchannels only through trial-and-error, with no frequency planning involved, and with no information exchange with

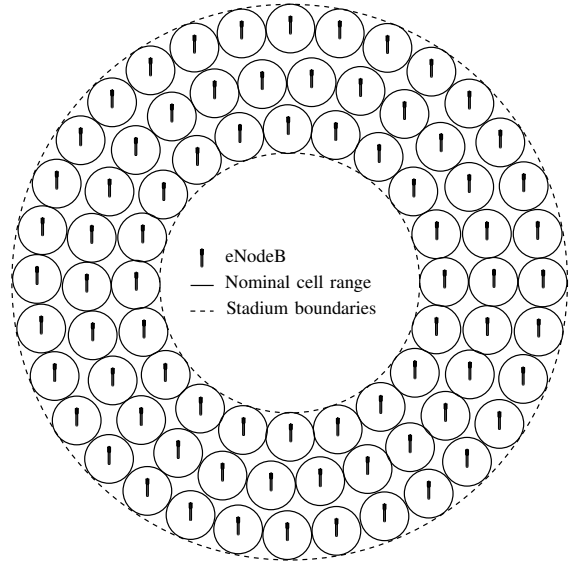


Figure 2. Stadium network architecture

other eNBs, e.g. [12]. In this way, frequency reuse patterns emerge autonomously using distributed artificial intelligence with no requirement for any prior knowledge of a given wireless environment.

### A. Reinforcement Learning

Reinforcement learning (RL) is a model-free type of machine learning which is aimed at learning the desirability of taking any available action in any state of the environment only through trial-and error [9]. This desirability of an action is represented by a numerical value known as the Q-value  $Q(s, a)$  - the expected cumulative reward for taking a particular action  $a$  in a particular state  $s$ .

The job of an RL algorithm is to estimate  $Q(s, a)$  for every action in every state, which are then stored in an array known as the Q-table. In some cases where an environment does not have to be represented by states, only the action space and a 1-dimensional Q-table  $Q(a)$  can be considered [14].

### B. Distributed Stateless Q-Learning

One of the most successful and widely used RL algorithms is Q-learning. In particular, a simple stateless variant of this algorithm, as formulated in [14], has been shown to be effective for several distributed DSA learning problems, e.g. [12][13][15].

Each eNB maintains a Q-table  $Q(a)$  such that every subchannel  $a$  has an expected reward or Q-value associated with it. The Q-value represents the desirability of assigning a particular subchannel to a file transmission. Upon each file arrival, the eNB either assigns a subchannel to its transmission or blocks it if all subchannels are occupied. It decides which subchannel to assign based on the current Q-table and the greedy action selection strategy described by the following equation:

$$\hat{a} = \underset{a}{\operatorname{argmax}}(Q(a)), a \in A', A' \subset A \quad (1)$$

where  $\hat{a}$  is the subchannel chosen for assignment out of the set of currently unoccupied subchannels  $A'$ ,  $Q(a)$  is the Q-value of subchannel  $a$ , and  $A$  is the full set of subchannels.

The Q-table is updated by the corresponding eNB each time it attempts to assign a subchannel to a file transmission in the form of a positive or a negative reinforcement. The update equation for stateless Q-learning, as defined in [14], is given below:

$$Q'(a) = (1 - \alpha)Q(a) + \alpha r \quad (2)$$

where  $Q(a)$  and  $Q'(a)$  represent the Q-value of the subchannel  $a$ , before and after the update respectively,  $r$  is the reward associated with the most recent trial and is determined by a reward function, and  $\alpha \in [0, 1]$  is the learning rate parameter which weights recent experience with respect to previous estimates of the Q-values.

The reward function returns one of the following values:

- $r = -1$  (negative reinforcement), if the transmission is blocked or interrupted due to low SINR on the selected subchannel.
- $r = 1$  (positive reinforcement), if SINR is above the allowed threshold throughout the whole transmission.

The choice of the learning rate value for this type of distributed Q-learning based DSA problems is thoroughly investigated in [13]. The best performance is achieved by using the Win-or-Learn-Fast (WoLF) variable learning rate principle described by Equation (3), where a lower value of  $\alpha$  is used for successful trials (when  $r = 1$ ), and a higher value of  $\alpha$  is used for failed trials ( $r = -1$ ). In this way, the learning agents are learning faster when “losing” and more slowly and cautiously when “winning”.

$$\alpha = \begin{cases} 0.01 & r = 1 \\ 0.05 & r = -1 \end{cases} \quad (3)$$

The values in the Q-tables are initialised to zero, so all eNBs start learning with equal choice among all available subchannels.

#### IV. COEXISTENCE WITH A PRIMARY LTE SYSTEM

Although, the distributed Q-learning algorithm described in the previous section has been shown to work effectively in self-organized cellular systems with dedicated spectrum, e.g. [12][13], it has not yet been applied in a scenario where the cognitive cellular system has to coexist with a primary network using the same spectrum. In this paper we assess three different methods of achieving this coexistence - a pure RL approach, a spectrum awareness database aided approach which requires external information about spectrum management of the primary network, and the cognitive radio type distributed spectrum sensing approach.

##### A. Reinforcement Learning

The RL approach does not require any modifications to the distributed Q-learning based DSA algorithm from the previous section. Its aim is to autonomously learn to avoid interference from both the primary and the secondary network at the same time. It achieves this through a simple, robust and widely applicable +1/-1 reward function described in Section III.

##### B. Spectrum Awareness

A standardized form of spectrum awareness (SA) used in LTE cellular systems is inter-cell interference coordination (ICIC), where the neighbouring eNBs frequently exchange short messages over the X2 interface, containing information which helps mitigate inter-cell interference between them [16]. For example, the format of ICIC signals in the LTE downlink is the Relative Narrowband Transmit Power (RNTP) indicator [17]. It consists of a bitmap which indicates on which resource blocks an eNB is planning to transmit at high power by setting their corresponding bits to 1, i.e. on which resource blocks it is likely to cause interference in adjacent cells. By exchanging such messages among neighbouring eNBs frequently, every eNB becomes “aware” of parts of the spectrum used in the neighbouring cells.

The SA based approach to facilitating secondary LTE spectrum sharing assumes that the cognitive cellular system has access to the ICIC signals of the eNBs from the primary system, thus providing a means to avoid interference between the two systems. We assume that the ICIC signals from the primary eNBs are received and processed by a simple SA server located at the stadium, which in turn broadcasts messages to all small cell eNBs stating which resource blocks have to be avoided, because they are used at high power by the primary system. This approach is similar to that employed in classical TV white space cognitive radio networks, which use geo-location databases to protect primary users of the spectrum, as well as increase the spectrum utilization efficiency of the cognitive radio systems, e.g. [4].

##### C. Spectrum Sensing

Figure 3 demonstrates how spectrum sensing functionality is normally embedded into RL based DSA algorithms in the context of cognitive radio [18]. It shows a flowchart of the Q-learning scheme described in the previous section, with a spectrum sensing module added to it. The additional functionality afforded by spectrum sensing is highlighted by shaded blocks and dotted lines.

Instead of simply assigning a subchannel with the highest Q-value in the Q-table, an eNB can check the interference level on it, and only assign it if it is below a certain allowed threshold. If interference on the selected subchannel is too high, an eNB can then check the interference on the next best subchannel and so on, whilst applying negative reinforcements to subchannels which fail these spectrum sensing checks using Equation (2). In this way, it is expected to improve the speed, due to increased number of negative reinforcements, and reliability of RL based DSA schemes. This scheme provides significantly more detailed and localised spectrum information than the spectrum awareness scheme, but is susceptible to the hidden terminal effect.

#### V. SIMULATION RESULTS AND DISCUSSION

This section presents the results of simulating the stadium temporary event scenario described in Section II. Four DSA schemes are applied to the secondary network:

- “Q-learning” - distributed Q-learning with no admission control from Section III.

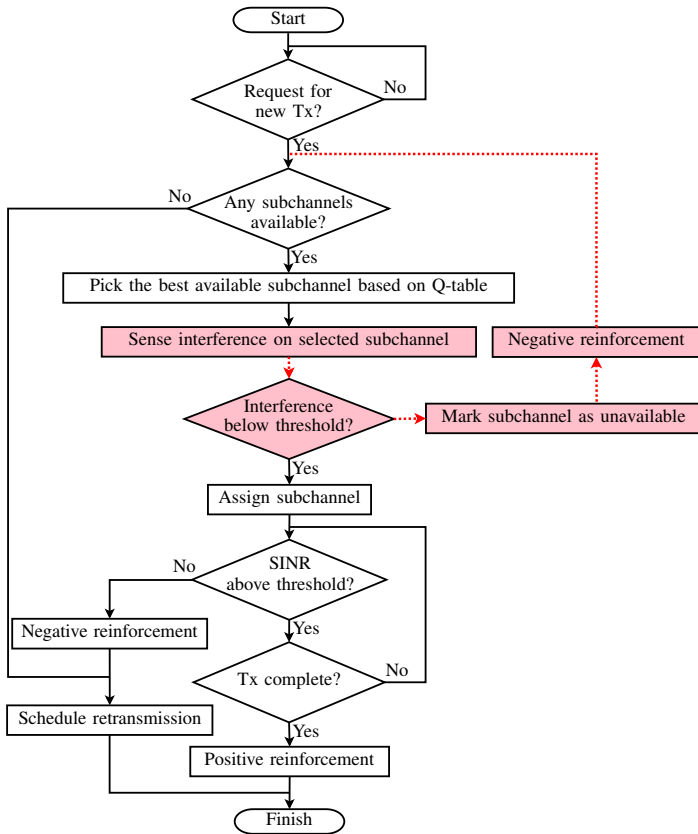


Figure 3. Flowchart of the distributed stateless Q-learning based DSA algorithm with spectrum sensing

- “Q-learning + spectrum awareness” - distributed Q-learning with spectrum awareness based admission control from Subsection IV-B, where all subchannels currently used by the primary network at 10 dBW Tx power are excluded from the available subchannel list of the secondary system.
- “Q-learning + spectrum sensing” - distributed Q-learning with spectrum sensing based admission control from Figure 3.
- “Spectrum sensing” - pure spectrum sensing based scheme, which uses the same admission control algorithm as the previous scheme, but with no Q-learning involved in it (achieved by setting  $\alpha$  to zero).

The interference threshold for the latter two schemes is 15 dB above the noise floor and 5 dB below received power at the UE receivers, based on a 5dB SINR admission threshold used in our previous work [12]. It was also found to achieve acceptable trade-off between deterioration in performance at low traffic loads (if higher threshold is used) and failure at higher traffic loads (if lower threshold is used).

The primary system is assumed to employ a dynamic ICIC scheme, where all three eNBs exchange their current spectrum usage as RNTP messages every 20 ms, and exclude the subchannels currently used by other two eNBs from their available subchannel list [16]. We assume that they always try to assign an available subchannel with the lowest index if any, e.g. they always scan the availability of the subchannels

in the same order from the 1st subchannel to the 25th. In this way, the primary network would make its spectrum usage more predictable for the cognitive cellular system, which is in the interests of both the primary and the secondary system.

### A. Simulation Setup

500 user equipments (UEs) are randomly distributed in the circular area from the stadium boundary (5 m from the radius of the last row) to 1.5 km away from the stadium centre point. 25% of the stadium capacity is filled with randomly distributed wireless subscribers, i.e.  $\approx 10,776$  UEs. The file arrival rates inside and outside the stadium are varied to obtain different offered traffic values, and all simulations last 1,000,000 transmissions. The other parameters and assumptions used in the simulation model are listed in Table I.

TABLE I. NETWORK MODEL PARAMETERS AND ASSUMPTIONS

Parameter	Value
Channel bandwidth	20 MHz (100 LTE physical resource blocks (PRBs))
Subchannel bandwidth	4 PRBs: $4 \times 180$ kHz [17]
Frequency band	2.6 GHz
UE receiver noise floor	94 dBm (290 K temperature, 20 MHz bandwidth, 7 dB noise figure)
Stadium propagation model	WINNER II B3 [19]
Outdoor propagation model	WINNER II C1 [19]
Propagation model between stadium and outdoors	Combined WINNER II C4 with C1 term [19]
Traffic model	3GPP FTP Traffic Model 1 [20], file size - 4.2 Mb ( $\approx 0.5$ MB)
Retransmission scheduling	Uniform random back-off between 0 and 960 ms [21]
Link model	3GPP Truncated Shannon Bound model [22]
Macro eNB Tx power	10 dBW
Assumptions	
Each UE is associated with an eNB with a minimum estimated downlink pathloss to it, based on the Reference Signal Received Power (RSRP)	
UEs inside the stadium are connected to the small cell network, UEs outside are connected to the macro eNBs	
Stadium network employs open loop power control, using a constant Rx power of -74 dBm (20 dB Signal-to-Noise-Ratio)	
The minimum Signal-to-Interference-plus-Noise (SINR) allowed to support data transmission is 1.8 dB [23]	
One subchannel (4 PRBs) is allocated to every data transmission	

### B. Simulation Results

Figure 4 shows the contour plots of the probability of retransmission ( $P(re-tx)$ ) of the secondary system, i.e. a ratio between the number of blocked/interrupted transmissions and the total number of transmissions, using four DSA schemes for the secondary system listed in the beginning of this section, at a wide range of traffic loads outside and inside the stadium.

Firstly, all schemes are affected by an increase in the traffic load of the primary system, which shows that the primary system produces harmful interference for the secondary system. Secondly, the distributed Q-learning scheme provides the best and most robust QoS in terms of  $P(re-tx)$ , compared to the other schemes. This difference becomes more significant, as the traffic load inside the stadium increases. Incorporating spectrum sensing or spectrum awareness into the Q-learning based DSA scheme, using the methodology described in Section IV, only deteriorates the performance of

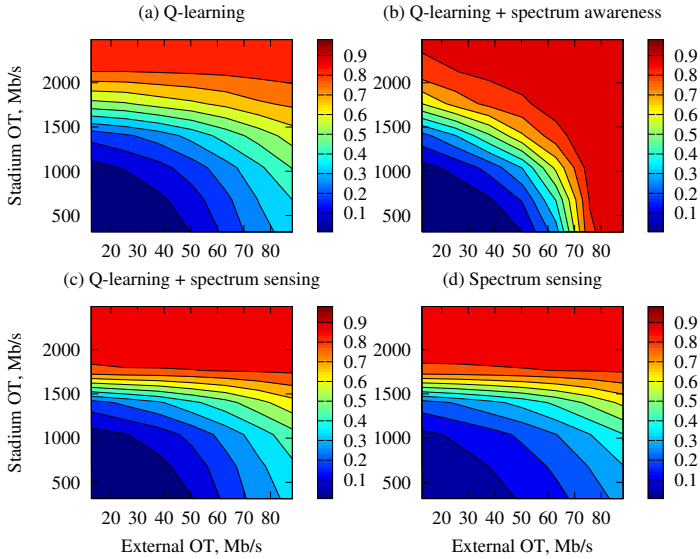


Figure 4. Probability of retransmission of the secondary system at different values of offered traffic (OT) outside (horizontal axis) and inside (vertical axis) the stadium

the secondary system by restricting the available resources, while Q-learning could still identify subchannels suitable for data transmission based on distributed machine intelligence.

Figure 5 shows that using spectrum awareness or spectrum sensing brings no extra benefit for protecting the primary system from interference. There,  $P(re - tx)$  of the primary system only depends on its own traffic load, and is independent from high traffic load variations inside the stadium. In fact, including spectrum awareness in the secondary system, at high traffic loads of the primary system, deteriorates the QoS in the latter. In these cases, most subchannels are occupied by the primary system, which forces all secondary system eNBs to choose among those few subchannels that are temporarily available, thus aggregating sufficient amount of interference on them to affect the primary system.

Figure 6 shows the secondary system throughput density, secondary system throughput divided by the area covered by the small cells, for the same set of simulations from Figures 4 and 5. These contour plots demonstrate that the Q-learning scheme achieves the highest system throughput density. The more centralised spectrum awareness approach has the poorest performance at higher primary system traffic loads, since it imposes a lot of restrictions on the spectrum usage of the secondary system. The fully distributed DSA schemes involving Q-learning and/or spectrum sensing are significantly more flexible in identifying spectrum reuse opportunities. Their performance is largely independent of the primary system, aided by the fact that interference between the two systems is attenuated by the stadium shell (captured by the propagation models listed in Table I), thus reducing the need for spectrum awareness based primary user protection.

Finally, Figure 7 compares secondary system throughput density at the maximum traffic load outside the stadium, when the whole LTE channel is in use by the primary system at most times. The secondary system manages to support 57 Gbps/km<sup>2</sup> using the distributed Q-learning algorithm from Section III,

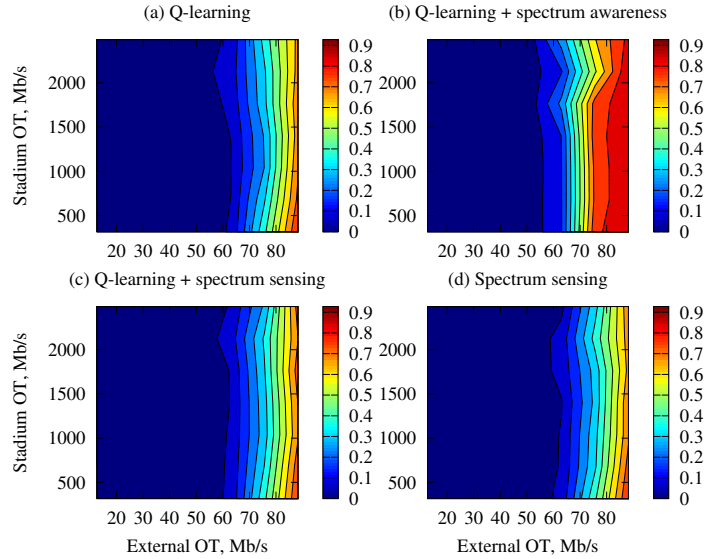


Figure 5. Probability of retransmission of the primary system at different values of offered traffic (OT) outside and inside the stadium

16% and 62% higher than that achieved if spectrum sensing or spectrum awareness are added respectively. Furthermore, as opposed to spectrum sensing and spectrum awareness based schemes, the Q-learning approach does not require any modifications to existing LTE network infrastructures and is easily implementable in commercially available small cell eNBs.

It must also be noted that the simulation scenario used in this study is more favourable towards the Q-learning approach. If both the primary and the secondary system were outdoors with no stadium shell attenuation and significantly more interference between them, the spectrum sensing or awareness based approaches would be more likely to contribute towards achieving coexistence between such networks. Another issue not discussed in this paper is the temporal performance of these

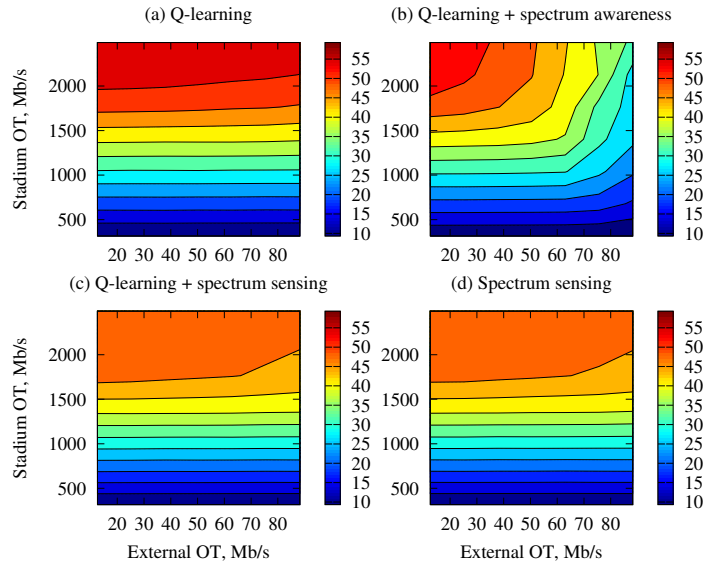


Figure 6. System throughput density (Gbps/km<sup>2</sup>) of the secondary system at different values of offered traffic (OT) outside and inside the stadium

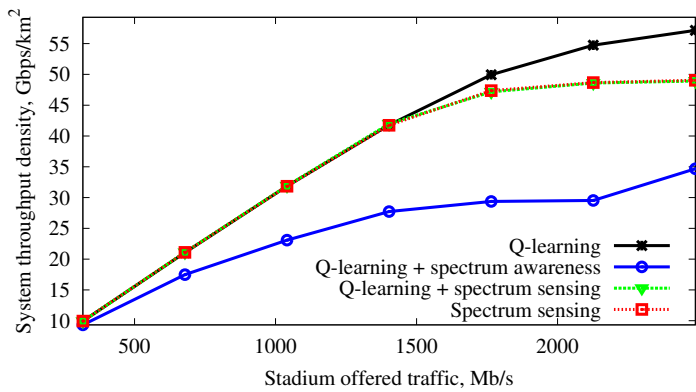


Figure 7. System throughput density of the secondary system at 88 Mb/s offered traffic outside the stadium, and at a range of traffic loads inside

DSA schemes. While pure Q-learning approach is known to have relatively poor performance at the start of the learning process and a significantly better steady-state performance [12], using spectrum sensing or spectrum awareness is likely to improve its initial behaviour and produce a more time-invariant QoS response. Both of these issues are the subject of our future research in this area.

## VI. CONCLUSION

In this paper we apply a distributed Q-learning based dynamic spectrum access (DSA) algorithm to a cognitive cellular system designed for providing ultra high capacity density with only secondary access to an LTE channel, simultaneously used by a primary network of macro eNodeBs. Large scale simulations of a stadium temporary event scenario show that the distributed Q-learning based DSA scheme provides robust quality of service (QoS) and extremely high system throughput densities ( $>55$  Gbps/km<sup>2</sup>) to the users of the stadium network, whilst successfully coexisting with the primary network on the same LTE channel. It is also shown that incorporating spectrum awareness or spectrum sensing based admission control into the Q-learning algorithm in our simulation scenario results in a decrease in QoS and system throughput density of the secondary network, with no effect on the primary system. Therefore, the simple distributed Q-learning based DSA algorithm, previously applied only in self-organizing cellular systems with dedicated spectrum, also provides an effective solution for spectrum sharing between high capacity density cognitive cellular systems and other LTE systems. Furthermore, it does not require any modifications to the current LTE standards and is easily implementable in commercially available small cell eNodeBs.

## ACKNOWLEDGMENT

This work has been funded by the ABSOLUTE Project (FP7-ICT-2011-8-318632), which receives funding from the 7th Framework Programme of the European Commission.

## REFERENCES

[1] H. Sun, A. Nallanathan, C.-X. Wang, and Y. Chen, "Wideband spectrum sensing for cognitive radio networks: a survey," *Wireless Communications, IEEE*, vol. 20, pp. 74–81, 2013.

[2] J. Sachs, I. Maric, and A. Goldsmith, "Cognitive cellular systems within the TV spectrum," in *IEEE Symposium on New Frontiers in Dynamic Spectrum*, 2010.

[3] C. Ghosh, S. Roy, and D. Cavalcanti, "Coexistence challenges for heterogeneous cognitive wireless networks in TV white spaces," *Wireless Communications, IEEE*, vol. 18, pp. 22–31, 2011.

[4] D. Gurney, G. Buchwald, L. Ecklund, S. Kuffner, and J. Grosspietsch, "Geo-location database techniques for incumbent protection in the TV white space," in *IEEE Symposium on New Frontiers in Dynamic Spectrum Access Networks (DySPAN)*, 2008.

[5] M. Matinmikko, H. Okkonen, M. Palola, S. Yrjola, P. Ahokangas, and M. Mustonen, "Spectrum sharing using licensed shared access: the concept and its workflow for LTE-advanced networks," *Wireless Communications, IEEE*, vol. 21, pp. 72–79, 2014.

[6] S. Hamouda, M. Zitoun, and S. TABBANE, "Win-win relationship between macrocell and femtocells for spectrum sharing in LTE-A," *Communications, IET*, vol. 8, pp. 1109–1116, 2014.

[7] G. Alnwaimi, T. Zahir, S. Vahid, and K. Moessner, "Machine Learning based Knowledge Acquisition on Spectrum Usage for LTE Femtocells," in *IEEE Vehicular Technology Conference (VTC-Fall)*, 2013.

[8] L. Reynaud, et al., "FP7-ICT-2011-8-318632-ABSOLUTE/D2.1 Use cases definition and scenarios description," 2014.

[9] R. Sutton and A. Barto, *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 1998.

[10] X. Chen, Z. Zhao, and H. Zhang, "Stochastic power adaptation with multiagent reinforcement learning for cognitive wireless mesh networks," *Mobile Computing, IEEE Transactions on*, vol. 12, pp. 2155–2166, 2013.

[11] J. Nie and S. Haykin, "A Q-learning-based dynamic channel assignment technique for mobile communication systems," *Vehicular Technology, IEEE Transactions on*, vol. 48, pp. 1676–1687, 1999.

[12] N. Morozs, D. Grace, and T. Clarke, "Case-based reinforcement learning for cognitive spectrum assignment in cellular networks with dynamic topologies," in *Military Communications and Information Systems Conference (MCC)*, 2013.

[13] N. Morozs, T. Clarke, D. Grace, and Q. Zhao, "Distributed Q-learning based dynamic spectrum management in cognitive cellular systems: Choosing the right learning rate," in *IEEE International Symposium on Computers and Communications (ISCC)*, 2014.

[14] C. Claus and C. Boutilier, "The dynamics of reinforcement learning in cooperative multiagent systems," in *Proceedings of the fifteenth national/tenth conference on Artificial intelligence/Innovative applications of artificial intelligence*, 1998.

[15] Q. Zhao, T. Jiang, N. Morozs, D. Grace, and T. Clarke, "Transfer learning: A paradigm for dynamic spectrum and topology management in flexible architectures," in *IEEE Vehicular Technology Conference (VTC Fall)*, 2013.

[16] S. Sesia, M. Baker, and I. Toufik, *LTE-The UMTS Long Term Evolution: From Theory to Practice*. John Wiley & Sons, 2011.

[17] 3GPP, "LTE; Evolved Universal Terrestrial Radio Access (E-UTRA); Physical layer procedures (3GPP TS 36.213 version 11.5.0 Release 11)," Dec. 2013.

[18] T. Jiang, D. Grace, and P. D. Mitchell, "Efficient exploration in reinforcement learning-based cognitive radio spectrum sharing," *Communications, IET*, vol. 5, pp. 1309–1317, 2011.

[19] P. Kyösti, et al., "IST-4-027756 WINNER II Deliverable D1.1.2: WINNER II channel models," 2008.

[20] 3GPP, "Technical Specification Group Radio Access Network; Evolved Universal Terrestrial Radio Access (E-UTRA); Further Advancements for E-UTRA physical layer aspects (3GPP TR 36.814 version 9.0.0 Release 9)," Dec. 2010.

[21] —, "LTE; Evolved Universal Terrestrial Radio Access (E-UTRA); Medium Access Control (MAC) protocol specification (3GPP TS 36.321 version 11.4.0 Release 11)," Jan. 2014.

[22] —, "LTE; Evolved Universal Terrestrial Radio Access (E-UTRA); Radio Frequency (RF) system scenarios (3GPP TR 36.952 version 11.0.0 Release 11)," Dec. 2012.

[23] T. Jiang, et al., "EU FP7 INFSO-ICT-248267 BuNGee Deliverable D4.1.2: Simulation Tool(s) and Simulation Results," 2012.