

# Cognitive Cellular Networks: A Q-Learning Framework for Self-Organizing Networks

Stephen S. Mwanje, Lars Christoph Schmelz, and Andreas Mitschele-Thiel

**Abstract**—Self-Organizing Networks (SON) aim at simplifying Network Management (NM) and optimizing network capital and operational expenditure through automation. Most SON Functions (SFs) are rule-based control structures which evaluate metrics and decide actions based on a set of rules. These rigid structures are, however, very complex to design since rules must be derived for each SF in each possible scenario. In practice rules only support generic behavior which cannot respond to the specific scenarios in each network or cell. Moreover, SON coordination becomes very complicated with such varied control structures. In this paper, we propose to advance SON towards Cognitive Cellular Networks (CCN) by adding cognition that enables the SFs to independently learn the required optimal configurations. We propose a generalized Q-learning framework for the CCN functions and show how the framework fits to a general SF control loop. We then apply this framework to two functions on Mobility Robustness Optimization (MRO) and Mobility Load Balancing (MLB). Our results show that the MRO function learns to optimize handover performance while the MLB function learns to distribute instantaneous load among cells.

**Index Terms**—SON; Cognitive Cellular Networks; MRO; MLB

## I. INTRODUCTION

OPTIMAL cell sizes in cellular networks are continuously decreasing, increasing the number and density of cells especially with the introduction of LTE. This has resulted in higher capital and operational expenses, and increased complexity of network operation. Self-Organizing Networks (SON) promise to minimize these challenges through automation of network operations. SON Functions (SFs) have been defined, e.g., in the LTE SON standard [1], with each SF representing a function that can be automated. Examples are Mobility Robustness Optimization (MRO), Mobility Load Balancing (MLB), Coverage and Capacity Optimization (CCO) or Inter-Cell Interference Coordination (ICIC). These functions have traditionally been developed as rule-based controllers which require the designers to have full understanding of the function's behavior. We propose to advance SFs to CCN functions and implement them as cognitive, Q-Learning (QL) based agents that act in the network and use the network's feedback to learn the effects of their actions.

The paper is organized as follows: Section II briefly summarizes the related works, Section III discusses QL and its application to CCN functions while Section III describes the

simulation scenario and environment. We apply the CCN concept to two candidate functions on MRO and MLB presented in sections IV and V, and conclude with a summary of the presented work in Section VI.

## II. SON FUNCTIONS: THE STATE OF THE ART

Each SF is characterized by a trigger that initiates the execution of an associated SON algorithm, which in turn configures a set of network parameters in order to optimize a particular metric. The basic SF (in Fig. 1a) is a control agent that: 1) observes the network to evaluate trigger conditions, 2) takes an action to optimize its metrics and 3) gets feedback on the effect of that action on the network. Accordingly, most SFs have been developed as rule-based controllers which select actions by applying defined rules on the observed metrics. Examples controllers are [2]–[5] for Handover (HO) optimization; [6], [7] for tilt optimization and [8], [9] for load balancing. Such an implementation, however, requires that the rules include all the possible scenarios, and that the rule designer fully understands the effects of each action in each of these scenarios. In reality, this is not possible even for a system expert. The biggest challenge is that these control loops create rigid structures that are very complex to design and hard to evolve as more functions are deployed in the network. Furthermore, they only allow for generic behavior which cannot respond to specific contexts in each network or cell.

To counter these challenges, we advance SON to Cognitive Cellular Networks (CCN) by adding cognition to SFs. The CCN concept advances SON beyond the rule based control loops towards artificial-intelligence based cognitive functions that autonomously learn the optimal configurations. Specifically, we propose to design CCN functions as QL agents that act and, using network's feedback, learn from the effects of their actions. QL has been applied in some SFs with positive results, e.g., in [10] [11] [12], but we advance this and propose a QL framework that is applicable to all optimization functions. We demonstrate the benefits of applying the framework with extended versions of the results in [13], [14].

## III. Q-LEARNING FRAMEWORK FOR SON FUNCTIONS

### A. Q-Learning (QL)

Multiple approaches have been used to develop SFs. Reinforcement Learning (RL), mainly QL [15], offers the best promise as it allows the network to learn and improve its solutions through experience. QL is a model-free RL algorithm which, using Temporal Difference (TD), solves learning problems even without models. A QL problem is the triple

Stephen Mwanje and L. C. Schmelz are with the Network Management Automation Research team at Nokia research, Munich, Germany (e-mail: stephen.mwanje; christoph.schmelz@nokia.com)

A. Mitschele-Thiel heads the Integrated Communications Group at the Technische Universität Ilmenau, Germany (e-mail: mitsch@tu-ilmenau.de)

Manuscript received October 8, 2015; revised January 14, 2016.





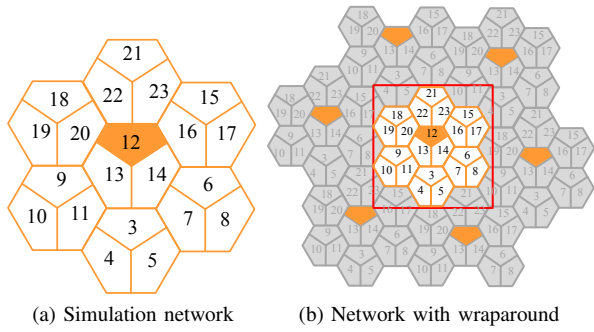


Fig. 3. Simulation network structure with and without wraparound

estimates  $N_B$  the number of bits that can be carried per symbol using realistic Modulation and Coding Schemes (MCS) as:

$$N_B = \begin{cases} 0; & S < 7.04 \\ -0.0001S^3 + 0.0074S^2 + 0.1397S + 0.6218; & -7.04 < S < 20.2 \\ N_B(20.2); & S > 20.2 \end{cases} \quad (4)$$

We consider the Guaranteed Bit Rate (GBR) traffic model where, as a minimum, a user must be allocated resources ensuring it achieves the specified rate. Then, the number of Physical Resource Blocks (PRBs),  $N_{PRB}$ , required for data transmission per scheduling interval  $T_s$  will be the ratio of the total data to be transmitted during  $T_s$  to the user's achievable rate per PRB. Consequently, the cell's offered load  $\rho$  is the ratio of the required PRB count for all the cell's users to the number of available PRBs. Given that the PRB has 7 symbols in a bandwidth  $B_{PRB}$  of 180 kHz [24], the offered load is

$$\rho = \frac{B_{PRB}}{B_{sys}} \cdot \sum \frac{GBR \cdot T_s}{7N_B} \quad (5)$$

Cell overload occurs when the offered load  $\rho$  exceeds a preset threshold  $\rho_{max}$ . Accordingly,  $\rho$  can be greater than 1, which represents the case where the total required PRBs exceed the available maximum PRBs within  $B_{sys}$ . Each cell allocates PRBs to its associated UEs using a simple round robin scheduler. The scheduler continuously allocates resources to users ensuring that each achieves the desired rate before allocating the next user, and this continues either until all resources are exhausted or until all users are allocated.

### B. Mobility and Handover Management

HOs between a serving cell  $s$  and a target cell  $t$  are triggered according to the A3 condition [26], which, using the Reference Signal Received Power (RSRP) for HO from  $s$  towards  $t$ , is

$$F_t + O_t^{s,t} - Hys > F_s + O_s^{s,t}. \quad (6)$$

$F_t, F_s$  are respectively the user's RSRP in dBm in  $t$  and  $s$  cells, without any offsets;  $O_t^{s,t}, O_s^{t,s}$  are the respective Cell Individual Offsets (CIOs) while Hys, which is uniform for the serving cell, is the Handover Hysteresis (Hys) in dB.

If A3 is fulfilled for a critical time called Time To Trigger (TTT), the UE initiates HO by sending a measurement report of the values  $F_s$  and  $F_t$  after being filtered by a Layer 1 averaging filter (L1) and a Layer 3 Infinite Impulse Response (IIR) filter (L3). L3 is implemented as specified in the LTE

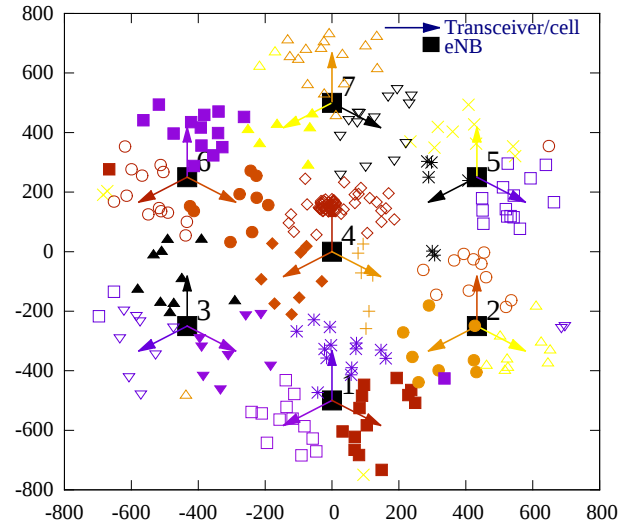


Fig. 4. Initial Network with deployed users colored according to serving cell

Radio Resource Control (RRC) specifications [26] with a filter coefficient of 6. Meanwhile, L1 averaging which is not standardized and can be vendor-specific, is implemented as a simple moving average of the most recent values. The averaging window, which must be updated every 200ms and is required to be long enough to average fast fading yet not so long to affect the results of L3, is selected to be 100 samples.

It is evident that HO outcomes mostly depend on Hys and TTT both having the effects of either delaying or advancing HOs. The possible outcomes are either a HO success, a Ping-Pong HO (PP) or a Radio Link Failure (RLF). These are modeled according to the timers defined by 3GPP in [26].

### C. Simulation Model

The simulation starts with cell deployment. It then executes multiple "batches" starting each with user deployment. This redeployment ensures that users experience as many of the prevalent radio conditions as possible, since in each batch, each user is placed at a different location and follows a different path during the execution of the batch.

Mobile devices are deployed following a random distribution, in a way that the average number per cell is the ratio of the total number of users to the number of cells. For MLB studies, a hot-spot-induced overload is artificially added by deploying static users (i.e. UEs with velocity = 0 m/s) in a "center cell", specifically, cell 12 in Figure 3a. An example distribution of the users after deployment is shown in Fig. 4.

Test studies showed that statistics for HO events are stable after at most 80 batches of 200 s each. MRO related results are therefore based on simulations of 120 batches each simulating 200 s of operation. Load statistics are, however, stable after at most 20 batches, and so MLB results are based on simulations of 30 batches each simulating 200 s of operation. The crucial simulation parameters are summarized in Table II. The next sections will show performance results of the two CCN based SFs when simulated with the described simulation tool.



here however, that such do not exist or their effects are so small and can be ignored. This is a justified assumption based on test results which showed that if HO are triggered adequately early ( $Hys=0$  dB and  $TTT=0$  s), RLFs will be eliminated (with excessively high PPs).

- 4) The selection of weights  $w_i$  is subjective. They have here been selected so as to equally balance effects of early HO ( $P, F_E$ ) against effects of late HO ( $F_L$ ), i.e.,  $w_3$  equal to  $w_1$  and  $w_2$  combined. Then, since RLFs are less desirable compared to PPs,  $w_2$  should be larger than  $w_1$ . The correspondingly selected weight vector is  $w = (0.2, 0.3, 0.5)$  and is the one used in all cases where HO performance is evaluated.

### B. HO Control Parameters Sensitivity

The core MRO goal is to dynamically select the optimum settings (OTP) even for a network with a dynamic mobility profile. To design an adaptive learning strategy, we investigate the sensitivity of the parameters to UE velocity. We do so by sweeping a selected range of the parameter space for four velocity scenarios. With UEs moving at constant velocity in each scenario, we observe that the OTP changes with velocity as shown in Fig. 5. Fig. 5a gives the linear variation of the HOAP with both Hys and TTT while Fig. 5b describes the detailed variation with TTT using a TTT log scale.

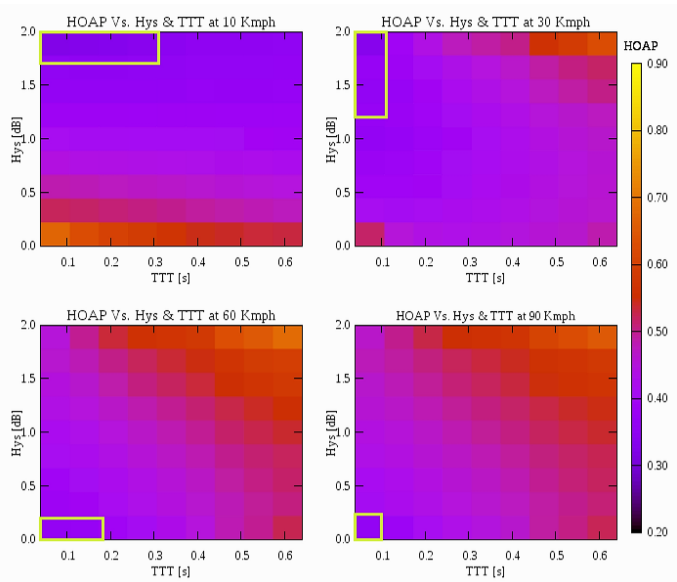
We observe in Fig. 5a that a very high Hys is unacceptable at all velocities, although a combination of moderately high Hys and low TTT could be acceptable. Similarly, a high TTT is only acceptable at low velocities and only in combination with low to medium Hys. Even then, the best settings at low velocity should be medium Hys with low-to-medium TTT, i.e., HO can moderately be delayed without great penalty since the risk of RLF is low yet even the possibility of PPs is low owing to the low velocity. This is evident in the 10 kmph case where for most TTTs the performance is good at  $Hys = 2$  dB.

As the velocity increases, the HO delay needs to reduce especially using the TTT. The HOAP is more susceptible to change in TTT, to the extent that the OTP continuously grazes the Hys axis, i.e., the performance changes with TTT but is fairly constant with Hys. At high velocity, even the Hys has major effect and so both parameters should be low. This is evident in the 60 and 90 kmph environments in Fig. 5a where the OTPs are restricted to the lower left corners of the grid, i.e., the part where TTT are within the range of 0-0.64 s. In Fig. 5b we observe that within this small range, although there is major variation in the HOAP with TTT for most Hys, this variation is blurred at points near the optimum point. In that case adjacent TTTs will have practically similar performance.

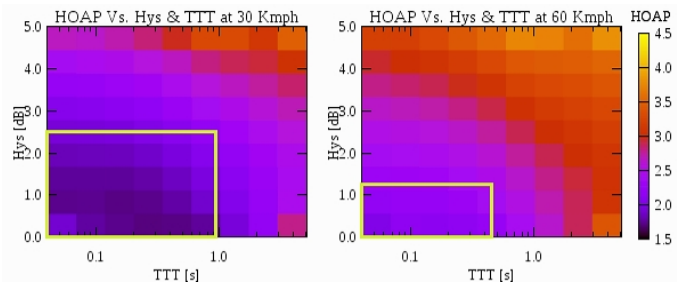
The most obvious conclusion from Fig. 5 is that the OTPs do not lie along any one diagonal for the different velocities as was assumed in [28] and [29]. Any MRO algorithm must thus scan the entire parameter-space or at least more than half the space in order to determine the required trigger point.

### C. QMRO: Q-Learning based MRO

QMRO wishes to determine the optimum  $Hys$ - $TTT$  action that minimizes the HOAP in any mobility state in a cell.



(a) Linear variation of HOAP with Hys and TTT



(b) Detailed (log scale) variation of HOAP with TTT for 30 and 60 kmph

Fig. 5. Handover Control Parameter Sensitivity

We note here that the actions only affect the performance of the cells and do not change the UEs' mobility states. Consequently, it is adequate to learn an action  $a$  in state  $x$  that maximizes the expected instantaneous reward  $r$  at time  $t$ . As derived in Equation 2, the corresponding Q-update algorithm for instantaneous rewards is Equation (8)

$$Q_{t+1}(x_t, a_t) = (1 - \alpha)Q_t(x_t, a_t) + \alpha[r_t(x_t, a_t)] \quad (8)$$

where  $\alpha$  is the learning rate previously defined in section III-A. The components of the QMRO algorithm are described below.

1) *QMRO State Space*: The required HO settings in a cell depend on the mobility of the UEs in the cell as showed in the parameter sensitivity analysis in section V-B. Thus, the states  $x$  are defined to be the degree of mobility in the cell evaluated as the average velocity over the *SON interval*. The average velocity being a continuous variable, we discretize the states  $x$  into bands for which the appropriate HO settings must be learned. Table III describes how velocities are grouped into mobility states and, based on the results in Fig. 5, the estimates of the likely default settings that would be used in a manual optimization process. We assume that the velocities are known or at least can be estimated by the cells. A simple estimate can be obtained as the ratio of the approximate cell size to the average time that a UE stays in the cell. It can also be more accurately estimated using the UEs Doppler power spectrum as proposed in [32]. Either way, with a good estimate, QMRO

TABLE III  
QMRO MOBILITY STATES AND THEIR DEFAULT ACTIONS

Average Velocity (kmph)	State (x)	Default Hys (dB)	Default TTT (s)
0-4	0	3.0	0.0-5.2
4-8	1	2.5	0.0-2.56
8-12	2	2.0	0.0-1.25
12-17	3	2.0	0.0-1.02
17-22	4	1.5-2.0	0.0-0.64
22-28	5	1.5-2.0	0.0-0.48
28-34	6	1.5-2.0	0.0-0.256
34-41	7	1.5	0.0-0.52
41-48	8	1.0	0.0-0.52
48-56	9	0.5	0-0.48
56-65	10	0.5	0.0-0.256
65-75	11	0.0-0.5	0.0-0.16
75+	12	0.0-0.5	0.0-0.16

can then learn the best configuration for the given cell.

2) *QMRO Action Space*: Actions are the Hys-TTT tuples signaled by the cells to their associated UEs. Without a SON solution, an operator configures a cell with default parameter settings obtained through trial and error, while with a local search based SON solution, a fixed set of settings similar to those shown in Table III are applied. However, with the observation that OTPs depend on speed, we need to change the settings based on the instantaneous speed in the cell.

It is evident from the parameter-sweep results in Fig. 5 that for all practical speeds, performance at Hys > 5 dB is almost always sub optimal. We thus consider Hys values only up to 6 dB. Meanwhile differences in HOAP for some TTT settings are unresolvable especially at low TTT values. For example for most Hys values at all velocities, the performance at TTT = 0.08s, 0.1s, 0.16s is practically the same. As such not all TTT values are considered, i.e., the possible TTT actions are the 11 values 0.04, 0.10, 0.128, 0.256, 0.32, 0.48, 0.512, 0.64, 1.02, 1.28, 2.56, 5.12 in s. The resulting action space (total number of actions) for each state is 143 possible combinations of the considered Hys and TTT.

3) *QMRO Reward function*: We desire to minimize RLFs without excessively increasing PPs and HOs. Since the learner is a rewards-maximizing agent, the reward  $r_{x,t}$  should be the negative HOAP evaluated over the SON interval. As stated earlier, the individual rates are normalized to the number of HO candidates, NH as given in Equation 9.

$$r_{x,t} = -(w_1P + w_2F_E + w_3F_L)/NH; \quad (9)$$

Meanwhile, the weight vector applied during learning may need to be adjusted to enforce particular results especially given the small evaluation period (the SON interval). For example, there may be instances in which no RLFs are observed during the SON interval which could tilt the result in favor of too many PPs. The vector is thus maintained as  $w = (0.2, 0.3, 0.5)$  for the typical results and changed to  $w = (0.4, 0.0, 0.6)$  when no RLFs are observed.

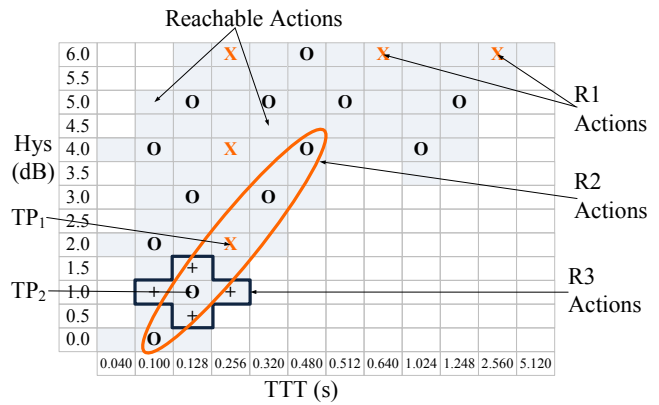


Fig. 6. QHO action space within the three learning regimes R1, R2, R3

4) *QMRO Cooperative Learning*: HO triggering could be affected by channel conditions that dictate the respective RSRP. However using  $A_3$  minimizes this dependence on the absolute RSRP values, since decisions are made based on RSRP differences among the cells. In that case, HO performance depends only on user mobility and the control parameter values. With mobility-based HO states, it is possible that a state observed in one cell reoccurs in another cell at some other time. As such, cells do not need to learn independent policies but can learn a single policy function based on the abstract mobility states. The result is a cooperative QL problem in which individual cells take actions but update a single Q-table that represents the shared learned policy.

The cooperative learning solution holds if all other crucial parameters are comparable among the cells. For example the cells in the considered network are assumed to be of similar size and applying comparable transmit powers. Other than this, the RSRP profiles at the cell edges may be different resulting in differing behaviors for different cells or cell pairs. Similarly if individual cells concurrently have users with differing behavioral patterns, the solution may fail. For example, the assumption that a state in one cell will be observed in another will not hold for a cell which covers a highway crossing through an office park. Such a cell concurrently has 2 user groups - the slow-moving office users and the fast-moving highway users, each of which will require different settings. Nevertheless, for the majority of networks that do not have such special conditions, cooperative learning as considered in these studies would be fully applicable.

#### D. Parameter Search Strategy and Optimization Algorithm

As earlier stated, each cell has up to 143 possible actions to consider for each velocity state. Thus even with cooperative learning, evaluating each action multiple times, would still require a long time to converge to the desired solutions.

1) *The Parameter Search Strategy*: To accelerate convergence, we subgroup the 143 actions so that for any state, we execute 3 learning regimes  $R1 - R3$  (shown in Fig. 6) that start with a global exploration and gradually move to a local. For  $R1$ , actions are selected from different regions of the grid in order to determine the area in which the desired action lies. From the parameter sensitivity analysis, combinations of low

Hys and high TTT are never optimal. As such this region is excluded from the possible candidates denoted by "R1 actions" in Fig. 6. The outcome trigger point of R1 ( $TP_1$ ) specifies the region in which the optimum point lies. This outcome is thus used to define the search space for the next regime R2.

During R2, actions along the diagonal that goes through  $TP_1$  are explored to obtain the approximate delay that is acceptable for the observed mobility state. In this case, subsequent actions differ in Hys by 1 dB to enable a large enough action space to be explored. As an example, if  $TP_1$  is obtained as  $TP_1 = (2.0dB, 0.256s)$ , at R2 the agent explores the region marked R2 Actions in Fig. 6. The obtained TP ( $TP_2$ ) is then used to define the search space for the next regime R3.

R3 refines the learned  $TP_2$  by exploring points near  $TP_2$ . It compares  $TP_2$  with its four neighbor points to the left, right, top or below. In Fig. 6, 'R3 Actions' shows the exploration region for R3 assuming that  $TP_2 = (5.0dB, 0.128s)$ .

2) *The SON Interval*: Each setting that is applied in a cell is monitored over a period of a SON interval. With different number of users in each cell and in the HO regions, different cells may have different counts of events within the same time period. Consequently, instead of setting the SON interval based on a fixed time period, it is based on a minimum number of HO events that must occur following the application of any action or configuration. This minimum number is the sum of the disjoint HO related events (i.e. HOs, RLFs and RLFLs). This number is set to 100 events although any value that ensures that comparable counts of all the necessary statistics (i.e. for NHs, PPs and RLFs) are observed would be appropriate.

3) *The QMRO Optimization Algorithm*: Given the QL elements as discussed above, the optimization algorithm is given in the procedure of Algorithm 2, i.e: For each possible state, the action set is initialized with R1 actions and the Q-table entries initialized to 0. Learning is then triggered to be executed after every SON interval  $t$ . Each cell  $c$  observes its environment over the interval  $t$  and at the end of  $t$ , the cell determines if an optimization is necessary, i.e., if the cell's velocity state has changed. During the learning phase,  $c$  selects an action as described in section V-D1, otherwise it selects the best action that would have been learned. It then signals that action to all its associated UEs and starts collecting the necessary performance statistics for the next interval( $t + 1$ ). At the end of interval  $t + 1$ ,  $c$  evaluates its HOAP and derives the reward  $r_t$  for the action at  $t$ . It then updates the learning agent (the Q-table) before repeating the process.

### E. Simulation Results and Discussion

MRO wishes to adjust the HO parameters in line with varying mobility states. We evaluate performance in the 5 different velocity scenarios in Table IV to prove that the algorithm is applicable to any network. The 3 'normal' scenarios (10, 30, 60 kmph), for example, represent 3 typical city districts - a city center, city edge and residential suburb. We then consider two extreme scenarios (3 & 120 kmph) which could respectively represent an office park and a highway.

In each mobility scenario, all UEs have independent randomly varying velocities. We implement this by allocating

### Algorithm 2: QMRO - The Q-Learning MRO Algorithm

Require: UE velocities during SON interval, action set  $A$

1. Set  $R_i=R1$ ; initialize action set  $A_{x,R1}$  for regime 1 in all states  $x$
- Repeat for each SON interval  $t$**
2. **if** HO action was taken at SON interval  $t - 1$  **do**
- determine HOAP and derive reward  $r_{t-1}(x_{t-1}, a_{t-1})$
3. update Q-table according to Equation 8
- end if**
4. determine current mobility state  $x_t$  (from table III)
5. **if** learning complete for state  $x$  **do**
- select  $a_{x,t} = a_x^{opt}$ , the best action for state  $x$
6. **else if** regime  $R_i$  exploration is incomplete **do**
- select  $a_{x,t}$  (sequentially after  $a_{x,t-1}$ ) from  $A_{x,R_i}$
7. **else do**
- select  $a_{x,t} = a_{x,R_i}^{opt}$ , the optimum value for state  $x$  at  $R_i$
8. **if** all learning regimes complete for state  $x$  **do**
- record all regimes complete for state  $x$
9. record  $a_{x,t}$  as best action in state  $x$
10. end learning, indefinitely use  $a_{x,t}$  in state  $x$
11. **else do**
- $R_i \leftarrow R_i + 1$
12. use  $a_{x,t}$  to set  $A_{x,R_i}$  i.e. reconfigure  $A$  for  $R_i$
13. **end if**
14. **end if**
15. Signal selected action  $a_{x,t}$  to all UEs in the cell.
16.  $t \leftarrow t + 1$ , monitor, collect statistics, continue at step 2
17. **end loop**

TABLE IV  
QMRO VELOCITY SCENARIOS

City Area	Initial velocity (kmph)	
	Mean	range
Office park	3	2 - 4
City Center	10	6 - 14
City Edge	30	18 - 42
City Suburb	60	36 - 84
Highway	120	72 - 168

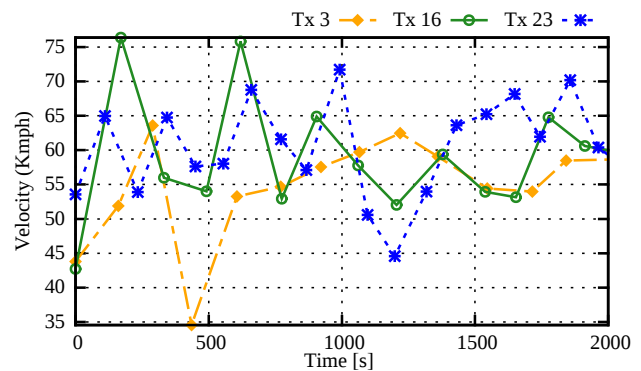


Fig. 7. Typical velocity profile in 3 cells in the 60kmph scenario

random velocities to the UEs at the start of the simulation and also randomly adjusting the velocities at the start of and during every batch, in each by up to  $\pm 40\%$ . For example, each of the 240 users in the city suburb (60 kmph) network has an individually assigned and continuously changing velocity, as shown in Fig. 7 for the average velocities in three selected cells over a period of 10 batches.

1) *Performance in terms of HOAP*: To evaluate QMRO, we compare its performance in each velocity scenario against the reference network, Ref. This represents the case when all the



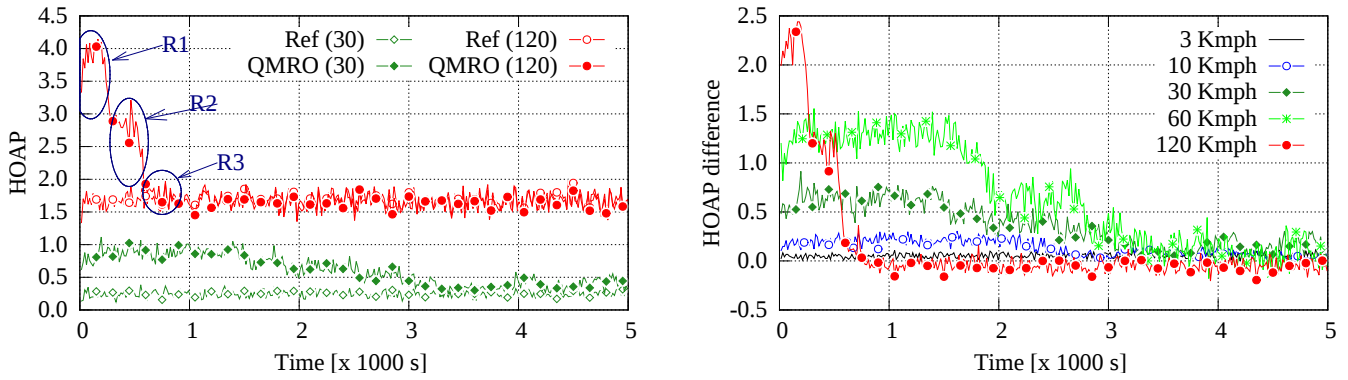


Fig. 8. QMRO Performance: Average network-wide HOAP for QMRO in comparison to the reference network

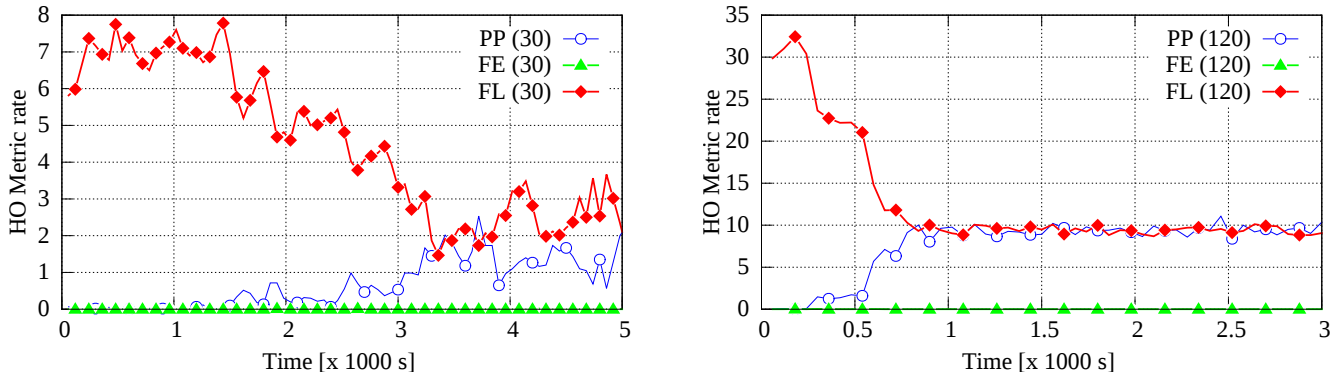


Fig. 9. QMRO Performance: Variation of average rates for all 3 core metrics in 30 and 120 kmph environment

cells in the network apply the best static settings as obtained from table III in the parameters sensitivity analysis in section V-B. The performance is evaluated in terms of the averages of the metric(s) values throughout the network although cell-specific results would demonstrate the same trends.

Fig. 8 summarizes the results in the different scenarios. Fig. 8a shows the comparison of the average HOAP values for QMRO and Ref in two typical cases, while Fig. 8b describes the differences in performance between QMRO and Ref for all the five velocity cases. We observe in both figures that QMRO initially performs poorly as it executes the first learning regime *R1*. The performance then improves in regimes *R2* and *R3* as QMRO focuses around the *OTP*, to the extent that it is eventually equivalent to that of Ref. Where user velocities are widely spread, QMRO actually performs better since it is able to set the right setting for each velocity range as opposed to a single setting for all velocities. This is evident for the 120 kmph case in Fig. 8b where, after learning, QMRO is consistently better than Ref. Also, in a given time interval, each user undertakes more HOs as the velocity increases, with the effect that cells have a shorter SON interval at a higher velocity i.e the cells reach the minimum event count much faster. This results into shorter convergence times as the velocity increases which is evident in Fig. 8b.

2) *QMRO Learning Trend*: Fig. 9 evaluates the time variation of the individual metrics ( $P$ ,  $F_E$  and  $F_L$ ) when applying QMRO in two velocity scenarios (30 and 120 kmph). We observe in both cases, that the agent learns to minimize RLFs

( $F_L$ ) by trading them with Ping-pongs ( $PPs$ ), which have less effect on the user's quality of experience. This is all while ensuring that RLFs ( $F_E$ ) remain low. In both velocity scenarios, QMRO suffers from many RLFs at the beginning as it considers settings across a large parameter space. Over time however, the agent continuously reduces  $F_L$  by trading such reduction with increase in  $PP$ . It then stops this trend as soon as it registers decreasing returns, i.e., when each extra reduction in  $F_L$  translates into an excessive increase in  $PPs$  or if it instead causes RLFs to occur.

The foregoing results demonstrate that given a good definition of states that appropriately capture the UEs' mobility and also given adequate learning time, a Cognitive function, e.g., the QL based MRO algorithm is able to learn the appropriate Hys-TTT settings for a given mobility environment.

## VI. QL FOR MOBILITY LOAD BALANCING

Users are rarely uniformly distributed in cellular networks, but this is critical if a serving cell  $s$  is overloaded at a time when free resources exist in neighbor cells. A solution is required to automatically redistribute the load among cells - thus Mobility Load Balancing (MLB). MLB seeks to minimize the number of users who would otherwise not be satisfied, in terms of their data rate, owing to the overload in the serving cell  $s$ , i.e., to reduce the "Number of unsatisfied users" ( $N_{us}$ ).

To lower the serving cell's load  $\rho_s$ , MLB moves some of the edge users in  $s$  towards one or more neighbor cells or so called target cells. Let us denote the set of all target cells as



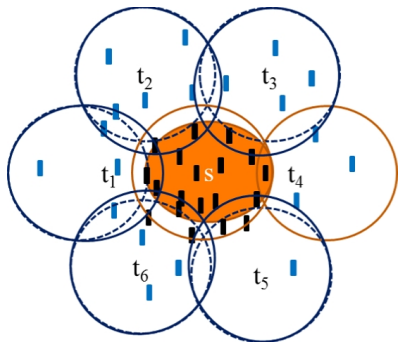


Fig. 10. Reactive change of CIOs

TABLE V  
CELL LOAD SCENARIOS ( $\Gamma$ )

$\rho_n \backslash \rho_s$	$< 0.45$	$[0.45, 0.60)$	$0.60+$
$[0.0-0.9)$	0	1	2
$[0.9-1.1)$	3	4	5
$1.1+$	6	7	8

$\rho_s$  - serving cell load  
 $\rho_n$  - Average T-cell Load

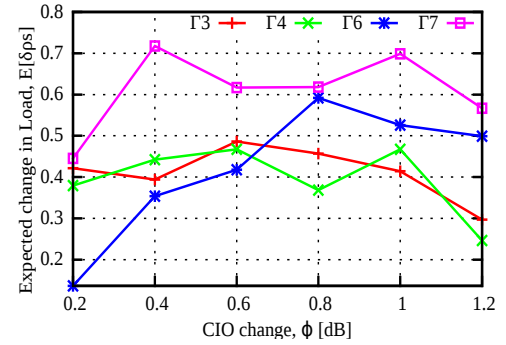


Fig. 11. Dependence of RLB gain on cells load

Fig. 11 shows the dependency of  $E[\Delta\rho_s]$  on  $\Gamma$  for different  $\phi$  values, evaluated in the network in Fig. 3a with 240 mobile users, and the ?centre cell? initially having 40 static users. We observe that the best  $\phi$ ,  $\phi_{opt}$ , is different for each  $\Gamma$ , e.g.,  $\phi_{opt}$  is 0.6, 1.0 and 0.4 for the three load scenarios  $\Gamma_3$ ,  $\Gamma_4$  and  $\Gamma_7$  respectively. This justifies the need for a learning solution that learns the required  $\phi_{opt}$  for each of the load scenario.

#### D. QLB: Learning the Optimum Actions

With RLB, a fixed  $\phi$  is found that guarantees good average performance, but not the best in each load scenario. For optimal performance, different  $\phi$  values would be required for the different scenarios. Moreover, for any change in CIO,  $\Delta\rho_s$  also depends on the user distribution ( $uD$ ) in cell  $s$ , i.e., more load can be offloaded from  $s$  if there are more cell edge users. We thus apply Q-Learning based Load Balancing (QLB) in order to learn the required CIO change for each combination of  $\rho_s$ ,  $\rho_n$  and  $uD$ . Here,  $uD$  describes the fraction of cell  $s$  users that are close to the cell edge, such that with more edge users, only a small CIO change is needed for these users to be handed over to neighbor cells.

QLB learns the action that instantaneously removes overload while ensuring that target cells are not overloaded as a result of its actions. Its Q-update algorithm is similar to Equation (8) with the corresponding components as follows:

1) *State-space*: Since  $\phi$  depends on  $\rho_s$ ,  $\rho_n$  and  $uD$ , each state is a vector  $[\rho_s, \rho_n, uD]$ . We defined 27 states as the 9 load scenarios in Table V for each of 3  $uD$  cases:  $uD < 20\%$ ;  $uD = [20 - 35]\%$  and  $uD \geq 35\%$ . Since the idea is to change the border to all neighbors,  $uD$  is evaluated generically within cell  $s$  and not specific to any one neighbor.

2) *Action-space*: Actions are the possible values that  $\phi$  can take, i.e., the values in dBs by which the CIO should be changed. Guided by RLB results, actions are selected as the discrete  $\phi$  values  $[0.2, 0.4, \dots, 1.0]$  dB.

3) *Rewards*: The rewards, as set in Equation 12, consider  $\Delta\rho_s$  (the achieved reduction in the serving cell load,  $\rho_s$ ) and the extra load created in neighbor cells.

$$r = \begin{cases} \Delta\rho_s + 1 & ; \Delta\rho_n = 2 \text{ and } \Gamma < 3 \\ \Delta\rho_s & ; \Delta\rho_n < 1 \\ \Delta\rho_s - 1 & ; \text{otherwise} \end{cases} \quad (12)$$

Positive  $\Delta\rho_s$  represents reduction in the offered  $\rho_s$  that results from users having moved to neighbor cells. Positive  $\Delta\rho_s$  is thus rewarded while the reverse is penalized. Since  $\rho_n$  is

#### Algorithm 3: QLB - The Q-Learning MLB Algorithm

Require: T-List(the list of  $s$  neighbors) and action set A

1. **if** LB action taken in previous SON interval  $i$  **do**
2. determine  $\Delta\rho_s$  **and** derive reward
3. update Q-table according to Equation 2
- end if**
4. evaluate  $\rho_s$
5. **if** overloaded **do**
6. determine load state  $l = [\rho_s, \rho_n, uD]$
7. **if** exploration complete **do**
8. select  $\phi = a_l^{opt}$ , the best value for state  $l$
9. **else do**
10. from A, select  $\phi = a_l^{i+1}$ , i.e., in sequence of last selected value  $a_l^i$  in state  $l$
- end if**
11. **for** each cell  $t$  in T-List with timer TOC expired
12. reduce  $O_t^s$  by  $\phi$  **and** increase  $O_s^t$  by  $\phi$
13. start timer  $T_{OC}$  for  $t$  LB HO towards  $s$
- end for**
- end if**
14. Restart SON interval timer , i.e.,  $t \leftarrow t + 1$

expected to increase as a result of adding users at the very edge of the cells, only  $\Delta\rho_n$  of more than 1 is penalized. In general, larger  $\Delta\rho_s$  values receive greater reward, but are accompanied by penalties for unrestrained actions taken in LB-states with high  $\rho_n$ . This allows high load  $t$  cells to overload just enough to propagate the load outwards but not too much to counterproductively cause further un-satisfaction after LB. However, special consideration is taken in cases where a large reduction in  $\rho_s$  can be achieved without overloading the target cell. In such cases, e.g., the change from scenario 9 to scenarios 1 or 2, the reward is increased by 1.

4) *QLB algorithm*: Each load state that is observed in one cell can reappear in any of the other cells in the network. As such, cells do not need to learn independent policies but learn a single shared policy in a cooperative learning process with a single Q-table which is updated by all cells.

During operation, each cell  $s$  observes its environment and at the end of a SON interval  $\tau = 5s$ ,  $s$  applies the procedure in Algorithm 3 to either take actions that reduce overload, learn from previous actions, or both. If overloaded,  $s$  selects an action according to Section VI-D2 and signals the new HO settings (with the revised CIOs) to all its associated UEs. Since CIO changes are symmetric for each  $s - t$  boundary,

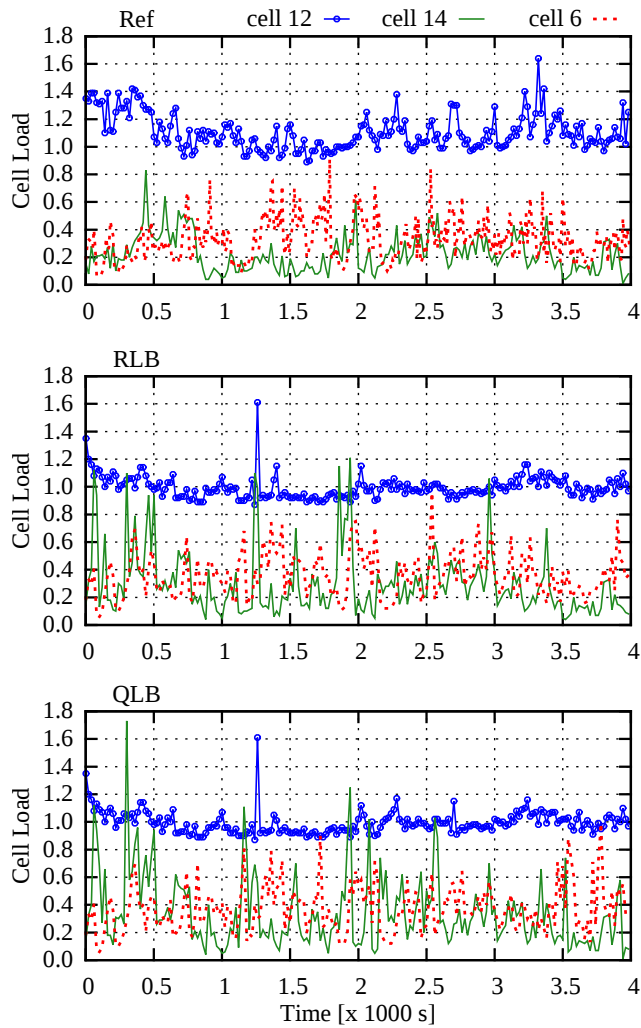


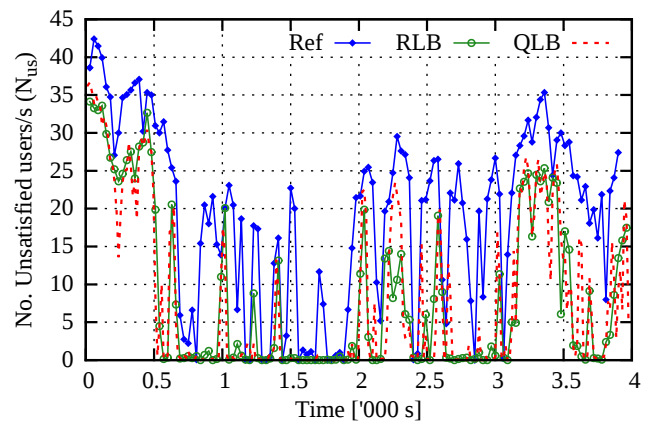
Fig. 12. QLB and RLB load redistribution (load moved from cell 12 to neighbor cells, e.g., 14 and 6)

$s$  also sends the new CIOs to its affected neighbor cells via the X2 interface. At the end of interval  $\tau + 1$ ,  $s$  evaluates the changes in load, derives the reward and updates the Q-table before repeating the process.

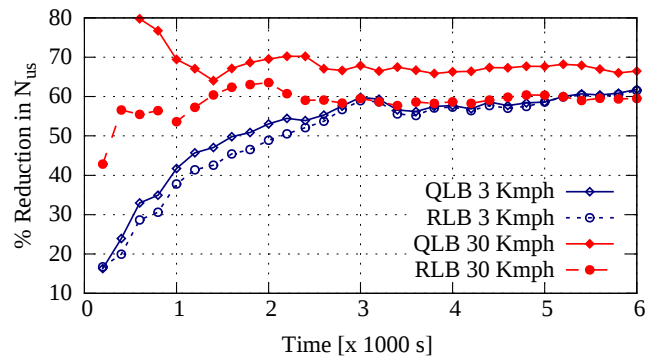
### E. QLB Simulation Results

QLB was evaluated using the LTE simulator of section IV. The results, shown in Figs. 12 and 13, consider two perspectives - load variation in individual cells to evaluate the dynamic performance of the solution, and the number of unsatisfied users ( $N_{us}$ ) which evaluates the global effect of the algorithm and its impact on user satisfaction. We compare the performance of QLB and RLB against *Ref*, the default scenario where the network is operated without any SON function.

1) *Learning towards load redistribution*: Fig. 12 shows the dynamic behavior of the load balancing solutions in terms of variations of cell load during the simulation, with the load evaluated every second for each cell. For clarity, the figure considers simulations at 3 kmph for which the dynamic behavior is slow enough to be analyzed. Otherwise, performance results at higher velocities are similar but only with much faster variations.



(a) Variation of  $N_{us}$  at 3 kmph



(b) QLB's reduction in  $N_{us}$

Fig. 13. Effect of QLB on user satisfaction

We observe that both solutions RLB and QLB lower the load in cell 12 by transferring it to neighbors cells (e.g., cell 14) and eventually to the outer cells (e.g., cell 6). After learning however, QLB responds better to overload compared to RLB. This can be seen for the period after 3000 s where cell 12 consistently has slightly lower load for QLB than for RLB. This is the direct consequence of QLB having learned the best CIO change for each load state, which is not the same with RLB. The drawback here is that in a low mobility network, the extra load added to the neighbor cells may take long to move outwards. In that case we need to evaluate performance in terms of the Number of unsatisfied users,  $N_{us}$ .

2) *QLB Effect on User satisfaction*: Fig. 13 compares the performance of RLB and QLB against the reference case (*Ref*) in terms of the  $N_{us}$  in the network for different velocity scenarios. Subfigure 13a shows the quasi-instantaneous variation of network-wide  $N_{us}$  in the 3 kmph network scenario for the same period considered in Fig. 12. With user satisfaction evaluated every 15s and each point capturing the total unsatisfaction events over the 15s interval, we see that both RLB and QLB reduce the  $N_{us}$  through the load redistribution.

Subfigure 13b evaluates the global benefit of the solutions for two velocity scenarios across the simulation. Considering the solutions over different batches of the simulation, we evaluate the gains of each of the two solutions measured in terms of the percentage reduction in  $N_{us}$  when compared to *Ref*. We observe that for both mobility scenarios, both RLB and QLB improve user Quality of Experience (QoE)

by reducing the  $N_{us}$ . The reduction in user dissatisfaction is comparable at low velocity since there is not much dynamism to be exploited by varying the CIO change. At higher velocity however, by adjusting CIOs for each instantaneous state that is observed, QLB achieves better user satisfaction.

The results above prove that a cognitive solution, in this a Q-learning based agent, can easily and successfully be applied towards a dynamic autonomous solution for MLB.

## VII. CONCLUSION

In this paper, we have proposed the Cognitive Cellular Networks (CCN) concept as a fully autonomous approach to Self-Organizing Networks (SON). Our contributions and results can be summarized as follows:

**Contributions:** We (1) proposed a Q-Learning (QL) framework as the method for implementing cognitive SON Functions (SFs); (2) justified how the QL framework can be used for any generic SF; and (3) discussed the application of this framework to two selected SFs: Mobility Robustness Optimization (MRO) and Mobility Load Balancing (MLB). We showed that, although all SFs use the same framework, special adjustments are required for each SF as dictated by its specific constraints, e.g., each SF required a different strategy on how to explore its action space. In general, however, the positive performance results for the developed solutions proved the benefit of cognitive approaches to SON and that QL provides a good framework for developing such functions.

In particular, Q-Learning based MRO (QMRO) is able to learn the best Handover Hysteresis (Hys) and Time To Trigger (TTT) settings for particular mobility states in the network. Applying cooperative learning, the cells learned a single policy function (single Q-table) that is thereafter exploited. Such an approach is applicable in any environment as demonstrated by the positive performance results obtained in the realistic network scenarios with User Equipments (UEs) having distinct and dynamically varying velocities.

Similarly, Q-Learning based Load Balancing (QLB) learned the best Cell Individual Offset (CIO) settings needed to reduce overload in different load conditions. Starting with Reactive Load Balancing (RLB) which adjusts the Handover (HO) boundary through the CIO, we observed that the required CIO change depends on the load state, characterized by the load in the serving and neighbor cells as well as the user distribution in the serving cell. QLB then learns the different CIO changes that are required for the different load states. Evaluating the solutions in multiple velocity environments showed that QLB achieves better results compared to the rule based reactive RLB. This is especially true in more dynamic environments, e.g., where users move at higher velocities. The subsequent effect of the load re-distribution is that the number of users, which would otherwise be unsatisfied as a result of low-data-rate induced overload, is reduced. **Convergence and complexity:** We have described in section III-E that with the same cost for each Q-update, the corresponding computational complexity is linear in the number of cells. In real systems, however, complexity may be a moving target which requires careful consideration of the specific case, i.e. the dynamics

of the environment. There will always be a trade-off between speed of adaption and complexity of computation. In the end, it is the deterioration of the system performance during the learning phase that stops us from employing too much learning rather than the complexity of any algorithms.

**Future works and extensions:** In both SFs, the learning strategy applied the simple approach of trying all the available actions a number of times and thereafter indefinitely applying the learned best action. Although, its convergence was improved by the combination of distributed exploration and centralized cooperative learning, some other variant of the approach would in practice be required since in a large network, cells operate in diverse environments that require specific handling. It may also be necessary to find methods through which states and actions can be automatically derived. This would reduce the design time and possibly the human subjectivity that is typically included in the solutions. On the contrary human intelligence could improve the performance, e.g., by limiting the applicable learning-time parameter space to a range known to have good performance.

In general, however, the results presented here confirm that the cognitive cellular network approach, and specifically QL, provides an outstanding approach to developing SON solutions especially where the dependence of metrics to the control parameters is not perfectly known.

## REFERENCES

- [1] 3GPP, "Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Self-configuring and Self-Optimizing Network (SON) use cases and solutions," 3GPP, Technical Report 36.902 version 9.3.1 Release 9, May 2011. [Online]. Available: <http://www.3gpp.org>
- [2] J. Alonso-Rubio, "Self-optimization for handover oscillation control in lte," in *IEEE/IFIP Network Operations and Management Symposium, Osaka, Japan.*, 2010.
- [3] I. Balan, T. Jansen, B. Sas, I. Moerman, and T. Kürner, "Enhanced weighted performance based handover optimization in LTE," in *Proceedings of FNMS, Warsaw, Poland.*, 2011.
- [4] G. Hui and P. Legg, "Soft metric assisted mobility robustness optimization in lte networks." Paris, France: International Symposium on Wireless Communication Systems (ISWCS 2012), 2012, pp. 1–5.
- [5] V. Capdevielle, A. F. A., and A. Fakhreddine, "Self-optimization of handover parameters in lte networks," in *Modeling and Optimization in Mobile, Ad Hoc and Wireless Networks (WiOpt), 2013 11th International Symposium on*, May 2013, pp. 133 – 139.
- [6] H. Eckhardt, S. Klein, and M. Gruber, "Vertical antenna tilt optimization for lte base stations," in *In Vehicular Technology Conference (VTC Spring), 2011 IEEE 73rd*, may 2011, pp. 1–5.
- [7] A. Gerdenitsch, S. Jakl, Y. Chong, , and M. Toeltsch, "A rule-based algorithm for common pilot channel and antenna tilt optimization in umts fdd networks," in *ETRI journal*, 2004, pp. 437–442.
- [8] M. Garcia-Lozano, S. Ruiz, , and J. Olmos, "Umts optimum cell load balancing for inhomogeneous traffic patterns," in *In Vehicular Technology Conference, 2004. VTC2004-Fall. 2004 IEEE 60th, Vol. 2*, sept 2004, pp. 909 – 913.
- [9] L. Du, J. Bigham, L. Cuthbert, C. Parini, , and P. Nahi, "Using dynamic sector antenna tilting control for load balancing in cellular mobile communications," in *In International Conference on Telecommunications, ICT2002, Beijing, Citeseer*, 2002, pp. 344– 348.
- [10] M. N. ul Islam and A. Mitschele-Thiel, "Reinforcement learning strategies for self-organized coverage and capacity optimization," in *Wireless Communications and Networking Conference (WCNC), 2012 IEEE*. IEEE, 2012, pp. 2818–2823.
- [11] M. Dirani and Z. Altman, "A cooperative reinforcement learning approach for inter-cell interference coordination in ofdma cellular networks," in *Modeling and Optimization in Mobile, Ad Hoc and Wireless Networks (WiOpt), 2010 Proceedings of the 8th International Symposium on*. IEEE, 2010, pp. 170–176.

