

# RAN Resource Slicing in 5G Using Multi-Agent Correlated Q-Learning

Hao Zhou, Medhat Elsayed and Melike Erol-Kantarci, *Senior Member, IEEE*

*School of Electrical Engineering and Computer Science*

*University of Ottawa*

Emails: {hzhou098, melsa034, melike.erolkantarci}@uottawa.ca

**Abstract**—5G is regarded as a revolutionary mobile network, which is expected to satisfy a vast number of novel services, ranging from remote health care to smart cities. However, heterogeneous Quality of Service (QoS) requirements of different services and limited spectrum make the radio resource allocation a challenging problem in 5G. In this paper, we propose a multi-agent reinforcement learning (MARL) method for radio resource slicing in 5G. We model each slice as an intelligent agent that competes for limited radio resources, and the correlated Q-learning is applied for inter-slice resource block (RB) allocation. The proposed correlated Q-learning based inter-slice RB allocation (COQRA) scheme is compared with Nash Q-learning (NQL), Latency-Reliability-Throughput Q-learning (LRTQ) methods, and the priority proportional fairness (PPF) algorithm. Our simulation results show that the proposed COQRA achieves 32.4% lower latency and 6.3% higher throughput when compared with LRTQ, and 5.8% lower latency and 5.9% higher throughput than NQL. Significantly higher throughput and lower packet drop rate (PDR) is observed in comparison to PPF.

**Index Terms**—5G RAN slicing, resource allocation, Q-learning, correlated equilibrium

## I. INTRODUCTION

The forthcoming 5G networks will provide support for vast amount of services and applications, where heterogeneous requirements for latency, bandwidth and reliability will coexist [1]. Three major traffic types are supported in 5G, namely enhanced Mobile Broad Band (eMBB), Ultra Reliable Low Latency Communications (URLLC), and massive Machine Type Communications (mMTC). The eMBB is regarded as an extension of LTE-Advanced services, which aims to provide high data rate for applications such as video streaming. The URLLC is proposed to provide a sub-millisecond latency and 99.999% reliability, which is critical for applications such as autonomous vehicles and remote surgery. The mMTC is designed to connect large number of Internet of Things devices, where data transmissions occur sporadically.

The stringent and heterogeneous QoS requirements of services have become a challenging problem in 5G, especially when different traffic types are multiplexed on the same channel. Considering the limited radio resources and increasing bandwidth demand, different methodologies are proposed for 5G radio resource allocation. A joint link adaptation and resource allocation policy is proposed in [2], which dynamically adjusts the block error probability of URLLC small payload transmissions based on cell load. A risk sensitive method is

used in [3] to allocate resources for the incoming URLLC traffic, while minimizing the risk of the eMBB transmissions and ensuring URLLC reliability. Puncturing technique is applied in [4] to guarantee minimum latency of URLLC, where eMBB traffic is scheduled at the beginning of slots, while URLLC traffic can puncture at any time with a higher priority.

A common feature of aforementioned works is that URLLC traffic is scheduled on top of eMBB traffic such as puncturing technique [3], [4], and a potential priority is applied to guarantee the latency and reliability of URLLC traffic [2]–[5]. As a result, the eMBB traffic will be affected with degraded throughput [2], [4], [5]. Meanwhile, another important problem is the increasing complexity of wireless networks, e.g., the evolving network architecture, dynamic traffic patterns and increasing devices numbers, which makes it harder to build a dedicated optimization model for resource allocation.

To this end, the emerging reinforcement learning (RL) techniques become a promising solution [6]. In [7], a Latency-Reliability-Throughput Q-learning algorithm is proposed for jointly optimizing the performance of both URLLC and eMBB users. [8] develops an RL method to select different scheduling rules according to the scheduler states, which aims to minimize the traffic delay and Packet Drop Rate (PDR). The random forest algorithm is applied in [9] to accomplish the Transmission Time Interval (TTI) selection for each service, and the result shows a lower delay and lower PDR for URLLC traffic while guaranteeing the eMBB throughput requirements. Furthermore, [10], [11] use deep reinforcement learning (DRL) scheme for resource allocation in 5G, in which neural networks are used to learn allocation rules.

In this paper, we propose a multi-agent reinforcement learning (MARL) based resource allocation algorithm, where the performance of URLLC and eMBB are jointly optimized. Different than aforementioned works, we apply the network slicing scheme to aggregate users with similar QoS requirements. Network slicing is an important feature in 5G [12]. Based on software defined network (SDN) and network function virtualization (NFV) techniques, physical network infrastructures are divided into multiple independent logical network slices. Each slice is presumed to support services with specific QoS requirements, and the whole network achieves a much higher flexibility and scalability. In the proposed correlated Q-learning based inter-slice RB allocation (COQRA) scheme, firstly, each slice is assumed to be an intelligent agent

to compete for limited RBs, and the model-free correlated Q-learning algorithm is applied for inter-slice resource allocation. Then resources (more specifically, RBs of the 5G New Radio (NR)) are distributed by each slice among its attached users by proportional fair algorithm, which is the intra-slice allocation [2]. Compared with Nash Q-learning (NQL) and Latency-Reliability-Throughput Q-learning (LRTQ) techniques [7], the results present a 5.8% and 32.4% lower latency for URLLC traffic, and 5.9% and 6.3% higher throughput for eMBB traffic. COQRA also achieves significantly higher throughput and a lower PDR than priority proportional fairness (PPF) algorithm.

The main contribution of this work is that we develop a MARL-based RAN resource slicing scheme for 5G NR. In the proposed multi-agent COQRA, each agent makes decisions autonomously, where they coordinate by exchanging Q-values among each other. Compared with other multi-agent coordination methods, such as Nash equilibrium, the correlated equilibrium is readily solved using linear optimization [13], which is critical for the fast response requirement of wireless network.

The rest of this paper is organized as follows. Section II presents related work. Section III defines the system model and problem formulation. Section IV introduces the proposed COQRA scheme and the baseline algorithms. Simulation results are presented in section V, and section VI concludes the paper.

## II. RELATED WORK

In the literature, various slicing based resource allocation methods have been investigated using both model-free and model-based algorithms. For instance, a QoS framework is proposed in [14] for network slicing, in which three types of slices are defined. [15] presents an RL method for resource optimization in 5G network slicing, and the results show an improvement in network utility and scalability. A QoS-aware slicing algorithm is proposed in [16] where the bandwidth is distributed based on utility function and priority rules. In [17], a network slicing and multi-tenancy based dynamic resource allocation scheme is presented, in which hierarchical decomposition is adopted to reduce complexity in optimization. Considering the multi-slice and multi-service scenarios, deep Q-learning is deployed in [18] for end to end network slicing resource allocation.

Allocation of limited bandwidth resources among slices has been the main challenge in RAN slicing [19]. The allocation should meet various QoS requirements under the constraint of limited bandwidth budget. Different from existing works, in this paper we solve the resource allocation problem by a MARL approach. We propose a multi-agent COQRA method to distribute RBs among slices. Correlated Q-learning has been applied in microgrid energy management in [20] to maximize the profit of agents. However, 5G network has much more stringent requirements for agents such as latency and PDR, which is different than the microgrid system.

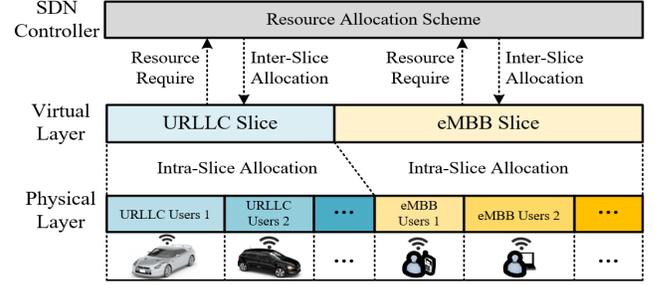


Fig. 1. Network slicing based two-step resource allocation.

## III. SYSTEM MODEL AND PROBLEM FORMULATION

As shown in Fig.1, we consider URLLC and eMBB slices where each slice serves several users. First, the slice manager collects QoS requirements of the users such as bandwidth and latency, then the collected information is sent to the SDN controller for the required resources. Based on the received requirements, SDN controller implements the inter-slice RB allocation to distribute RBs between slices. Then, the users are scheduled within the allocated RBs for that particular slice. We consider numerology 0 where one RB contains 12 sub-carriers in frequency domain.

Here we assume each slice manager is an intelligent agent making decisions autonomously. For the eMBB agent, it needs to maximize the throughput, as denoted by:

$$\max \sum_{j=1}^J \sum_{e=1}^{E_j} b_{j,e,t}, \quad (1)$$

where  $b_{j,e,t}$  is the throughput of  $e^{th}$  eMBB user in  $j^{th}$  Base station (BS) at time  $t$ , and  $E_j$  is the number of eMBB users of  $j^{th}$  BS.

URLLC agent needs to minimize the delay as follows:

$$\min \sum_{j=1}^J \sum_{u=1}^{U_j} d_{j,u,t}, \quad (2)$$

where  $d_{j,u,t}$  is the delay of  $u^{th}$  URLLC user in  $j^{th}$  BS at time  $t$ , and  $U_j$  is the number of URLLC users of  $j^{th}$  BS.

The packet delay  $d$  mainly consists of three components:

$$d = d^{tx} + d^{que} + d^{rtx}, \quad (3)$$

where  $d^{tx}$  is the transmission delay,  $d^{que}$  is the queuing delay, and  $d^{rtx}$  is the HARQ re-transmission delay.

The transmission delay of a packet depends on the link capacity between the UE and the BS:

$$d^{tx} = \frac{L_u}{C_{u,j}}, \quad (4)$$

where  $L_u$  is the packet size of  $u^{th}$  UE, and  $C_{u,j}$  is the link capacity between  $u^{th}$  UE and the BS it belongs to.

The link capacity is calculated as follows:

$$C_{u,j} = \sum_{r \in N_{u,RB}^{RB}} c_r \log \left( 1 + \frac{p_{j,r} x_{j,r,u} g_{j,r,u}}{c_r N_0 + \sum_{j' \in J-j} p_{j',r} x_{j',r',u'} g_{j',r',u'}} \right), \quad (5)$$

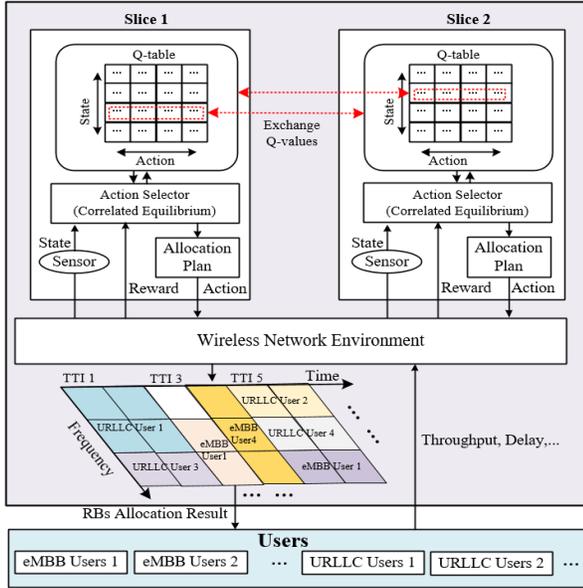


Fig. 2. Proposed COQRA architecture for intelligent resource management among slices.

where  $N_u^{RB}$  is the set of RBs that the  $u^{th}$  UE uses,  $c_r$  is the bandwidth of  $r^{th}$  RB,  $p_{j,r}$  is the transmission power of  $r^{th}$  RB in  $j^{th}$  BS,  $x_{j,r,u}$  is a binary variable to indicate whether this RB is distributed to  $u^{th}$  UE,  $g_{j,r,u}$  is the channel gain between BS and UE,  $N_0$  is the unit noise power density, and  $j' \in J_{-j}$  is the BS set except  $j^{th}$  BS.

#### IV. CORRELATED Q-LEARNING BASED RESOURCE ALLOCATION (COQRA)

##### A. COQRA architecture

The architecture of the proposed multi-agent COQRA is illustrated in Fig.2. Each slice is an independent agent, and it observes its own state from the environment. The agent exchanges Q-values with its peers to make decisions, and the action selection is determined by correlated equilibrium. Then, the selected actions are sent to wireless environment, and eMBB and URLLC users are scheduled RBs within the slice resources according to the proportional fairness algorithm [2]. Users transmit packets based on allocated bandwidth, and the experienced throughput and delay are sent back to agents. Finally, the reward is received, and slice managers make next decisions based on new state and updated Q-values.

##### B. Markov decision process and Q-learning

In this section, based on the system model in Section III, we will define the Markov decision process (MDP) to describe agents and introduce the learning scheme. Here we define each slice manager as a intelligent agent, which will interact with the environment and make decisions autonomously.

We assume each agent has its own state, action and reward signal. The state  $s^u$  for URLLC slice manager agent (from hereon, referred to as URLLC agent) is the number of packets in its queue, and the action  $a^u$  is the number of RBs it

allocates. The state and action for the eMBB slice manager agent (from hereon, referred to as eMBB agent) is defined similarly, namely  $s^e$  and  $a^e$ . Thus, the Q-space for both agents are  $Q^u = \{s^u, a^u\}$  and  $Q^e = \{s^e, a^e\}$ .

The reward function for eMBB agent is given in (6), where obtaining higher throughput leads to higher reward.

$$r_{j,t}^{eMBB} = \frac{2}{\pi} \arctan\left(\sum_{j=1}^J \sum_{e=1}^{E_j} b_{j,e,t}\right). \quad (6)$$

The reward function for URLLC agent at time  $t$  is:

$$r_{j,t}^{URLLC} = \begin{cases} 1 - \max_{u \in H_t^u} (d_u^{que})^2, & |H_t^u| \neq 0, \\ 0, & |H_t^u| = 0, \end{cases} \quad (7)$$

where  $d_u^{que}$  is the queuing delay for  $u^{th}$  URLLC user,  $|H_t^u|$  denotes the length of the queue for URLLC users at time slot  $t$ . In (7), a lower queuing delay means a higher reward, which depends on the number of RBs that the agent gets. URLLC agent competes for more RBs to reduce the queuing delay. Meanwhile, to guarantee the PDR performance, we apply a penalty if any packet is dropped.

In Q-learning, one agent always aims to maximize the long-term accumulated reward. For one agent  $i$ , the state value is:

$$V_i^\pi(s_i) = \mathbb{E}_\pi\left(\sum_{n=0}^{\infty} \theta^n r_i(s_{i,n}, a_{i,n}) | s_i = s_{i,0}\right), \quad (8)$$

where  $\pi$  is the policy,  $s_{i,0}$  is the initial state,  $r_i(s_{i,n}, a_{i,n})$  is the reward of taking action  $a_{i,n}$  at state  $s_{i,n}$ ,  $\theta$  is the reward discount factor.  $V_i^\pi(s_i)$  represents long-term expected accumulated reward at state  $s_i$ .

Then we define the state-action value  $Q_i^\pi(s_i, a_i)$  to describe the expected reward of taking action  $a_i$  under state  $s_i$ :

$$Q_i^\pi(s_i, a_i) = (1 - \alpha)Q_i^\pi(s_i, a_i) + \alpha(r(s_i, a_i) + \gamma \max_{a'_i} Q_i^\pi(s'_i, a'_i)) \quad (9)$$

where  $s_i$  and  $s'_i$  are current and next state, and  $a_i$  and  $a'_i$  are current and next action, and  $\alpha$  is the learning rate. By updating the Q-values, the agent will learn the best action sequence, and achieve a long term best reward.

When there is only one agent, the  $\epsilon$ -greedy policy is generally applied to balance the exploration and exploitation.

$$a_i = \begin{cases} \text{random action}, & \text{rand} \leq \epsilon, \\ \arg \max(Q_i^\pi(s_i, a_i)), & \text{otherwise,} \end{cases} \quad (10)$$

where  $\epsilon$  is the exploration probability and  $0 < \epsilon < 1$ , and  $\text{rand}$  indicates a random number between 0 and 1.

On the other hand, in a multi-agent environment, the action of one agent will affect both the environment and the reward of the other agents. Therefore, we propose a correlated Q-learning based resource allocation approach to address the multi-agent 5G environment.

### C. Correlated equilibrium

Given the limited bandwidth resources, in our multi-agent environment, each slice manager agent will compete for more RBs to optimize their own goal, which may lead to a conflict. We use the correlated equilibrium to balance the reward of each agent, and maintain a good overall performance for the whole multi-agent system. In correlated equilibrium, agents exchange Q-values to share information with each other, and the joint action is chosen according to the following equations:

$$\begin{aligned} & \max \sum_{\vec{a} \in A} Pr(\vec{s}, \vec{a}) Q(\vec{s}, \vec{a}) \\ & \text{sub.to} \sum_{\vec{a} \in A} Pr(\vec{s}, \vec{a}) = 1 \\ & \sum_{\vec{a}_{-i} \in A_{-i}} Pr(\vec{s}, \vec{a}_i) (Q(s, \vec{a}_i) - Q(\vec{s}, \vec{a}_{-i}, a_i)) \geq 0 \\ & 0 \leq Pr(\vec{s}, \vec{a}) \leq 1 \end{aligned} \quad (11)$$

where  $\vec{s}$  is the system state of eMBB and URLLC agents  $\vec{s} = (s^e, s^u)$ ,  $\vec{a} = (a^e, a^u)$  is the joint action,  $Pr(\vec{s}, \vec{a})$  is the probability of choosing action combination  $\vec{a}$  under state  $\vec{s}$ ,  $a_i$  is the action of agent  $i$ ,  $\vec{a}_{-i}$  is the action combination except agent  $i$ , and  $A_{-i}$  is the set of  $\vec{a}_{-i}$ . The correlated equilibrium is described as a linear program in (11). The objective is to maximize the total expected reward of all agents, and the constraints guarantee a probability distribution of action combination in which each agent chooses an optimal action.

Based on correlated equilibrium, an improved  $\epsilon$ -greedy policy is applied for action selection:

$$\pi_i(s) = \begin{cases} \text{random action,} & \text{rand} \leq \epsilon, \\ \text{equation(11),} & \text{rand} > \epsilon. \end{cases} \quad (12)$$

Exploration is performed whenever  $\text{rand} \leq \epsilon$ , i.e., random action is selected. Otherwise the exploitation is implemented by correlated equilibrium. The COQRA scheme is summarized in Algorithm 1. In COQRA, the two-step resource allocation method avoids the complexity of processing all UE requirement by a central controller, which will reduce the computational complexity of learning algorithm by a smaller action space.

### D. Nash equilibrium

In this section, we introduce the NQL algorithm, which is generally used in multi-agent problems. Compared with correlated equilibrium, the Nash equilibrium is a iterative based coordination method. There could be more than one Nash equilibrium or no equilibrium in some cases. We use NQL as a baseline algorithm. In NQL, agents select actions by:

$$U_i(\vec{a}_{-i}, a_i) \geq U_i(\vec{a}), \vec{a} \in A, \vec{a}_{-i} \in A_{-i}, \quad (13)$$

where  $\vec{a}$  is the action combination,  $\vec{a}_{-i}$  is the action combination except agent  $i$ ,  $A$  is the set of  $\vec{a}$ , and  $A_{-i}$  is the set of  $\vec{a}_{-i}$ .  $U_i$  is the utility function for agent  $i$ , which refers to Q-values in this paper. At Nash equilibrium, agents are less likely to change their actions as this will lead to lower observed

utility. Similarly, we apply an improved  $\epsilon$ -greedy policy to select actions as (12). The NQL scheme is summarized in Algorithm 2. We assume the equilibrium is randomly selected if more than one equilibrium are found.

### E. LRTQ and PPF algorithms

To further investigate the performance of the proposed method, two more baseline algorithms are used in this paper. LRTQ was proposed in [7]. LRTQ is also a learning-based resource allocation method, but it only defines one reward function for all users. PPF algorithm is applied in [2]. In PPF, the RBs are first allocated to URLLC users to guarantee low latency, then the remaining RBs are distributed among eMBB users. Note that network slicing is not implemented in both

---

#### Algorithm 1 COQRA based Resource Allocation

---

- 1: **Initialize:** Q-learning and wireless network parameters
  - 2: **for**  $TTI = 1$  to  $T$  **do**
  - 3:   **if**  $\text{rand} < \epsilon$  **then**
  - 4:     Choose  $a_t^u$  and  $a_t^e$  randomly.
  - 5:   **else**
  - 6:     Agents exchange Q-values under their current state. Find correlated equilibrium using eq. (11) and choose action  $a_t^u, a_t^e$
  - 7:   **end if**
  - 8:   Complete the inter-slice resource allocation.
  - 9:   Implement intra-slice allocation. Schedule users on respective slice resources using the proportional fair algorithm.
  - 10:   Agents calculate reward based on received QoS metrics.
  - 11:   Update agent state  $s^u, s^e$  and Q-table  $Q^u$  and  $Q^e$ .
  - 12: **end for**
- 

---

#### Algorithm 2 NQL based Resource Allocation

---

- 1: **Initialize:** Q-learning and wireless network parameters
  - 2: **for**  $TTI = 1$  to  $T$  **do**
  - 3:   **if**  $\text{rand} < \epsilon$  **then**
  - 4:     Choose  $a_t^u$  and  $a_t^e$  randomly.
  - 5:   **else**
  - 6:     **for** Each agent **do**
  - 7:       Search its own Nash equilibrium using eq. (13).
  - 8:     **end for**
  - 9:     Match each agent's equilibrium to get the system Nash equilibrium.
  - 10:     Agents choose actions according to system Nash equilibrium.
  - 11:   **end if**
  - 12:   Complete the inter-slice resource allocation.
  - 13:   Implement intra-slice allocation, using the proportional fair algorithm.
  - 14:   Agents calculate reward based on received QoS metrics.
  - 15:   Update agent state  $s^u, s^e$  and Q-table  $Q^u$  and  $Q^e$ .
  - 16: **end for**
-

TABLE I  
PARAMETERS SETTINGS

Parameters	Value
Traffic Model	URLLC Traffic: 80% Poisson distribution and 20% constant bit rate.
	eMBB Traffic: Poisson distribution.
	Packet Size: 32 Bytes.
Propagation Model	$128.1+37.6\log(\text{distance}(\text{km}))$ .
	Log-Normal shadowing with 8 dB.
Transmission settings	Transmission power: 40 dBm (Uniform distributed)
	Tx/Rx antenna gain: 15 dB.
	3GPP urban network.
Q-learning	Learning rate: 0.9
	Discount factor: 0.5
	Epsilon value: 0.05

of these baseline algorithms. These algorithms perform RB allocation decisions only.

## V. PERFORMANCE EVALUATION

### A. Parameter settings

We use MATLAB to implement our proposed algorithm. We consider five gNBs with 500 meter inter-site distance, each serving one eMBB and one URLLC slice. Each eMBB slice is serving 5 users, and URLLC slice has 10 users, which is randomly distributed in the cell. The bandwidth for a cell is 20 MHz, and there are 100 RBs. Each RB contains 12 sub-carriers, and each sub-carrier has 15kHz. 100 RBs are divided

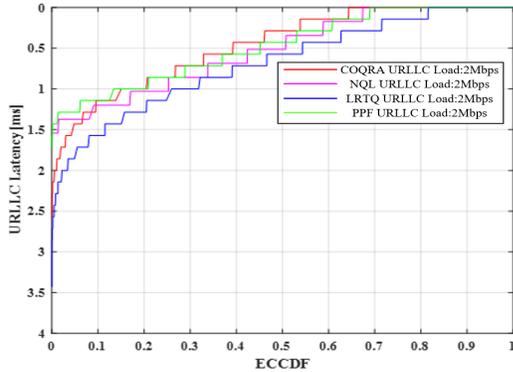


Fig. 3. URLLC latency distribution[ms] under load=2Mbps per cell.

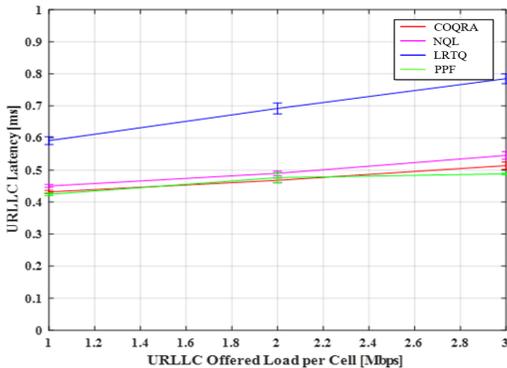


Fig. 4. URLLC latency [ms] with varying offered load [Mbps] per cell.

into 13 RB groups, where the first 12 groups contain 8 RBs each, and the last group contains 4 RBs. The simulation period is 5000 TTIs, and each TTI contains 2 OFDM symbols (5G mini-slot length 0.143ms) allocations. The  $\epsilon$ -greedy policy is implemented in first 3000 TTI, and the rest 2000 TTI is pure exploitation. Other parameters are shown in Table I. Each scenario is repeated 10 runs to get an average value with 95% confidence interval.

### B. Simulation Results

First, we set eMBB load to 0.5 Mbps per cell, and consider that URLLC load changes from 1 Mbps to 3 Mbps per cell. The latency distribution of four algorithms are shown in Fig.3. Meanwhile, Fig.4 presents the average URLLC latency against varying URLLC offered loads. The results show that COQRA, NQL and PPF have a comparable latency distribution, while the LRTQ has a relatively higher latency. The PPF has the lowest overall latency for URLLC traffic, and the main reason is that URLLC traffic has a priority in this method. Whenever URLLC packet arrives, it will be directly scheduled over eMBB traffic in the PPF algorithm. Meanwhile, the COQRA achieves 4.4% and 5.8% lower latency than NQL when URLLC load is 2 Mbps and 3 Mbps, respectively. Compared with LRTQ, the COQRA has a 27.1% lower latency under 2 Mbps load, and 32.4% lower latency under 3 Mbps load.

Furthermore, the eMBB throughput under different URLLC load is shown in Fig.5. The result shows that COQRA, NQL and LRTQ have a close performance of throughput, while the PPF has a much lower value. When URLLC load is 3

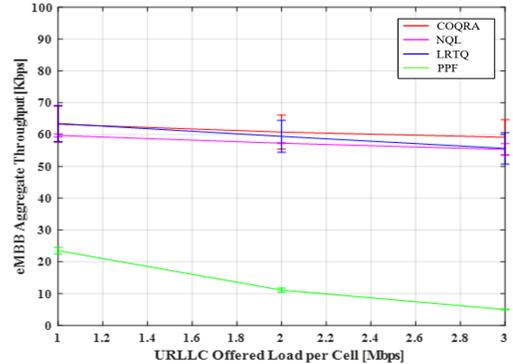


Fig. 5. eMBB throughput with varying URLLC loads [Mbps].

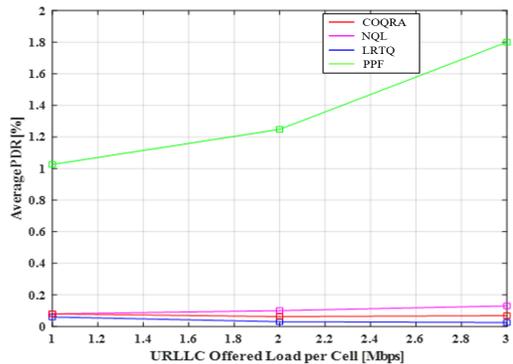


Fig. 6. Average PDR under varying URLLC loads [Mbps].

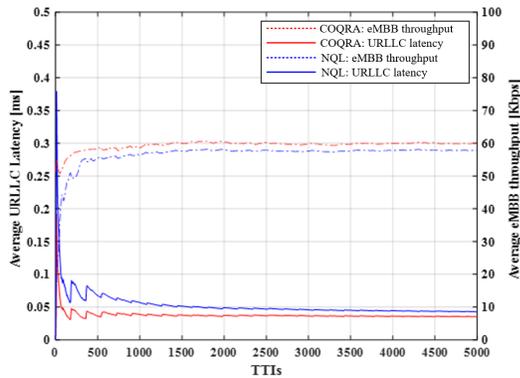


Fig. 7. Convergence performance of COQRA and NQL.

Mbps, the COQRA method has a 5.9% higher throughput over NQL method, and 6.3% higher than LRTQ. Although the PPF has a good latency performance for URLLC, the eMBB throughput is almost 90% lower than other three algorithms. This result can still be explained by the priority settings in PPF, which means the eMBB throughput will decrease with increasing prioritized URLLC load. On the other hand, the COQRA, NQL and LRTQ benefit from the jointly optimizing scheme, and they maintain a good throughput performance. In Fig.6, we compare the PDR of four algorithms. We show that COQRA, NQL and LRTQ maintain a much lower PDR than PPF method. In learning algorithms, agents will get a negative reward when dropping packets. However, PPF fails to maintain low PDR under all traffic loads, where a worst case PDR of 1.8% is observed. Finally, we compare the convergence performance of COQRA and NQL, and a faster convergence is observed for COQRA in Fig.7. The reason is that COQRA has a more efficient way to find the equilibrium. Overall, COQRA outperforms all baseline methods in terms of latency, throughput, packet loss and convergence time.

## VI. CONCLUSION

5G and beyond 5G networks will serve heterogeneous users of multiple slices which calls for new ways of network slicing and resource allocation. Machine learning techniques provide a promising alternative to the existing schemes. In this paper, we propose a Radio Access Network (RAN) slicing based resource allocation method for 5G, namely correlated Q-learning based inter-slice RB allocation (COQRA), to allocate radio resources to eMBB and URLLC slices. The proposed algorithm is compared with Nash Q-learning method, Latency-Reliability-Throughput Q-learning method and priority proportional fairness algorithm. Simulation results show that the proposed COQRA scheme achieves the best overall performance. In the future works, we plan to enhance the scalability of COQRA such that it can be used for intra-slice allocations.

## ACKNOWLEDGMENT

We thank Dr. Shahram Mollahasani for his generous help and useful discussions. This work is funded by the NSERC CREATE and Canada Research Chairs programs.

## REFERENCES

- [1] M. Shafi, A. Molisch, P. Smith, T. Haustein, P. Zhu, P. Silva, F. Tufvesson, A. Benjebbour, and G. Wunder, "5G: A Tutorial Overview of Standards, Trials, Challenges, Deployment, and Practice," in *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 6, pp. 1201-1221, Jun. 2017.
- [2] G. Poci, K. Pedersen, P. Mogensen, "Joint Link Adaptation and Scheduling for 5G Ultra-Reliable Low-Latency Communications," in *IEEE Access*, vol. 6, DOI: 10.1109/ACCESS.2018.2838585.
- [3] M. Alsenwi, N. Tran, M. Bennis, A. Bairagi, and C. Hong, "eMBB-URLLC Resource Slicing: A Risk-Sensitive Approach," in *IEEE Communications Letters*, vol. 23, no. 4, pp. 740-743, Apr. 2019.
- [4] S. R. Pandey, M. Alsenwi, Y. K. Tun, and C. S. Hong, "A Downlink Resource Scheduling Strategy for URLLC Traffic," in *Proc. of IEEE Int. Conf. on Big Data and Smart Computing*, pp. 1-6, Feb. 2019.
- [5] A. A. Esswie and K. I. Pedersen, "Multi-User Preemptive Scheduling For Critical Low Latency Communications in 5G Networks," in *IEEE Symposium on Computers and Communications (ISCC)*, pp. 136-141, Jun. 2018.
- [6] M. Elsayed and M. Erol-Kantarci, "AI-Enabled Future Wireless Networks: Challenges, Opportunities, and Open Issues," in *IEEE Vehicular Technology Magazine*, vol. 14, no.3, pp. 70-77, Sep.2019.
- [7] M. Elsayed and M. Erol-Kantarci, "AI-Enabled Radio Resource Allocation in 5G for URLLC and eMBB Users," in *Proc. of 2019 IEEE 2nd 5G World Forum (5GWF)*, pp. 590-595, Sep. 2019.
- [8] I. Comsa, S. Zhang, M. Aydin, P. Kuonen, Y. Lu, R. Trestian, and G. Ghinea, "Towards 5G: A Reinforcement Learning-Based Scheduling Solution for Data Traffic Management," in *IEEE Trans. on Network and Service Management*, vol. 15, no. 4, pp. 1661-1675, Dec. 2018.
- [9] J. Zhang, X. Xu, K. Zhang, B. Zhang, X. Tao, and P. Zhang, "Machine Learning Based Flexible Transmission Time Interval Scheduling for eMBB and uRLLC Coexistence Scenario," in *IEEE Access*, vol. 7, DOI: 10.1109/ACCESS.2019.2917751.
- [10] Y. Huang, S. Li, C. Li, Y. Hou, and W. Lou, "A Deep-Reinforcement-Learning-Based Approach to Dynamic eMBB/URLLC Multiplexing in 5G NR," in *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6439-6456, Jul. 2020.
- [11] H. Chergui and C. Verikoukis, "Offline SLA-Constrained Deep Learning for 5G Networks Reliable and Dynamic End-to-End Slicing", in *IEEE Journal on Selected Areas in Communications*, Vol. 38, Num. 2, 2020.
- [12] F. Rezazadeh, H. Chergui et al., "Continuous Multi-objective Zero-touch Network Slicing via Twin Delayed DDPG and OpenAI Gym", in *Proc. of IEEE Global Communications Conference*, pp. 1-6, Dec. 2020.
- [13] H. Zhou, and M. Erol-Kantarci, "Correlated Deep Q-learning based Microgrid Energy Management," in *Proc. of 2020 IEEE 25th International Workshop on CAMAD*, pp. 1-6, Sep. 2020.
- [14] Z. Shu, and T. Taleb, "A Novel QoS Framework for Network Slicing in 5G and Beyond Networks Based on SDN and NFV," in *IEEE Network*, vol. 34, No. 3, pp. 256-263, May 2020.
- [15] Y. Shi, Y. Sagduyu, and T. Erpek, "Reinforcement Learning for Dynamic Resource Optimization in 5G Radio Access Network Slicing," in *Proc. of 2020 IEEE International Workshop on CAMAD*, pp. 1-6, Sep. 2020.
- [16] R. Schmidt, C. Chang, and N. Nikaen, "Slice Scheduling with QoS-Guarantee Towards 5G," in *Proc. of 2019 IEEE Global Communications Conference*, pp. 1-7, Dec. 2019.
- [17] S. Oladejo, and O. Falowo, "Latency-Aware Dynamic Resource Allocation Scheme for Multi-Tier 5G Network: A Network Slicing-Multitenancy Scenario," in *IEEE Access*, DOI: 10.1109/ACCESS.2020.2988710, Apr. 2020.
- [18] T. Li, X. Zhu, and X. Liu, "An End-to-End Network Slicing Algorithm Based on Deep Q-Learning for 5G Network," in *IEEE Access*, DOI: 10.1109/ACCESS.2020.3006502, Jul. 2020.
- [19] M. Maule, P. Mekikis, K. Ramantas, J. Vardakas, and C. Verikoukis, "Dynamic partitioning of radio resources based on 5G RAN Slicing," in *Proc. of 2020 IEEE Global Communications Conference*, pp. 1-7, Dec. 2020.
- [20] H. Zhou and M. Erol-Kantarci, "Decentralized Microgrid Energy Management: A Multi-agent Correlated Q-learning Approach," in *Proc. of IEEE Int. Conf. on Communications, Control, and Computing Technologies for Smart Grids*, pp.1-7, Nov. 2020.