

# ADAPTIVE CONTENTION WINDOW DESIGN USING DEEP Q-LEARNING

Abhishek Kumar\*, Gunjan Verma†, Chirag Rao†, Ananthram Swami†, and Santiago Segarra\*

\*Rice University, USA

†US Army’s CCDC Army Research Laboratory, USA

## ABSTRACT

We study the problem of adaptive contention window (CW) design for random-access wireless networks. More precisely, our goal is to design an intelligent node that can dynamically adapt its minimum CW (MCW) parameter to maximize a network-level utility knowing neither the MCWs of other nodes nor how these change over time. To achieve this goal, we adopt a reinforcement learning (RL) framework where we circumvent the lack of system knowledge with local channel observations and we reward actions that lead to high utilities. To efficiently learn these preferred actions, we follow a deep Q-learning approach, where the Q-value function is parametrized using a multi-layer perceptron. In particular, we implement a rainbow agent, which incorporates several empirical improvements over the basic deep Q-network. Numerical experiments based on the NS3 simulator reveal that the proposed RL agent performs close to optimal and markedly improves upon existing learning and non-learning based alternatives.

**Index Terms**— Wireless network, random access, contention window, reinforcement learning, deep Q-learning.

## 1. INTRODUCTION

Wireless networks have become an essential part of our lives. We use them for a wide range of purposes including video streaming, web services, and file sharing. Moreover, the number of wireless devices is constantly increasing, thus precipitating the urgent need for smart and efficient access to the limited available spectrum. Though oftentimes preferred due to their simplicity, conventional random access protocols such as IEEE 802.11 lead to suboptimal spectrum utilization and, more importantly, the performance of these protocols drops significantly under dynamic and uncertain settings [1–3].

In this context, there has been growing interest over the last few years in improving the management of wireless networks under uncertain scenarios. A testament to this trend is the recently completed DARPA spectrum collaboration challenge [4]. In this competition, multiple wireless networks are present in the same geographical area and, without full information about the protocols employed by each network, they must interact with each other to efficiently utilize the shared spectrum. Furthermore, recent works on heterogeneous networks [5–7] analyze the efficient spectrum management in the case where nodes in the *same* network may follow different access protocols.

Research was sponsored by the Army Research Office and was accomplished under Cooperative Agreement Number W911NF-19-2-0269. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Office or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein. Emails: {ak109, segarra}@rice.edu, {gunjan.verma.civ, chirag.r.rao.civ, ananthram.swami.civ}@mail.mil.

Following this trend, we consider the scenario where nodes *do* know the protocol followed by others but the uncertainty stems from not knowing key parameters governing this protocol. In particular, we consider a binary exponential backoff protocol with the minimum contention window (MCW) being the key unknown parameter to be designed. The proposed uncertain scenario is of empirical relevance in cases where nodes may act unfairly by reducing their MCW to acquire more channel access at the expense of others. In this setting, we would rely on intelligent agents to adapt their own MCW to restore fairness [3]. A similar situation might arise when a new node joins a network and, without full knowledge of others’ parameters, the new node seeks to learn their dynamics and predict the optimal MCW value at every time. We address this setting of *parameter uncertainty* by modeling the problem through a reinforcement learning (RL) framework and adapting state-of-the-art deep learning techniques for its solution.

**Related work.** The problem of optimal CW design has been studied extensively over the past decades [8–14]. For example, authors in [8] and [11] set the optimal MCW as a linear function of the estimated number of active nodes. Furthermore, control-based [9] and game-theoretic [10] approaches have been developed to avoid the need for estimating the number of nodes. However, unlike our setting, all these approaches assume that the nodes behave in a cooperative manner by either choosing the same MCW or by not deviating from a prespecified behavior. Another class of approaches relies on modifications in the increase and decrease of the CW instead of finding an optimal MCW [2, 15, 16]. For instance, under [15], a node multiplies its current CW by a constant factor if packets collide and decreases it by subtracting a constant value if the transmission is successful. Our approach differs from this body of work since it does not deviate from the current protocol structure and is learning-based as opposed to being controlled by preset rules.

Motivated by the recent success of machine learning applied to wireless networks [17–22], learning-based solutions have been proposed for optimal distributed CW design [3, 23–26]. Germane to our approach is [24], where (non-deep) Q-learning is applied. However, since the Q-function is learned as a table instead of parametrized by a neural network, the authors deviate from the current protocol and consider a different (much smaller) state-action space to enable learning. Another relevant approach is presented in [3], which considers a similar setting to ours but proposes a supervised learning strategy based on a random forest algorithm. As such, it performs well in a static scenario but fails to generalize to the cases where other nodes vary their CW parameters, as we illustrate through experiments.

**Contribution.** The contributions of this paper are twofold:

- i) We formulate the problem of adaptive MCW design as an RL problem and implement a deep Q-learning architecture for its solution.
- ii) Through NS3 simulations, we compare the proposed solution with established baselines and demonstrate its near-optimal performance across several scenarios.

## 2. SYSTEM MODEL AND PROBLEM STATEMENT

We consider a wireless network with  $N$  nodes transmitting packets to a central access point. Nodes follow a binary exponential backoff mechanism as mentioned in the IEEE 802.11 protocol for channel access. The exponential backoff mechanism reduces collisions in the channel by specifying the amount of time slots that each node must wait before a transmission [1]. More precisely, before every transmission, a node draws the number of slots to wait from a uniform distribution in  $(0, c - 1)$ , where  $c$  is equal to the MCW  $\omega$  if the previous transmission was successful. After the first collision,  $c$  is increased to  $2\omega$  to promote longer waiting times and reduce the probability of further collisions. In general,  $c$  is chosen as  $2^j \omega$  after  $j$  successive collisions, and saturates at a maximum permissible pre-specified value. The network is assumed to be operating in a high-load regime where every node has packets to transmit at all times. Furthermore, all nodes are assumed to be in range of all other nodes and there are no hidden nodes [1].

The network is observed in discrete intervals  $k = 1, 2, \dots$ , each of duration  $T$  seconds. Within each time interval  $k$ , we assume the MCW of all nodes  $\omega^k = [\omega_0^k, \omega_1^k, \dots, \omega_{N-1}^k]^\top$  to be constant. As described in Section 1, we consider the problem of CW design from the point of view of an intelligent node (here node 0) that seeks to optimize  $\omega_0^k$  without knowledge of the MCW of other nodes and, even more challenging, in a setting where these values might change from one interval  $k$  to the next. We model the MCW of every other node  $\omega_i^k$  for  $i \neq 0$  as a discrete-time stochastic process that can take values in a finite state space  $\Omega$ . In selecting the right MCW, node 0 must rely exclusively on local information. In particular, for each interval  $k$ , apart from knowing its own  $\omega_0^k$ , node 0 observes the fraction of time  $f^k$  that it transmits without collision and the fraction of time  $b^k$  that other nodes transmit without collision. Based on this, we define the observation vector  $\mathbf{o}^k = [f^k, b^k]^\top$ . Notice that  $\mathbf{o}^k$  is a random vector that depends on  $\omega^k$  through the (stochastic) exponential backoff mechanism. Lastly, to discuss *optimal* CW design we need to focus on a notion of utility. In our case, we associate interval  $k$  with the fairness utility given by

$$u(\mathbf{o}^k, N) = 1 - \left| \frac{f^k}{b^k + f^k} - \frac{1}{N} \right|. \quad (1)$$

Intuitively, with  $N$  nodes in the network and under the high-load assumption, the fair share of node 0 would be to transmit  $1/N$ -th of the total transmit time during interval  $k$ , while the actual proportion of transmit time is given by  $f^k / (b^k + f^k)$ . The discrepancy between these two numbers can be seen as a fairness loss and the utility  $u(\mathbf{o}^k, N)$  subtracts this loss from a perfect utility of 1. Notice that in computing (1), the total number of nodes  $N$  is assumed to be known to (or accurately estimated by) node 0. With this notation in place, we formally state our optimal CW design problem.

**Problem 1.** *Given  $N$  and the history  $\{\mathbf{o}^k, \omega_0^k\}_{k=1}^{K-1}$ , determine the MCW  $\omega_0^K$  to maximize  $\mathbb{E}(\sum_{k=K}^{\infty} \alpha^{k-K} u(\mathbf{o}^k, N))$ .*

At every discrete-time instant  $K$ , our intelligent node is faced with an instance of Problem 1: it knows the size of the network  $N$ , its past decisions  $\{\omega_0^k\}_{k=1}^{K-1}$ , and how these past decisions worked out in terms of realized proportions of channel occupancy  $\{\mathbf{o}^k\}_{k=1}^{K-1}$ . With this information, node 0 seeks to maximize future fairness. In general, it seeks to maximize not only the immediate fairness at time  $K$  but rather an infinite-horizon discounted sum, where  $\alpha \in (0, 1)$  establishes the relative importance of future utilities. Since  $\mathbf{o}^k$  is a random vector, future utilities are random variables, as is their discounted sum. Hence, the objective of node 0 is to maximize the

expected value of this discounted sum. In the next section, we frame Problem 1 from an RL perspective and propose to solve it via deep Q-learning. Finally, notice that even though in this paper we focus on fairness, a similar problem can be postulated to maximize throughput or other measures of interest by modifying the utility in (1).

## 3. DEEP Q-LEARNING FOR CW DESIGN

A first naive approach towards solving Problem 1 would be to try to explicitly compute the expected value of the discounted utilities. The first obstacle in achieving this is that we do not have access to the underlying stochastic process regulating the evolution of  $\omega_i^k$  for  $i \neq 0$ . Thus, we would first need to adopt a model for this stochastic process – e.g., a Markov process with  $|\Omega|^{N-1}$  states representing the  $|\Omega|$  possible choices of the  $N - 1$  nodes – and then estimate from data the parameters of this model such as the transition probability matrix of the Markov process. Apart from being computationally expensive, we would then face a second obstacle related to partial observability, namely that node 0 does not have access to the true state of the network  $\omega^k$  but rather to a real-valued observation  $\mathbf{o}^k$  whose distribution depends on  $\omega^k$ . Inferring the transitions between these hidden states further increases the complexity, indicating that a classical model-based approach is not well-suited to solve Problem 1.

We turn our attention to a learning-based approach under the framework of RL, where the lack of system knowledge is overcome through methodical interactions with the system. In particular, inspired by its success on other problems pertaining to wireless networks [17–21], we focus on Q-learning. RL problems are modelled using Markov decision processes where an agent learns by interacting with an environment. At every time step  $k$ , the agent takes an action  $a_k$  that affects the state of the environment. The state changes from  $s_k$  to  $s_{k+1}$  and the agent gets a reward  $r(s_k, a_k)$ . The action taken by the agent at each state  $a_k = \pi(s_k)$  is given by the policy  $\pi$ , a mapping from the set of states  $\mathcal{S}$  to the set of actions  $\mathcal{A}$ . The goal of Q-learning is to find an optimal policy  $\pi^*$  that maximizes the long-term expected accumulated discounted reward. To do this, let us define the optimal Q-function  $Q^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  where the value  $Q^*(s, a)$  corresponds to the long-term reward of selecting action  $a$  in state  $s$  and, from that point onward, following the optimal strategy  $\pi^*$ . From this definition, it follows that

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q^*(s, a). \quad (2)$$

In our setting, the RL agent is the intelligent node 0, the action space is given by the MCW choices, the reward is given by the utility in (1), and we model the state as an  $M$ -memory buffer (or history) of the local observations at node 0, i.e.  $\{\mathbf{o}^k, \omega_0^k\}_{k=K-M+1}^K$ . Having drawn this analogy, it follows that if we can compute  $Q^*$  then the solution of Problem 1 is given by the optimal strategy  $\pi^*$  as in (2).

The classical way of finding  $Q^*$  is through a value iteration method based on the Bellman equation [27]

$$Q(s, a) \leftarrow Q(s, a) + \eta \left( r(s, a) + \alpha \max_{a'} Q(s', a') - Q(s, a) \right), \quad (3)$$

where  $\eta$  is a pre-defined learning step-size,  $\alpha$  is the discount factor in Problem 1, and  $s'$  is the realized state after  $s$ . It has been shown [27] that under certain conditions on the step  $\eta$  and the frequency of visits to the different states, the generic Q-function in (3) is guaranteed to converge to  $Q^*$ , from where  $\pi^*$  can be obtained [cf. (2)], seemingly solving Problem 1. However, the convergence to  $Q^*$  can be extremely slow in practice.

Classical Q-learning performs well in settings where the state-action space is small, since the table-like iterative procedure in (3) updates each of the  $|\mathcal{S}| \times |\mathcal{A}|$  pairs separately. By contrast, in our setting we have an agent (node 0) whose action space can potentially be large (different choices for the MCW) and, more critically, the observations  $\mathbf{o}^k$  from the environment are real-valued. To extend the benefits of Q-learning to our more challenging scenario we rely on a *deep* Q-learning framework [28], where we parameterize the Q-function as a neural network and update its parameters as opposed to independent entry-wise updates of the Q-function. More precisely, with  $\theta$  representing the parameters of our neural network – usually called a deep Q-network (DQN) – we seek to determine the optimal parameters  $\theta^*$  such that  $Q_{\theta^*}(s, a) \approx Q^*(s, a)$ .

The specific implementation of deep Q-learning for the solution of Problem 1 is illustrated in Fig. 1. At a generic time interval  $K$ , the communication of packets to the access point is governed by the choice of the MCW of every user. In the example, every user has selected the same MCW at time  $K$  (depicted by their blue states) but this is not known by node 0. Instead, node 0 observes the operation of the wireless network  $\mathbf{o}^K$  and, of course, knows its own  $\omega_0^K$ . The concatenation of these observations for the last  $M$  time intervals is used as the input to a trained neural network (Rainbow DQN). The output of this neural network consists of a vector containing the Q-values corresponding to every possible action (MCW) that node 0 can select. Node 0 will select as  $\omega_0^{K+1}$  the MCW that attains the largest Q-value at the output of the DQN [cf. (2)], while  $\omega_i^{K+1}$  for  $i \neq 0$  evolves following an unknown stochastic process. This same process is repeated and node 0 relies again on the neural network to select  $\omega_0^{K+2}$ , and so on. As we illustrate in Section 4 for several settings, the discounted sum of utilities yielded by the trained rainbow DQN with optimal parameters  $\theta^*$  is comparable with the best attainable strategy when full knowledge of the system is available.

In order to find the optimal parameters  $\theta^*$  for our neural network in the first place, we rely on the established strategy of experience replay training; see Section 4 for implementation details. In a nutshell, inspired by the fixed-point solution of the classical iteration in (3), we update the values of  $\theta$  via gradient descent to minimize the following loss

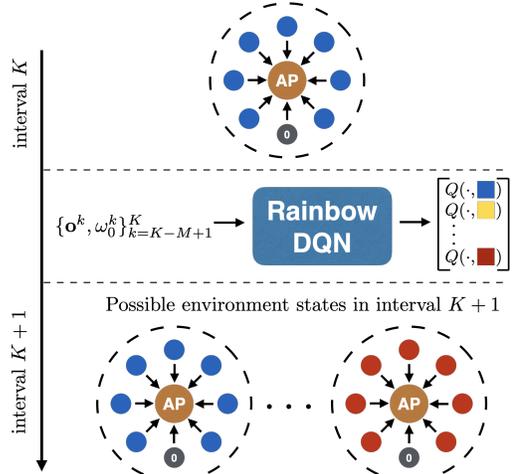
$$\mathcal{L}(\theta) = \mathbb{E}_{s'} \left( Q_{\theta}(s, a) - \left( r(s, a) + \alpha \max_{a'} Q_{\theta}(s', a) \right) \right)^2. \quad (4)$$

It should be noted that, in our implementation, instead of using a simple multi-layer perceptron to approximate the optimal Q-function, we rely on recent modifications that have shown empirical improvements in terms of training stability and amount of data needed for training. Specifically, the rainbow agent [29] incorporates six improvements (double DQN, prioritized replay, duelling networks, multi-step learning, distributional RL, and noisy nets) that led to superior performance over vanilla DQN in gaming environments. We have also witnessed those benefits in the implementation of deep Q-learning for CW design.

#### 4. NUMERICAL EXPERIMENTS

To evaluate the performance of the proposed approach, we simulate a wireless network using the NS3 simulator. Nodes in the network are modelled to follow the IEEE 802.11b protocol with a constant data rate of 1 Mbps. Furthermore, a constant packet size of 1400 bytes is used and the observation interval length  $T$  is fixed as 20 seconds. Unless otherwise stated, the considered networks have  $N = 10$  nodes.

We use two coupled feed forward neural networks as our rainbow DQN, each having four layers with 32 nodes per layer and



**Fig. 1:** Schematic view of our proposed deep RL agent. Node 0 relies on a rainbow DQN to select the MCW for the next time interval based on historical local observations.

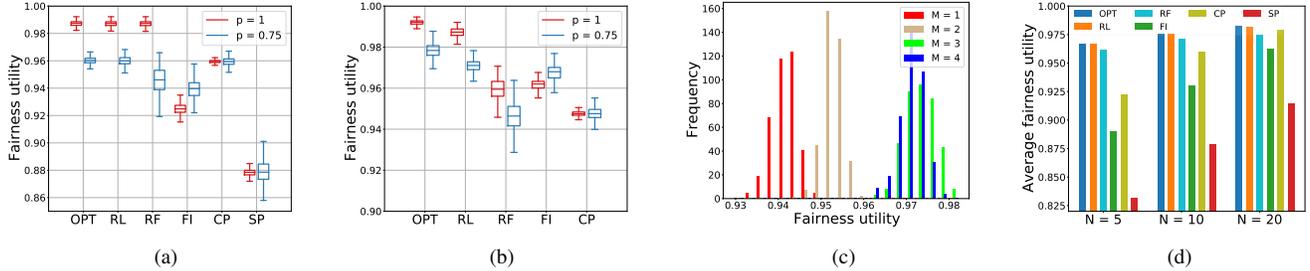
ReLU non-linearities.<sup>1</sup> For all experiments, the discount factor is set as  $\alpha = 0.9$  and the DQN is trained for 5000 episodes. At the start of each episode, node 0 chooses an action uniformly from the set of actions and the trajectory evolves for a specified number of time steps. We fix the number of time steps in each episode to be equal to 50. Rainbow DQN is implemented using the PyTorch [30] library in Python and training is done using mini-batch gradient descent with a batch size of 32. The buffer size for experience replay is set to 10000.

In evaluating the performance of our proposed RL approach, we compare it with the following five baselines:

- i) **Optimal design (OPT):** This benchmark assumes full knowledge of the underlying stochastic process and computes explicitly the expected value in Problem 1. Although unrealistic in practice, it sets the upper bound of what is achievable by other methods.
- ii) **Random forest (RF) [3]:** The classifier is trained using the local observations  $\{\mathbf{o}^k, \omega_0^k\}_{k=K-M+1}^K$  as input and  $\omega_0^{K+1}$  as target labels. The training labels  $\omega_0^{K+1}$  are selected as the optimal actions that maximize the fairness utility at the next time step. We fix the number of trees as 15 and the depth of each tree as 5.
- iii) **Fairness index approach (FI) [2]:** Node 0 computes a local fairness index and follows a threshold-based policy to update its MCW. We set this threshold as  $C = 1.5$ ; see [2] for details.
- iv) **Optimal constant policy (CP):** Just like for OPT, we assume complete knowledge of the system. However, we restrain the design of the MCW to be constant, i.e., it cannot change across time steps. Improvements upon CP illustrate the value of an adaptive design.
- v) **Standard protocol (SP):** Node 0 follows a static policy as mentioned in the IEEE 802.11 protocol with its MCW fixed at 32.

Before proceeding with the evaluation, it should be noted that FI is not a predictive method but rather a reactive one in the sense that, whenever the system deviates from fairness, node 0 would react and try to restore a fair setting. Thus, it is bound to underperform in a fast-changing setting. For a fairer comparison, we consider the case where the update speed of FI is three times that of other methods,

<sup>1</sup>Code to replicate the numerical experiments here presented can be found at <https://github.com/kumarabhish3k/Rainbow-DQN-for-Contention-Window-design.git>.



**Fig. 2:** Performance evaluation of the proposed RL method. (a) Boxplots for fairness utility over 500 episodes of 50 time steps for the six methods under consideration. A Markov process with  $p \in \{0.75, 1\}$  is considered for the evolution of the MCW of other nodes. (b) Counterpart of (a) for the more complex stochastic process described in the text. SP is not shown due to its low relative performance, achieving mean fairness utilities of 0.752 for both values of  $p$ . (c) Histograms of the utility attained by the RL agent for different lengths of the memory buffer  $M$ . The respective mean utilities yielded by the increasing values of  $M$  are 0.941, 0.953, 0.973, and 0.971. (d) Average network utility over 500 episodes for varying number of nodes  $N \in \{5, 10, 20\}$ .

i.e., FI updates its MCW thrice in every time slot as opposed to once for the other methods (and zero times for CP and SP).

We first consider a simple scenario where the MCWs of all other nodes follow a Markov process with only two states. In state 1,  $\omega_i = 32$  for all  $i \neq 0$  and, in state 2,  $\omega_i = 128$  for all  $i \neq 0$ . At each discrete time step, the system switches states with probability  $p$ , and we analyze the cases where  $p \in \{0.75, 1\}$ . The action space of node 0 (potential MCW choices) is  $\mathcal{A} = \{32, 48, 64, 96, 128\}$ . Since the underlying process is Markovian, we consider a memory-less implementation with  $M = 1$  (cf. Fig. 1). Fig. 2a reveals that, when  $p = 1$ , RL and RF achieve optimal performance even in the absence of knowledge of the underlying Markov process. The advantage of RL over the supervised RF cannot be appreciated in this deterministic process. The lack of adaptability of CP and SP makes them lag behind and this fast-changing setting is not well suited for FI even with updates three times faster. This latter effect is attenuated when  $p = 0.75$  since the system tends to remain longer in a given state. More importantly, for  $p = 0.75$ , RL is still able to learn the optimal policy while the performance of RF is degraded. As a consequence of the stochasticity of the process, similar observations  $\{\omega^k, \omega_0^k\}$  can lead to drastically different optimal values for  $\omega_0^{*k+1}$ , thus confusing the supervised RF classifier.

As our second evaluation setting, we consider the case where the MCWs of others follow a more complex process. More precisely, we consider five states  $\{s_j\}_{j=1}^5$  where  $\omega_i = 2^{j-1}32$  for all  $i \neq 0$  at state  $j$ . At each time step, the system transitions to the next state with probability  $p$ , however, the next state follows an increasing trajectory from  $s_1$  to  $s_5$  followed by a decreasing trajectory from  $s_5$  to  $s_1$ , and so on. In this sense, the next state does not only depend on the current state but also on the previous one. To incorporate memory in the predictive methods RL and RF, we consider  $M = 4$ . Furthermore, we consider the action space  $\mathcal{A} = \{2^{j-1}32\}_{j=1}^5 \cup \{2^{j-1}48\}_{j=1}^4$ . Fig. 2b reveals that RL achieves the closest utility to optimal among all methods considered for both values of  $p$ . As expected, the scenario where  $p = 0.75$  is more challenging except for the reactive method FI that benefits from slower changes in MCWs. Also notice that the performance gap between OPT and RL is reasonable since the former has perfect knowledge of the underlying system whereas the latter must rely exclusively on local observations to gather understanding of the underlying dynamics. The absence of this gap in Fig. 2a can be attributed to the simplicity of the model and the smaller size of  $\mathcal{A}$ .

Selecting  $M > 1$  is essential for the RL agent to meaningfully

learn from experience in this more complex setting and, for the cases where  $p < 1$ , having  $M > 2$  can also be beneficial. This is illustrated in Fig. 2c, where we present the histograms of the fairness utility attained by RL for different values of  $M$  over 500 episodes for the case  $p = 0.75$ . There is a big increase in performance when going from  $M = 1$  to  $M = 2$  in accordance with the memory of the underlying system. Moreover, a similar jump is observed when going from  $M = 2$  to  $M = 3$ . This is as expected since, for this value of  $p$ , the transition to the next state might be delayed underscoring the value of longer memory. The marginal gain for  $M = 4$  decreases and, in fact, a small loss is observed which can be attributed to the fact that a larger model must be fit (larger dimension of input to the rainbow DQN) without apparent modeling gain of the environment.

Finally, we consider the Markov process analyzed in our first experiment but for  $p = 0.9$  and vary the number of nodes  $N \in \{5, 10, 20\}$ . In Fig. 2d we portray the utility attained by the six considered methods averaged over 500 episodes as a function of the network size. It can be observed that RL achieves utility closest to optimal for all network sizes. As one might expect, the effect that an intelligent node has in promoting fairness gets diluted in larger networks, thus, resulting in a smaller dynamic range of utilities across methods when  $N = 20$ . In particular, the constant conservative strategy of CP that selects a large MCW ( $\omega_0^k = 128$  for all  $k$ ) attains a small suboptimality gap, but does not match the RL performance. Note that RL outperforms the standard protocol for all values of  $N$ . Indeed, the consistent near-optimality of the proposed RL approach across network sizes is encouraging and motivates a range of related future work.

## 5. CONCLUSIONS AND FUTURE WORK

We formulated the design of CWs in wireless networks as an RL problem. Within this framework, we leveraged a state-of-the-art DQN architecture to select the optimal value of the MCW in a sequential manner exclusively from local observations. Through NS3 simulations, we showed that the proposed RL agent achieves close-to-optimal behavior and outperforms other learning-based and rule-based methods. Ongoing work includes: i) the implementation of RL for parameter learning on other access protocols, ii) the combination of RL with graph neural networks to accommodate the case where every node is intelligent and can be trained based on its local observations, and iii) the incorporation of more realistic topologies that go beyond star connections to a common access point.

## 6. REFERENCES

- [1] G. Bianchi, "Performance analysis of the IEEE 802.11 distributed coordination function," *IEEE J. Sel. Areas Commun.*, vol. 18, no. 3, pp. 535–547, Mar. 2000.
- [2] B. Bensaou, Y. Wang, and C. C. Ko, "Fair medium access in 802.11 based wireless ad-hoc networks," in *IEEE Wrkshp. Mobile and Ad Hoc Netw. and Comp. (MoobiHOC)*, 2000, pp. 99–106.
- [3] A. H. Y. Abyaneh, M. Hirzallah, and M. Krunz, "Intelligent-CW: AI-based framework for controlling contention window in WLANs," in *IEEE Int. Symp. on Dynamic Spectrum Access Netw. (DySPAN)*, 2019, pp. 1–10.
- [4] Defense Advanced Research Projects Agency, "Spectrum Collaboration Challenge (SC2)," <https://www.darpa.mil/program/spectrum-collaboration-challenge>.
- [5] Y. Yu, T. Wang, and S. C. Liew, "Deep-reinforcement learning multiple access for heterogeneous wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 6, pp. 1277–1290, 2019.
- [6] Y. Yu, S. C. Liew, and T. Wang, "Carrier-sense multiple access for heterogeneous wireless networks using deep reinforcement learning," in *IEEE Wireless Commun. and Netw. Conf. Wrkshp. (WCNCW)*, 2019, pp. 1–7.
- [7] J. Tan, L. Zhang, Y.-C. Liang, and D. Niyato, "Deep reinforcement learning for the coexistence of LAA-LTE and WiFi systems," in *IEEE Int. Conf. on Commun.*, 2019, pp. 1–6.
- [8] I. Syed and B.-h. Roh, "Adaptive backoff algorithm for contention window for dense IEEE 802.11 WLANs," *Mobile Info. Syst.*, vol. 2016, 2016.
- [9] Q. Xia and M. Hamdi, "Contention window adjustment for IEEE 802.11 WLANs: a control-theoretic approach," in *IEEE Int. Conf. on Commun.*, 2006, vol. 9, pp. 3923–3928.
- [10] L. Chen, S. H. Low, and J. C. Doyle, "Random access game and medium access control design," *IEEE/ACM Trans. Netw.*, vol. 18, no. 4, pp. 1303–1316, 2010.
- [11] S. Chun, D. Xianhua, L. Pingyuan, and Z. Han, "Adaptive access mechanism with optimal contention window based on node number estimation using multiple thresholds," *IEEE Trans. Wireless Commun.*, vol. 11, no. 6, pp. 2046–2055, 2012.
- [12] K. Hong, S. Lee, K. Kim, and Y. Kim, "Channel condition based contention window adaptation in IEEE 802.11 WLANs," *IEEE Trans. Commun.*, vol. 60, no. 2, pp. 469–478, 2012.
- [13] D.-J. Deng, C.-H. Ke, H.-H. Chen, and Y.-M. Huang, "Contention window optimization for IEEE 802.11 DCF access control," *IEEE Trans. Wireless Commun.*, vol. 7, no. 12, pp. 5129–5135, 2008.
- [14] A. Ksentini, A. Nafaa, A. Gueroui, and M. Naimi, "Determinist contention window algorithm for IEEE 802.11," in *IEEE Int. Symp. on Personal, Indoor and Mobile Radio Commun.*, 2005, vol. 4, pp. 2712–2716.
- [15] V. Bharghavan, A. Demers, S. Shenker, and L. Zhang, "MACAW: a media access protocol for wireless LANs," *ACM SIGCOMM Computer Commun. Review*, vol. 24, no. 4, pp. 212–225, 1994.
- [16] N.-O. Song, B.-J. Kwak, J. Song, and M. Miller, "Enhancement of IEEE 802.11 distributed coordination function with exponential increase exponential decrease backoff algorithm," in *IEEE Veh. Tech. Conf. (VTC)*, 2003, vol. 4, pp. 2775–2778.
- [17] O. Naparstek and K. Cohen, "Deep multi-user reinforcement learning for dynamic spectrum access in multichannel wireless networks," in *IEEE Global Commun. Conf. (GLOBECOM)*, 2017, pp. 1–7.
- [18] S. Wang, H. Liu, P. H. Gomes, and B. Krishnamachari, "Deep reinforcement learning for dynamic multichannel access in wireless networks," *IEEE Trans. on Cognitive Commun. and Netw.*, vol. 4, no. 2, pp. 257–265, 2018.
- [19] N. C. Luong, D. T. Hoang, S. Gong, D. Niyato, P. Wang, Y.-C. Liang, and D. I. Kim, "Applications of deep reinforcement learning in communications and networking: A survey," *IEEE Commun. Surveys & Tutorials*, vol. 21, no. 4, pp. 3133–3174, 2019.
- [20] H. Ye, G. Y. Li, and B.-H. F. Juang, "Deep reinforcement learning based resource allocation for V2V communications," *IEEE Trans. Vehicular Tech.*, vol. 68, no. 4, pp. 3163–3173, 2019.
- [21] Y. S. Nasir and D. Guo, "Multi-agent deep reinforcement learning for dynamic power allocation in wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 10, pp. 2239–2250, 2019.
- [22] A. Chowdhury, G. Verma, C. Rao, A. Swami, and S. Segarra, "Unfolding wmmse using graph neural networks for efficient power allocation," *arXiv preprint arXiv:2009.10812*, 2020.
- [23] Y. Edalat and K. Obraczka, "Dynamically tuning IEEE 802.11's contention window using machine learning," in *ACM Int. Conf. on Modeling, Analysis and Simulation of Wireless and Mobile Syst.*, 2019, pp. 19–26.
- [24] A. Pressas, Z. Sheng, F. Ali, and D. Tian, "A Q-learning approach with collective contention estimation for bandwidth-efficient and fair access control in IEEE 802.11p vehicular networks," *IEEE Trans. Vehicular Tech.*, vol. 68, no. 9, pp. 9136–9150, 2019.
- [25] Z. Ali, L. Giupponi, J. Mangues-Bafalluy, and B. Bojovic, "Machine learning based scheme for contention window size adaptation in LTE-LAA," in *IEEE Int. Symp. on Personal, Indoor, and Mobile Radio Commun. (PIMRC)*, 2017, pp. 1–7.
- [26] S. Amuru, Y. Xiao, M. van der Schaar, and R. M. Buehrer, "To send or not to send-learning MAC contention," in *IEEE Global Commun. Conf. (GLOBECOM)*, 2015, pp. 1–6.
- [27] R. S. Sutton and A. G. Barto, *Introduction to Reinforcement Learning*, vol. 135, MIT press Cambridge, 1998.
- [28] V. Mnih, K. Kavukcuoglu, D. Silver, et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [29] M. Hessel, J. Modayil, H. Van Hasselt, et al., "Rainbow: Combining improvements in deep reinforcement learning," *arXiv preprint arXiv:1710.02298*, 2017.
- [30] A. Paszke, S. Gross, F. Massa, et al., "Pytorch: An imperative style, high-performance deep learning library," in *Adv. in Neur. Inf. Process. Syst.* 32, pp. 8024–8035, 2019.